

Statistical Inference

Anthony Davison

©2022

<http://stat.epfl.ch>

Introduction	2
Some Basic Concepts	40
Hypothesis Testing	56
Data Complications	86
Likelihood Theory	96
Bootstrap Inference	128
Bayesian Inference	165
Bayesian Computation	198
Hierarchical Models	223

Starting point

- ☐ We start with a concrete question, e.g.,
 - Does the Higgs boson exist?
 - Is fraud taking place at this factory?
 - Are these two satellites likely to collide soon?
 - Do lockdowns reduce Covid transmission?
- ☐ We aim
 - to use **data**
 - to provide **evidence** bearing on the question,
 - to draw a **conclusion** or sometimes reach a **decision**.
- ☐ Here we mostly discuss how to express the evidence, but the choice and quality of the data, and how they were obtained, affect the evidence and the clarity of any decision.
- ☐ The data typically display both **structure** and **haphazard variation**, so any conclusion reached is uncertain, i.e., is an **inference**.
- ☐ A **statistical inference** uses probability models to express the variation in the data.

stat.epfl.ch

Autumn 2022 – slide 3

Theory

- ☐ The goal of statistical theory is to clarify the foundations of statistical inference, which need not involve deep mathematics.
- ☐ Some reasons for studying theory:
 - to clarify underlying issues;
 - to guide strategy for applied work;
 - to elucidate the common structure of diverse applied procedures and thus to formulate general approaches;
 - to transfer ideas between different application domains;
 - to give a basis for developing new methods (including adaptations for specific applied problems)
 - to give a basis (or bases) for comparing competing procedures; and
 - because it's interesting.

stat.epfl.ch

Autumn 2022 – slide 4

Mathematics, asymptotics and all that

- Statistical science is a mathematical science (like physics, ...): problems are often formulated in mathematical terms, but success or failure is judged in consequences for the real world, not in an abstract realm.
- Very often we argue as sample size n (or some other measure of information) becomes large, because
 - exact calculations are impossible or too onerous to be worthwhile;
 - asymptotic arguments clarify the essential structure of a problem (and in particular the necessary assumptions), and thus help generalisation;
 - a procedure that fails when $n = \infty$ will be suspect for finite n .
- Such arguments are often used to generate approximations useful for some fixed finite n , so we should always ask 'will this approximation be adequate in this context?' Often this involves
 - Monte Carlo studies with a realistic sample size;
 - general knowledge (e.g., Laplace approximations generally work better than Edgeworth approximation);
 - previous experience (hard to teach).

stat.epfl.ch

Autumn 2022 – slide 5

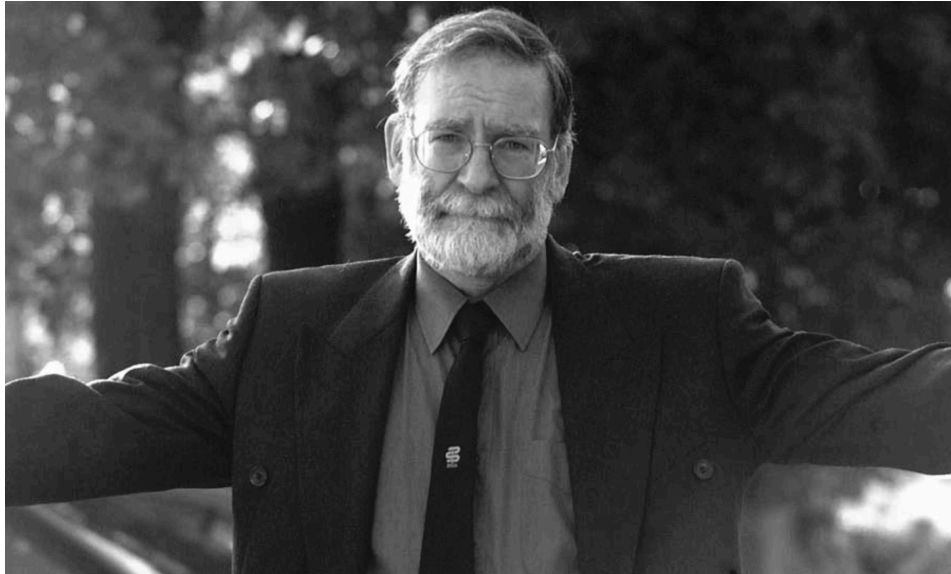
Types of statistical analysis

- Conventionally divided into
 - **descriptive statistics/exploratory data analysis** — informal (often graphical) methods used to gain insight into data, used for initial data analysis and for presenting conclusions;
 - **statistical inference** — more formal approaches using probability models to clarify the support for particular statements about the underlying situation;to which we nowadays add
 - **algorithmic methods** — machine learning algorithms, generally complex and computationally demanding, often used for prediction/decision-making.
- Even the first and the last are studied using probability models (e.g., formulation of a boxplot, 'is that difference significant?', analysis of neural networks).
- Need to be clear what type of analysis is being conducted, otherwise conclusions may be horribly biased.

stat.epfl.ch

Autumn 2022 – slide 6

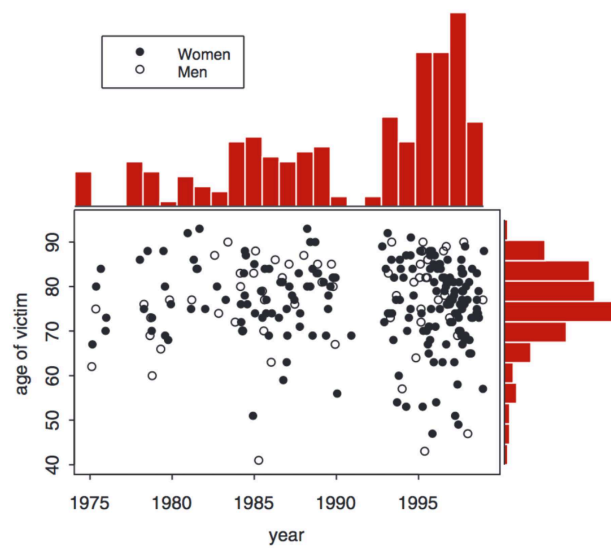
Dr Harold Shipman



stat.epfl.ch

Autumn 2022 – slide 7

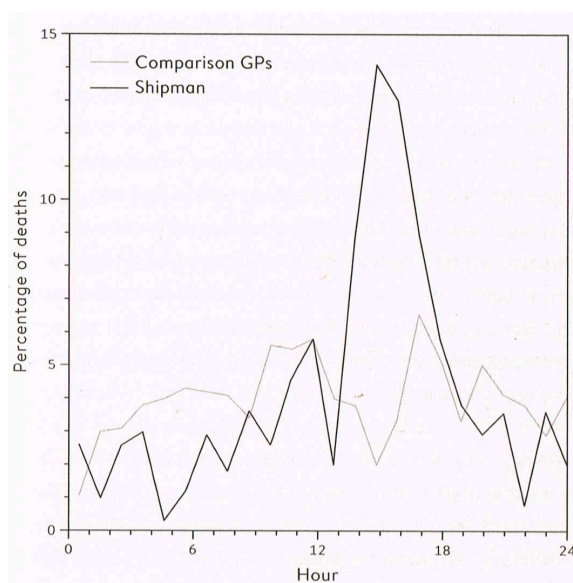
Descriptive statistics



stat.epfl.ch

Autumn 2022 – slide 8

More descriptive statistics

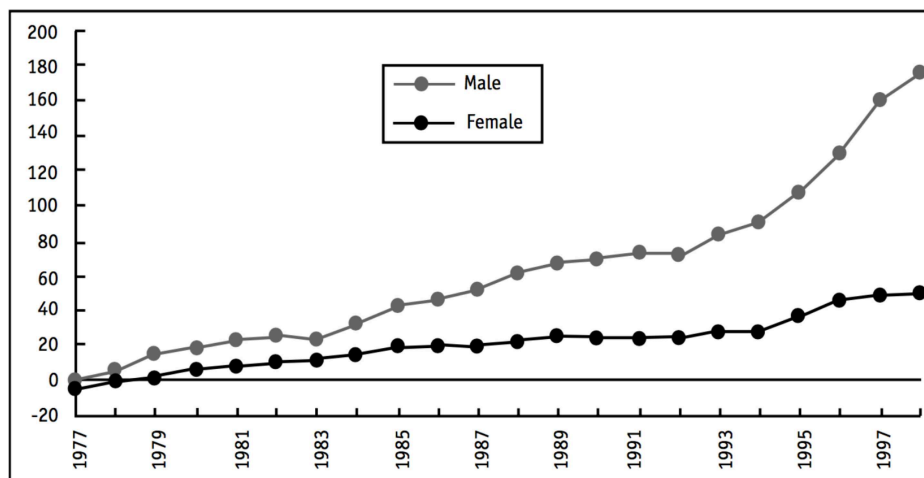


stat.epfl.ch

Autumn 2022 – slide 9

Excess mortality statistics

Figure 2. Cumulative excess death certificates signed by Shipman, for people older than 64 and who died at home or in his practice



stat.epfl.ch

Autumn 2022 – slide 10

Data

- Theoretical discussion generally takes observed data as given, but
 - to get the data we may need to **plan an investigation**, perhaps **design an experiment** largely controlled by the investigator — not considered here but often crucial to obtaining strong data and hence secure conclusions; or
 - to use data from an **observational study** (the investigator has little or no control over data collection).
- In both cases the data used may be selected from those available, and especially if we have ‘found data’ we must ask
 - why am I seeing these data?
 - what exactly was measured, and how?
 - can the observations actually shed light on the problem?
 - will using a function of the available data give more insight?
- Below we assume that these questions have suitable answers, and we can continue with our analysis.

stat.epfl.ch

Autumn 2022 – slide 11

Example: Satellite conjunction

- Goal is to say whether two ‘space objects’ will collide, based on repeated noisy measurements of their velocities v and positions p .
- In a ‘close conjunction’ the problem is usually simplified:
 - only the latest measurements are used,
 - the relative position and velocity of the second object relative to the first are considered, giving $p_2 - p_1$ (km), $v_2 - v_1$ (km/s),
 - the relative motion is taken as rectilinear (Newton I), the objects are represented as spheres and their radii are added to give a ‘combined hard-body radius’ (HBR),so we ask

based on noisy measurements of $p_2 - p_1$ and $v_2 - v_1$, what is the evidence that the true path of object 2 will pass inside the HBR?

- A simple probability model for the relative position and velocity is

$$\begin{pmatrix} p_2 - p_1 \\ v_2 - v_1 \end{pmatrix} \sim \mathcal{N}_6 \left\{ \begin{pmatrix} \mu \\ \nu \end{pmatrix}, \begin{pmatrix} \Omega & \Omega' \\ \Omega'^T & \Omega'' \end{pmatrix} \right\}$$

where μ and ν represent the unknown true relative position and velocity vectors and the 3×3 matrices Ω , Ω' and Ω'' are treated as known.

stat.epfl.ch

Autumn 2022 – slide 12

Multivariate normal distribution

A random variable $X_{n \times 1}$ with real components has the **multivariate normal distribution**, $X \sim \mathcal{N}_n(\mu, \Omega)$, if $a^T X \sim \mathcal{N}(a^T \mu, a^T \Omega a)$ for any constant vector $a_{n \times 1}$, and then

- the mean vector and covariance matrix of X are

$$E(X) = \mu_{n \times 1}, \quad \text{var}(X) = \Omega_{n \times n},$$

where Ω is symmetric semi-positive definite with real components;

- for any matrix $A_{m \times n}$ and vector $b_{m \times 1}$ of constants,

$$AX + b \sim \mathcal{N}_m(A\mu + b, A\Omega A^T);$$

- X has a density on \mathbb{R}^n iff Ω is positive definite (i.e., has rank n), and then

$$f(x; \mu, \Omega) = \frac{1}{(2\pi)^{n/2} |\Omega|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Omega^{-1} (x - \mu) \right\}, \quad x \in \mathbb{R}^n; \quad (1)$$

- If $X^T = (X_1^T, X_2^T)$, where X_1 is $m \times 1$, and μ and Ω are partitioned correspondingly, then the marginal and conditional distributions of X_1 are also multivariate normal:

$$X_1 \sim \mathcal{N}_m(\mu_1, \Omega_{11}), \quad X_1 | X_2 = x_2 \sim \mathcal{N}_m \left\{ \mu_1 + \Omega_{12} \Omega_{22}^{-1} (x_2 - \mu_2), \Omega_{11} - \Omega_{12} \Omega_{22}^{-1} \Omega_{21} \right\}.$$

stat.epfl.ch

Autumn 2022 – slide 13

Example: Satellite conjunction II

- When the velocity measurements can be treated as noiseless compared to the position measurements, i.e., $\nu = v_2 - v_1$ (equivalently $\Omega' = \Omega'' = 0$), then we can write

$$X = p_2 - p_1 \sim \mathcal{N}_3(\mu, \Omega), \quad \psi_{\min} = \text{HBR} > 0,$$

and object 2 traverses the line $\{\mu + t\nu : t \geq 0\}$.

- We use a matrix $A_{2 \times 3}$ to project X and μ into the plane \mathcal{P} perpendicular to ν and passing through the origin, and obtain the two-dimensional model

$$\begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} = Y = AX \sim \mathcal{N}_2(\xi, D^{-1}), \quad D^{-1} = A\Omega A^T = \begin{pmatrix} d_1^{-1} & 0 \\ 0 & d_2^{-1} \end{pmatrix},$$

say, with known d_1, d_2 (as Ω is known and A depends on the known ν), and unknown

$$\xi = \begin{pmatrix} \xi_1 \\ \xi_2 \end{pmatrix} = \begin{pmatrix} \psi \cos \lambda \\ \psi \sin \lambda \end{pmatrix}, \quad \psi > 0, 0 \leq \lambda < 2\pi.$$

- We observe a value x of X and ask

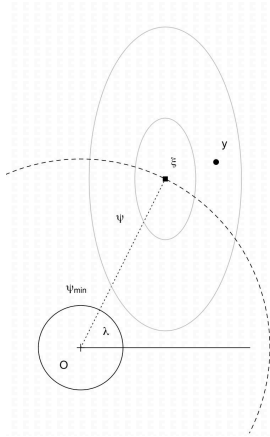
based on the projection $y = Ax$ of x into \mathcal{P} , what is the evidence that $\xi = A\mu$ lies inside the HBR, or equivalently that $\psi = \|\xi\| < \psi_{\min}$?

stat.epfl.ch

Autumn 2022 – slide 14

Example: Satellite conjunction III

- We seek inferences about the scalar $\psi = \|\xi\|$ based on the bivariate observation y .
- A natural estimate of ψ is $\|y\| = (y_1^2 + y_2^2)^{1/2}$, but geometrically it is obvious that $P(\|y\| > \psi) > 1/2$, i.e., $\|y\|$ will tend to be larger than ψ .
- The sample size is $n = 1$, and any asymptotics are as $d_1, d_2 \rightarrow \infty$, i.e., $\text{var}(Y) \rightarrow 0$.



Some definitions

- We use the term **(probability) density** for both the PMF of a discrete random variable, and the PDF of a continuous random variable.
- A **statistical model** for data y is a probability density $f(y)$ defined for $y \in \mathcal{Y}$.
- A **parametric model** $f \equiv f(y; \theta)$ is determined by **parameters** $\theta \in \Theta \subset \mathbb{R}^d$. If no such θ exists, f is **nonparametric**.
- Sometimes we use the term **family of models** to stress that there are many possibilities, $\{f(y; \theta) : \theta \in \Theta\}$.
- Generally $\theta = (\psi, \lambda)$ splits into
 - **parameters of interest (interest parameters)** ψ , usually scalar, that we focus on,
 - **nuisance parameters** λ , usually vector, needed to complete the model but not of main interest,
 though different elements of θ may become of interest during an investigation.
- By convention we (try to) use
 - letters like c, d, \dots for (known) constants,
 - Roman letters for random variables X, Y, \dots and their realisations x, y, \dots ,
 - Greek letters $\mu, \nu, \psi, \lambda, \Omega, \Delta, \dots$ for unknown parameters.

Some shortcuts

- $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f$ means that the X_j are independent and all have the density f , and we call the X_j a **random sample of size n from f** .
- $X_1, \dots, X_n \stackrel{\text{ind}}{\sim} f_1, \dots, f_n$ means that the X_j are independent and $X_j \sim f_j$.
- $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} (\mu, \sigma^2)$ means that the X_j are independent with mean μ and variance σ^2 (we assume that $0 < \sigma^2 < \infty$). Here we do not assume that the X_j are normal.
- $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} (\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2)$ means that the X_j are independent with means μ_j and variances σ_j^2 (where $0 < \sigma_j^2 < \infty$).
- We write $\mathcal{N}(0, 1)$ for the standard normal distribution, which has density and distribution functions

$$\phi(u) = \frac{1}{\sqrt{2\pi}} e^{-u^2/2}, \quad u \in \mathbb{R}, \quad \Phi(z) = \int_{-\infty}^z \phi(u) du, \quad z \in \mathbb{R}.$$

and write z_p for the p quantile of the standard normal distribution, i.e., $\Phi(z_p) = p$, where $0 < p < 1$.

Example: Satellite conjunction IV

- If $Y \sim \mathcal{N}_n(\mu, \Omega)$, then $\theta = (\mu, \Omega) \in \Theta \subset \mathbb{R}^n \times \mathcal{S}$, where \mathcal{S} denotes the set of all $n \times n$ symmetric positive definite matrices, so this is a parametric model.
- In the satellite example, we assume that the observed position y is a realisation of

$$Y \sim \mathcal{N}_2 \left\{ \begin{pmatrix} \psi \cos \lambda \\ \psi \sin \lambda \end{pmatrix}, \begin{pmatrix} d_1^{-1} & 0 \\ 0 & d_2^{-1} \end{pmatrix} \right\},$$

where $d_1, d_2 > 0$, $\theta = (\psi, \lambda) \in \Theta = \mathbb{R}_+ \times [0, 2\pi)$, the density function is

$$\frac{(d_1 d_2)^{1/2}}{2\pi} \exp \left[-\frac{1}{2} \{ d_1 (y_1 - \psi \cos \lambda)^2 + d_2 (y_2 - \psi \sin \lambda)^2 \} \right], \quad y_1, y_2 \in \mathbb{R}.$$

- Given the constants d_1, d_2 and observed data y_1, y_2 ,
 - the interest parameter is the length ψ of ξ ,
 - the nuisance parameter is the angle λ , and
 we want to know the evidence that $\psi < \psi_{\min}$.
- A natural estimate of ψ is $\|y\| = (y_1^2 + y_2^2)^{1/2}$, but this will tend to be too long.

Objectives

- ☐ Given a specified family of models $f(y; \theta)$ with $\theta = (\psi, \lambda)$ and observed data y , the objectives of a statistical analysis might be to
 - give intervals (or more generally sets) of values within which ψ is ‘likely to lie’;
 - assess the consistency of y with a particular ψ_0 ;
 - predict as-yet unobserved random variables from the system that generated y ;
 - use y to choose one of a specified set of decisions — requires the specification of the consequences of the decisions.
- ☐ We might also seek to check whether the family of models adequately represents the data.

stat.epfl.ch

Autumn 2022 – slide 19

Probability

- ☐ Two distinct roles of probability in statistical analysis:
 - as a description of **variation** in data (‘aleatory probability’, ‘chance’), using a probability model and treating the observed data y as an outcome of that model;
 - to formulate **uncertainty** (‘epistemic probability’) about the reality modelled in terms of the random experiment, based on y .

stat.epfl.ch

Autumn 2022 – slide 20

Probability models for variation

- ☐ Two broad types of probability model:
 - **substantive** — based on fundamental subject-matter theory (e.g., particle physics, Mendelian genetics, Navier–Stokes equations);
 - **empirical** — a convenient, adequately realistic, representation of data variation;
 - and of course there is a spectrum between them.
- ☐ The satellite example is partly substantive (using Newtonian mechanics for the ideal trajectory $\mu + t\nu$) and partly empirical (normal distributions for measurement error).
- ☐ We aim that
 - primary questions/issues are encapsulated in the parameter(s) of interest;
 - secondary aspects can be taken into account, often via nuisance parameters;
 - variation in the data is realistically modelled, leading to reasonable statements of uncertainty;
 - any special feature of the data collection process is represented;
 - different approaches to analysis can if necessary be compared.
- ☐ Such models are always provisional and should if possible be checked against data.

stat.epfl.ch

Autumn 2022 – slide 21

Uncertainty

- Essentially three bases for statements of uncertainty:
 - a **Bayesian (inverse probability) inference** expresses it via a prior probability density and uses Bayes' theorem to update this in light of the data;
 - a **frequentist (sampling theory) inference** compares y with the set \mathcal{S} of other data that might have been observed in a hypothetical sampling experiment;
 - in a designed experiment, clinical trial, sample survey or similar the investigator uses **randomisation** to generate a distribution against which y is compared.
- There are many variants of the first two approaches.
- In particular, a frequentist should choose \mathcal{S} thoughtfully:
 - what doctors would be a suitable comparison group for Shipman?

Example 1 (Measuring machines) A physical quantity θ can be measured with two machines, both giving normal observations Y such that $E(Y) = \theta$. A measurement from machine 1 has variance 1, and one from machine 2 has variance 100. A machine is chosen by tossing a fair coin, giving $M = 1, 2$ with equal probabilities.

If we observe $m = 1$ and $y = 2$, then clearly we can ignore the fact that we might have observed $m = 2$, i.e., we should take $\mathcal{S} = \{(y, 1) : y \in \mathbb{R}\}$ rather than $\mathcal{S} = \{(y, m) : y \in \mathbb{R}, m \in \{1, 2\}\}$.

stat.epfl.ch

Autumn 2022 – slide 22

Bayesian inference

- Our observed data y^o are assumed to be a realisation from a density $f(y | \psi)$.
- If we can summarise information about ψ , separately from y^o , in a **prior density** $f(\psi)$, then we can use Bayes' theorem to obtain the **posterior density**

$$f(\psi | y^o) = \frac{f(y^o | \psi)f(\psi)}{\int f(y^o | \psi)f(\psi) d\psi},$$

and base all our uncertainty statements on this.

- For example, if ψ_p satisfies $P(\psi \leq \psi_p | y^o) = p$ for any $p \in (0, 1)$, we could give a **$(1 - 2\alpha)$ posterior credible interval** $\mathcal{I}_{1-2\alpha} = (\psi_\alpha, \psi_{1-\alpha})$ such that

$$P(\psi \in \mathcal{I}_{1-\alpha} | y^o) = 1 - 2\alpha;$$

here ψ is regarded as random and y^o as fixed.

- Likewise if there is a nuisance parameter, we require a prior density $f(\psi, \lambda)$ and compute the **marginal posterior density of ψ** ,

$$f(\psi | y^o) = \frac{\int f(y^o | \psi, \lambda)f(\psi, \lambda) d\lambda}{\iint f(y^o | \psi, \lambda)f(\psi, \lambda) d\lambda d\psi}.$$

stat.epfl.ch

Autumn 2022 – slide 23

Example

Example 2 In the satellite problem a natural prior density for $\xi \in \mathbb{R}^2$ is the **improper prior**

$$f(\xi_1, \xi_2) \propto c > 0, \quad \xi_1, \xi_2 \in \mathbb{R},$$

used to express total ignorance about ξ , and if $\mathcal{D} = \{(\xi_1, \xi_2) : \xi_1^2 + \xi_2^2 \leq \psi_{\min}\}$, then

$$P(\psi < \psi_{\min} \mid y^o) = \iint_{\mathcal{D}} \frac{(d_1 d_2)^{1/2}}{2\pi} \exp \left[-\frac{1}{2} \{d_1(y_1^o - \xi_1)^2 + d_2(y_2^o - \xi_2)^2\} \right] d\xi_1 d\xi_2$$

also equals (in polar coordinates)

$$\int_0^{\psi_{\min}} \int_0^{2\pi} \frac{(d_1 d_2)^{1/2}}{2\pi} \exp \left[-\frac{1}{2} \{d_1(y_1^o - \psi \cos \lambda)^2 + d_2(y_2^o - \psi \sin \lambda)^2\} \right] \psi d\lambda d\psi$$

corresponding to the prior

$$f(\psi, \lambda) \propto \psi > 0, \quad \psi > 0, 0 \leq \lambda < 2\pi,$$

on the polar coordinates.

How does this prior express ignorance about ξ ?

stat.epfl.ch

Autumn 2022 – slide 24

Note to Example 2

- The posterior density for ξ_1, ξ_2 is proportional to

$$c \times \frac{(d_1 d_2)^{1/2}}{2\pi} \exp \left[-\frac{1}{2} \{d_1(y_1^o - \xi_1)^2 + d_2(y_2^o - \xi_2)^2\} \right],$$

and the constant of proportionality is just the integral of this expression with respect to ξ_1, ξ_2 . By symmetry of the quadratic form in the exponent, this is just c , so the posterior density for ξ_1, ξ_2 is just

$$\frac{(d_1 d_2)^{1/2}}{2\pi} \exp \left[-\frac{1}{2} \{d_1(\xi_1 - y_1^o)^2 + d_2(\xi_2 - y_2^o)^2\} \right], \quad \xi_1, \xi_2 \in \mathbb{R},$$

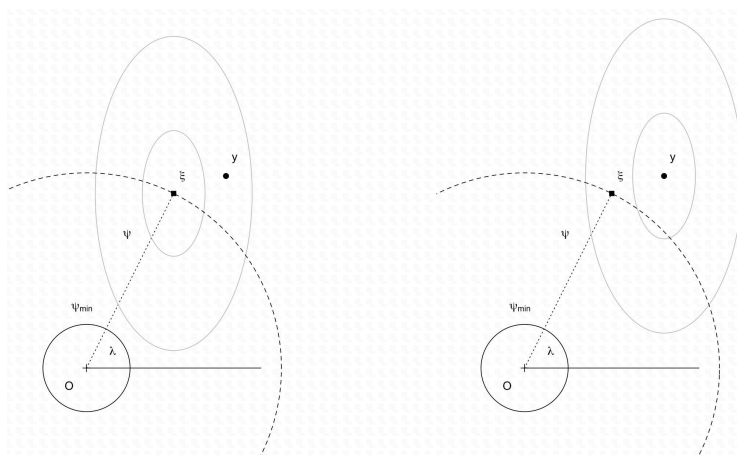
i.e., $\xi \mid y^o \sim \mathcal{N}_2(y^o, D^{-1})$. Compare the densities on the next slide to see how the sampling model for $y \mid \xi$ on the left is updated by Bayes' theorem to the posterior density of $\xi \mid y = y^o$ on the right.

- The Jacobian for changing from the Cartesian coordinates ξ_1, ξ_2 to the polar coordinates $\psi = (\xi_1^2 + \xi_2^2)^{1/2}$, $\lambda = \tan^{-1}(\xi_2/\xi_1)$ is ψ , which gives the second expression.
- As the prior is uniform the prior probability that ξ lies in any set is the ratio of the measure of that set to the measure of \mathbb{R}^2 , and thus is zero for any bounded set, and thus in particular for any disk around the origin.

stat.epfl.ch

Autumn 2022 – note 1 of slide 24

Conjunction



Left: sampling model $y \mid \xi \sim \mathcal{N}_2(\xi, D^{-1})$. Right: posterior density $\xi \mid y \sim \mathcal{N}_2(y, D^{-1})$ based on the constant prior.

stat.epfl.ch

Autumn 2022 – slide 25

Comments

- ☐ Bayesian inference
 - requires the specification of a prior distribution on unknowns, separate from the data;
 - implies that we regard prior information as equivalent to data, putting uncertainty and variation on the same footing;
 - reduces inference to computation of probabilities, so in principle is simple and direct.
- ☐ Specifying prior 'ignorance' in an objective way is problematic and can lead to paradoxes, especially in high-dimensional settings.
- ☐ (Approximate) Bayesian computation can be performed using
 - conjugate prior distributions (exact computations in simple cases),
 - integral approximations (e.g., Laplace's method),
 - deterministic methods (e.g., variational approximation),
 - simulation, especially Markov chain Monte Carlo.

stat.epfl.ch

Autumn 2022 – slide 26

Sampling theory

- Sampling theory inference treats the observed data y^o as a realisation from some model $f(y; \theta)$, and calculates probabilities using hypothetical samples from $f(y; \theta)$.
- We assess the plausibility of different values of θ by asking how well they explain y^o , often using a pivot.

Definition 3 If Y has density $f(y; \theta)$, then a **pivot (or pivotal quantity)** $Q = q(Y, \theta)$ is a function of Y and θ that has a known distribution (i.e., does not depend on θ). Often it is convenient if Q is monotone in θ for each Y .

Example 4 If $M = \max(Y_1, \dots, Y_n)$, where $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} U(0, \theta)$, show that $Q = M/\theta$ is a pivot.

stat.epfl.ch

Autumn 2022 – slide 27

Note to Example 4

Q_1 is a function of the data and the parameter, and

$$P(M \leq x) = F_Y(x)^n = (x/\theta)^n, \quad 0 < x < \theta,$$

so

$$P(Q_1 \leq q) = P(M/\theta \leq q) = P(M \leq \theta q) = (\theta q/\theta)^n = q^n, \quad 0 < q < 1.$$

which is known and does not depend on θ . Hence Q_1 is a pivot.

stat.epfl.ch

Autumn 2022 – note 1 of slide 27

Confidence intervals

Definition 5 Let $Y = (Y_1, \dots, Y_n)$ be data from a parametric statistical model with scalar parameter θ . A **confidence interval (CI)** (L, U) for θ with lower confidence bound L and upper confidence bound U is a random interval that contains θ with a specified probability, called the **(confidence) level** of the interval.

- $L = l(Y)$ and $U = u(Y)$ are statistics that can be computed from the data. They do not depend on θ .
- In a continuous setting (so $<$ gives the same probabilities as \leq), and if we write the probabilities that θ lies below and above the interval as

$$P(\theta < L) = \alpha_L, \quad P(U < \theta) = \alpha_U,$$

then (L, U) has confidence level

$$P(L \leq \theta \leq U) = 1 - P(\theta < L) - P(U < \theta) = 1 - \alpha_L - \alpha_U.$$

- Often we seek an interval with equal probabilities of not containing θ at each end, with $\alpha_L = \alpha_U = \alpha/2$, giving an **equi-tailed** $(1 - \alpha) \times 100\%$ **confidence interval**.
- We often take standard values of α , such that $1 - \alpha = 0.9, 0.95, 0.99, \dots$

stat.epfl.ch

Autumn 2022 – slide 28

Construction of a CI

- We use pivots to construct CIs:
 - we find a pivot $Q = q(Y, \theta)$ involving θ ;
 - we obtain the quantiles q_{α_U} , $q_{1-\alpha_L}$ of Q ;
 - then we transform the equation

$$P\{q_{\alpha_U} \leq q(Y, \theta) \leq q_{1-\alpha_L}\} = (1 - \alpha_L) - \alpha_U$$

into the form

$$P(L \leq \theta \leq U) = 1 - \alpha_L - \alpha_U,$$

where the bounds L , U depend on Y , $q_{1-\alpha_L}$ and q_{α_U} , but not on θ .

- Going from quantiles of Q to confidence limits for θ is known as **inverting the pivot**.
- In many cases, the bounds are of a standard form (see below).

Example 6 In Example 4, find a CI based on Q .

stat.epfl.ch

Autumn 2022 – slide 29

Note to Example 6

The p quantile of $Q_1 = M/\theta$ is given by $p = P(Q_1 \leq q_p) = q_p^n$, so $q_p = p^{1/n}$. Thus

$$P\{\alpha_U^{1/n} \leq M/\theta \leq (1 - \alpha_L)^{1/n}\} = 1 - \alpha_L - \alpha_U,$$

and a little algebra gives that

$$P\{M/(1 - \alpha_L)^{1/n} \leq \theta \leq M/\alpha_U^{1/n}\} = 1 - \alpha_L - \alpha_U,$$

so

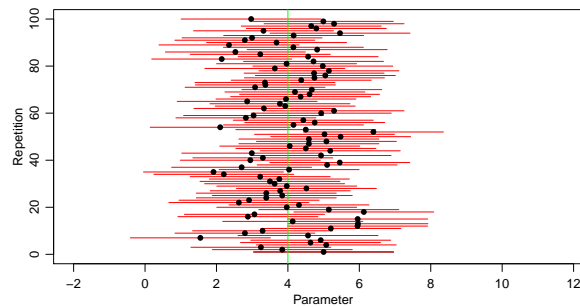
$$L = M/(1 - \alpha_L)^{1/n}, \quad U = M/\alpha_U^{1/n}.$$

stat.epfl.ch

Autumn 2022 – note 1 of slide 29

Interpretation of a CI

- ☐ (L, U) is a random interval that contains θ with probability $1 - \alpha$.
- ☐ We imagine an infinity of possible datasets from the experiment that resulted in (L, U) .
- ☐ Our particular CI is regarded as randomly chosen from the corresponding infinity of CIs.
- ☐ Although we do not know whether our particular CI contains θ , the event $\theta \in (L, U)$ has probability $1 - \alpha$ across these datasets.
- ☐ In the figure below, the parameter θ (green line) is contained (or not) in realisations of the 95% CI (red). The black points show the corresponding estimates.



stat.epfl.ch

Autumn 2022 – slide 30

One- and two-sided intervals

- ☐ A **two-sided confidence interval** (L, U) is generally used, but **one-sided confidence intervals**, $(-\infty, U)$ or (L, ∞) , are sometimes required instead.
- ☐ For one-sided CIs, we take $\alpha_U = 0$ or $\alpha_L = 0$, giving (L, ∞) or $(-\infty, U)$.
- ☐ For a one-sided $(1 - \alpha) \times 100\%$ interval, we compute a two-sided interval with $\alpha_L = \alpha_U = \alpha$, and then replace the unwanted limit by $\pm\infty$ (or another value if required in the context).

stat.epfl.ch

Autumn 2022 – slide 31

Comments

- We assume that y° is just one of many possible datasets $y \in \mathcal{S}$ that might have been generated from $f(y; \theta)$, and the probability calculations are with respect to \mathcal{S} .
- We choose the **reference set** \mathcal{S} to ensure that the probability calculation is **relevant** to the data actually observed. For example, if y° has n observations, we usually insist that every element of \mathcal{S} also has n observations.
- The repeated sampling principle ensures that (if we use an exact pivot) inferences are **calibrated**, for example, a $(1 - \alpha)$ confidence interval (L, U) satisfies

$$P(L < \theta \leq U) = 1 - \alpha,$$

for every $\theta \in \Theta$ and every $\alpha \in (0, 1)$. Hence if such an interval is used repeatedly, then the probability it does not contain θ is exactly α .

- Calibration guarantees that the procedure, if repeated, has the stated error probability, and any particular interval either does or does not contain θ .
- Bayesians object that inferences should only be based on the dataset y° actually observed, so the reference set \mathcal{S} is irrelevant.

Example 7 What would the confidence intervals look like in Example 1? How would the image on slide 30 change? What hypothetical repetitions form the reference set?

stat.epfl.ch

Autumn 2022 – slide 32

Randomisation

- To compare how **treatments** affect a **response**, they are **randomised** to experimental **units**:
 - **treatments** are clearly-defined procedures, one of which is applied to each unit;
 - a **unit** is the smallest division of the raw material such that two different units might receive two different treatments;
 - the **response** is a well-defined variable measured for each unit-treatment combination.
- Examples are agricultural trials, industrial experiments, clinical trials, ...
- The experiment is 'under the control' of the investigator, making strong inferences possible.
- Main goals of randomisation:
 - avoidance of systematic error (eliminating bias);
 - estimation of baseline variation (e.g., by use of replication and/or blocking);
 - realistic statement of uncertainty of final conclusions;
 - providing a basis for exact inferences using the randomisation distribution.

stat.epfl.ch

Autumn 2022 – slide 33

Example: Shoe data

- Shoe wear in an paired comparison experiment in which materials A (expensive) and B (cheaper) were randomly assigned to the soles of the left (L) or right (R) shoe of each of $m = 10$ boys.
- The $m = 10$ differences d_1, \dots, d_m have average $\bar{d} = 0.41$.

Boy	Material		Difference d
	A	B	
1	13.2 (L)	14.0 (R)	0.8
2	8.2 (L)	8.8 (R)	0.6
3	10.9 (R)	11.2 (L)	0.3
4	14.3 (L)	14.2 (R)	-0.1
5	10.7 (R)	11.8 (L)	1.1
6	6.6 (L)	6.4 (R)	-0.2
7	9.5 (L)	9.8 (R)	0.3
8	10.8 (L)	11.3 (R)	0.5
9	8.8 (R)	9.3 (L)	0.5
10	13.3 (L)	13.6 (R)	0.3

stat.epfl.ch

Autumn 2022 – slide 34

Example: Shoe data II

- A unit is a foot, a treatment is the type of sole, and the response is the amount of wear.
- This is **paired comparison** experiment, as there are **blocks** of two similar units, each of which is given one treatment at random, according to the scheme

Treatment for boy j	Left foot	Right foot
A	l_j	r_j
B	$\psi + l_j$	$\psi + r_j$

- We observe either $(\psi + l_j, r_j)$ or $(l_j, r_j + \psi)$ so the difference D_j of B and A for boy j is $\psi + l_j - r_j$ or $\psi + r_j - l_j$. These are equally likely, so we can write $D_j = \psi + I_j c_j$, where
 - ψ is the unknown (extra wear) effect of B compared to A,
 - $I_j = 1$ if the left shoe of boy j has material B and otherwise equals -1 , and
 - $c_j = l_j - r_j$ is the unobserved baseline difference in wear between the left and right feet of boy j .
- If we observe $(\psi + l_j, r_j)$ for boy j , then we cannot observe $(l_j, \psi + r_j)$, which is said to be **counterfactual**.

stat.epfl.ch

Autumn 2022 – slide 35

Example: Shoe data III

- There are 2^m equally-likely treatment allocations, and the observed \bar{d} is a realisation of the random variable

$$\bar{D} = \frac{1}{m} \sum_{j=1}^m D_j = \frac{1}{m} \sum_{j=1}^m \psi + I_j c_j = \psi + \frac{1}{m} \sum_{j=1}^m I_j c_j,$$

where $I_j = \pm 1$ with equal probabilities, so

$$E(I_j) = 0, \quad \text{var}(I_j) = 1.$$

- Hence $E(\bar{D}) = \psi$ and $\text{var}(\bar{D}) = m^{-2} \sum_{j=1}^m c_j^2$, which is unknown because the c_j are unknown, is estimated by (exercise)

$$S^2 = \frac{1}{m(m-1)} \sum_{j=1}^m (D_j - \bar{D})^2.$$

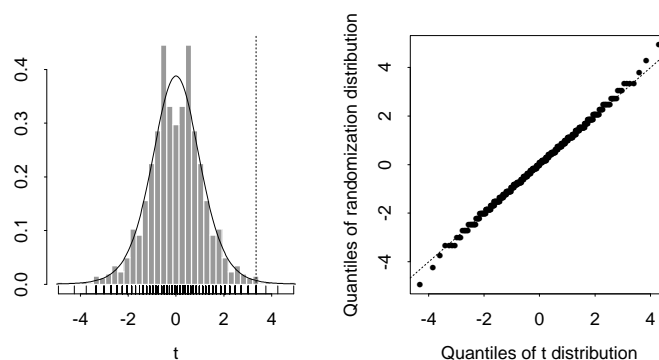
- \bar{D} and S^2 can be computed from the observed data, so the standardized quantity $Z = (\bar{D} - \psi)/S$ is an approximate pivot.
- If there was no difference between B and A (i.e., $\psi = 0$), then $T = \bar{D}/S$ would be symmetrically distributed, as positive and negative values of \bar{D} would be equally likely.

stat.epfl.ch

Autumn 2022 – slide 36

Example: Shoe data IV

Randomization distribution of $T = \bar{D}/S$ for the shoes data, i.e., setting $\psi = 0$, together with a t_9 distribution. Left: histogram and rug for the values of T , with the t_9 density overlaid; the observed value is given by the vertical dotted line. Right: probability plot of the randomization distribution against t_9 quantiles.



stat.epfl.ch

Autumn 2022 – slide 37

Comments

- **Systematic error** is reduced by randomisation,
 - but if material A had by chance been allocated to all the left feet, then we might have re-randomised;
 - we could have used a design in which A appeared on left feet exactly 5 times.
- **Baseline variation** was reduced by blocking, i.e., using two treatments for each boy, and is estimated by S^2 , based only on the observed values D_1, \dots, D_m .
- S^2 also allows a statement of **uncertainty** for \bar{D} and hence for estimates of ψ .
- If $\psi = 0$, then the observed value of \bar{D} is highly unlikely: just 3 values of \bar{D} exceed $\bar{d} = 0.41$, so if $\psi = 0$ then **exact calculation** gives

$$P(\bar{D} \geq \bar{d}) = 7/2^{10} \doteq 0.007,$$

which seems unlikely enough to suggest that $\psi > 0$.

- Normal distribution theory suggests that $Z \sim t_9$, and the QQ-plot shows that this would work well even here. The symmetry induced by randomisation justifies the widespread use of normal errors in designed experiments.

stat.epfl.ch

Autumn 2022 – slide 38

Wrapping up

- Statistical inference involves (a family of) **probability models** from which observed data are assumed to be drawn.
- These models express **variation** inherent in the data, but we also wish to express our **uncertainty** about the underlying situation.
- Uncertainty is formulated using
 - a **Bayesian approach**, which requires that 'prior information' on unknown quantities be expressed as a probability distribution, or
 - a **repeated sampling (frequentist) approach**, which invokes hypothetical repetitions of the data-generating mechanism, or
 - a **randomisation approach**, in which the model and hypothetical repetitions are controlled by the investigator.
- The last is the strongest approach, but it is not always applicable.

stat.epfl.ch

Autumn 2022 – slide 39

Likelihood

- Given observed data y thought to come from a parametric model $f_Y(y; \theta)$ for which $\theta \in \Theta$, the **likelihood** and the **log likelihood** are

$$L(\theta) = f_Y(y; \theta), \quad \ell(\theta) = \log f_Y(y; \theta), \quad \theta \in \Theta;$$

we regard these as functions of θ for fixed y . The log likelihood is often more convenient to work with because if y consists of independent observations y_1, \dots, y_n , then

$$\ell(\theta) = \log f_Y(y; \theta) = \log \prod_{j=1}^n f(y_j; \theta) = \sum_{j=1}^n \log f(y_j; \theta), \quad \theta \in \Theta,$$

so laws of large numbers and other limiting results apply directly to $n^{-1}\ell(\theta)$.

- The posterior density based on data y and prior $f(\theta)$ is proportional to $L(\theta) \times f(\theta)$.

Likelihood quantities

- The **maximum likelihood estimate (MLE)** $\hat{\theta}$ satisfies

$$\ell(\hat{\theta}) \geq \ell(\theta) \quad \text{or equivalently} \quad L(\hat{\theta}) \geq L(\theta), \quad \theta \in \Theta.$$

- Often $\hat{\theta}$ is unique and satisfies the **score (or likelihood) equation**

$$\nabla \ell(\theta) = \frac{d\ell(\theta)}{d\theta} = 0,$$

interpreted as a $d \times 1$ vector equation if θ is a $d \times 1$ vector.

- The **observed information** and **expected (Fisher) information** are defined as

$$J(\theta) = -\nabla^2 \ell(\theta) = -\frac{d^2 \ell(\theta)}{d\theta d\theta^T}, \quad I(\theta) = E\{J(\theta)\};$$

these are $d \times d$ matrices if θ has dimension d and otherwise are scalars.

- To evaluate $I(\theta)$ we replace y by the random variable Y and take expectations.

Example 8 (Poisson sample) If $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} \text{Pois}(\theta)$, for $\theta > 0$, find the log likelihood, the MLE, the observed information and the Fisher information.

Note to Example 8

- The log likelihood is

$$\ell(\theta) = \sum_{j=1}^n \log f(y_j; \theta) = \sum_{j=1}^n \log(\theta^{y_j} e^{-\theta} / y_j!) \equiv s \log \theta - n\theta, \quad \theta > 0,$$

where $s = \sum_j y_j$ and \equiv means that we have dropped additive constants from the log likelihood.

- It follows that $\ell'(\theta) = s/\theta - n$ and $\ell''(\theta) = -s/\theta^2$, so provided $s > 0$ we have $\hat{\theta} = s/n = \bar{y}$ and $J(\theta) = s/\theta^2$. If $s = 0$ then a sketch of $\ell(\theta)$ shows that $\hat{\theta} = 0 = s/n$, which is on the boundary of the parameter space. Also $\ell''(\hat{\theta})$ is only defined as a limit from the right, and equals zero.
- As $E(S) = nE(Y_j) = n\theta$, the Fisher information is computed as

$$E(S/\theta^2) = n/\theta.$$

stat.epfl.ch

Autumn 2022 – note 1 of slide 42

Invariance

- Seek invariance to (smooth) 1–1 transformations of data and/or parameter.
- If $Z = z(Y)$ is a 1–1 function of a continuous variable Y and the transformation does not depend on θ , then $f_Z(z; \theta) = f_Y\{y^{-1}(z); \theta\} |dy/dz|$, so (in an explicit notation)

$$\ell(\theta; z) = \log f_Z(z; \theta) \equiv \ell(\theta; y) = \log f_Y(y; \theta),$$

where \equiv means that an additive constant not depending on θ has been dropped — hence likelihood inference is the same whether we use Y or Z .

- Likewise a smooth 1–1 transformation from θ to $\phi(\theta)$ will give

$$\tilde{f}(y; \phi) = \tilde{f}\{y; \phi(\theta)\} = f(y; \theta),$$

where the tilde denotes the density expressed using ϕ . Clearly

$$\tilde{f}(y; \hat{\phi}) = \tilde{f}\{y; \phi(\hat{\theta})\} = f(y; \hat{\theta}), \quad J(\hat{\theta}) = \frac{\partial \phi^T}{\partial \theta} \tilde{J}(\phi) \frac{\partial \phi}{\partial \theta^T} \Big|_{\phi=\phi(\hat{\theta})},$$

so the respective maximum likelihood estimates satisfy $\hat{\phi} = \phi(\hat{\theta})$.

- If possible inferences on ψ should be invariant to **interest-respecting (or interest-preserving) transformations**

$$\psi, \lambda \mapsto \eta = \eta(\psi), \zeta = \zeta(\psi, \lambda).$$

stat.epfl.ch

Autumn 2022 – slide 43

Sufficiency

- A statistic $S = s(Y)$ is **sufficient (for θ)** under a model $f_Y(y; \theta)$ if the conditional density $f_{Y|S}(y | s; \theta)$ is independent of θ for any θ and s .
- This implies that

$$f_Y(y; \theta) = f_S(s; \theta) f_{Y|S}(y | s), \quad \ell(\theta; s) \equiv \ell(\theta; y),$$

so we can regard s as containing all the sample information about θ : if we consider Y to be generated in two steps,

- first generate S from $f_S(s; \theta)$, and
- then generate Y from $f_{Y|S}(y | s)$,

we see that if the model holds, then the second step gives no information about θ , so we could stop after the first step.

- The conditional distribution $f_{Y|S}(y | s)$ allows assessment of the model without reference to θ .

stat.epfl.ch

Autumn 2022 – slide 44

Minimal sufficiency

- If $S = s(Y)$ is sufficient and $T = t(Y)$ is any other function of Y , then (S, T) contains at least as much information as S , and is also sufficient.
- To define a ‘smallest’ sufficient function of Y , we define a **minimal sufficient statistic** to be a function of any other sufficient statistic. This is unique up to 1–1 maps.
- Below we always use minimal sufficient statistics, identified by finding the ‘lowest-dimensional function of the data’ that allows us to plot the log likelihood.

Example 9 (Poisson sample) Let $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} \text{Pois}(\theta)$, for $\theta > 0$.

- Obtain a sufficient statistic and discuss how to assess the fit of the model.
- How does the sufficient statistic change if the sample size n is a realisation of a geometric random variable with success probability θ ?

Example 10 (Location model) If $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} g(y - \theta)$, with g a known continuous density, find a sufficient statistic.

Example 11 (Uniform model) If $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} U(\theta)$, sketch the likelihood, and find a sufficient statistic and the MLE of θ .

stat.epfl.ch

Autumn 2022 – slide 45

Note to Example 9

- Here $Y = (Y_1, \dots, Y_n)$ so

$$f_Y(y; \theta) = \prod_{j=1}^n \frac{\theta^{y_j}}{y_j!} e^{-\theta} = m(y) \theta^s e^{-n\theta},$$

where $m(y) = 1/\prod y_j!$, and clearly we need only know n and $s = \sum_j y_j$, and of these only $s = s(y)$ depends on y . Hence $S = \sum_j Y_j$, which has a Poisson distribution with mean $n\theta$, looks like a suitable sufficient statistic. We have

$$f_{Y|S}(y | s; \theta) = \frac{f_Y(y; \theta)}{f_S(s; \theta)} = \frac{m(y) \theta^s e^{-n\theta}}{(n\theta)^s e^{-n\theta}/s!} = \frac{s!}{y_1! \dots y_n!} n^{-s}, \quad y \in \mathcal{Y}_s,$$

where $\mathcal{Y}_s = \{y : y_1, \dots, y_n \in \{0, \dots, s\}, \sum_j y_j = s\}$. This conditional distribution does not depend on θ , so S is sufficient for θ .

- The conditional distribution is in fact a multinomial distribution with denominator s and probability vector $(1/n, \dots, 1/n)$, and would be used for testing the fit of the Poisson model.
- When N is geometric with success probability $\theta \in (0, 1)$, the likelihood becomes

$$f(y, n; \theta) = f(y | n; \theta) f(n; \theta) = \prod_{j=1}^n \frac{\theta^{y_j}}{y_j!} e^{-\theta} \times (1 - \theta)^{n-1} \theta = m(y) (1 - \theta)^{n-1} \theta^{s+1} e^{-n\theta},$$

and now we see that we would need both n and s to sketch the likelihood. Hence the sufficient statistic is now (n, s) , which is two-dimensional, although there is only one parameter. This is because the value of n now contains information about θ : very large n will suggest that θ is small.

stat.epfl.ch

Autumn 2022 – note 1 of slide 45

Note to Example 10

- The density g is continuous, so all the y_j are distinct with probability one. The joint density is therefore

$$f(y; \theta) = \prod_{j=1}^n g(y_j - \theta) = n! \prod_{j=1}^n g(y_{(j)} - \theta), \quad y_{(1)} < \dots < y_{(n)},$$

where $s = (y_{(1)}, \dots, y_{(n)})$ are the sample order statistics. The labels on the original data are simply a permutation of the n labels on the order statistics, but the values are the same, so

$$f(y | s; \theta) = \frac{f(y; \theta)}{f(s; \theta)} = \frac{1}{n!}, \quad y \in \mathcal{Y}_s,$$

where \mathcal{Y}_s is the set of permutations of (y_1, \dots, y_n) with order statistics s ; clearly $|\mathcal{Y}_s| = n!$, because there are no ties.

- Here $|s| = n$ in general. In special cases (e.g., the normal distribution, $g = \phi$) we can find a sufficient statistic of lower dimension.

stat.epfl.ch

Autumn 2022 – note 2 of slide 45

Note to Example 11

- The density is $f(y; \theta) = \theta^{-1}I(0 < y < \theta)$, so since the observations are independent, the likelihood is

$$L(\theta) = \prod_{j=1}^n \theta^{-1}I(0 < y_j < \theta) = \theta^{-n}I(0 < y_1, \dots, y_n < \theta) = \theta^{-n}I(\theta > m), \quad \theta > 0,$$

where $m = \max(y_1, \dots, y_n)$; note that $\prod_j I(0 < y_j < \theta) = I(m < \theta)$. Viewed as a function of θ this is maximised at $\hat{\theta} = m$, which is therefore the MLE.

- Here the maximum is NOT found by differentiation of the likelihood, which is not differentiable at $\hat{\theta}$.
- We need m and n to compute the likelihood, and since n is supposed constant, the maximum m is minimal sufficient.

stat.epfl.ch

Autumn 2022 – note 3 of slide 45

Formal results

- For a less informal treatment we note that
- any statistic $T = t(Y)$ taking values $t \in \mathcal{T}$ partitions the sample space \mathcal{Y} into equivalence classes $\mathcal{C}_s = \{y' \in \mathcal{Y} : t(y') = s\}$;
 - the partition \mathcal{C}_s corresponding to T is sufficient if and only if the distribution of Y within each \mathcal{C}_s does not depend on θ ; and
 - a minimal sufficient statistic gives the coarsest possible sufficient partition.
- We use the following two results to identify the (minimal) sufficient statistic.

Theorem 12 (Factorisation) *A statistic $S = s(Y)$ is sufficient for θ in a model $f(y; \theta)$ if and only if there exist functions g and h such that*

$$f(y; \theta) = g\{s(y); \theta\} \times h(y).$$

Theorem 13 *Let Y have density $f(y; \theta)$ and let $S = s(Y)$ be such that the ratio*

$$\frac{f(y; \theta)}{f(z; \theta)}$$

is free of θ if and only if $s(y) = s(z)$. Then S is minimal sufficient for θ .

stat.epfl.ch

Autumn 2022 – slide 46

Note to Theorem 12

- The result is 'if and only if', so we need to argue in both directions.
- If S is sufficient, then the factorisation

$$f(y; \theta) = f\{s(y); \theta\} \times f(y | s) = g\{s(y); \theta\} \times h(y)$$

holds.

- To prove the converse, suppose for simplicity that Y is discrete and that there is a factorisation. Then S has density

$$f(s; \theta) = \sum_{y' \in \mathcal{Y}: s(y')=s} g\{s(y'); \theta\} h(y') = g(s; \theta) \sum_{y' \in \mathcal{Y}: s(y')=s} h(y'),$$

where the sum is in fact over $y' \in \mathcal{C}_s$. Thus the conditional density of Y given $S = s = s(y)$ is

$$f(y | s; \theta) = \frac{g\{s(y); \theta\} h(y)}{g(s; \theta) \sum_{y' \in \mathcal{C}_s} h(y')} = \frac{h(y)}{\sum_{y' \in \mathcal{C}_s} h(y')},$$

which does not depend on θ . Hence S is sufficient.

- The continuous case is similar, but the presence of a Jacobian makes the argument a bit messier.

stat.epfl.ch

Autumn 2022 – note 1 of slide 46

Note to Theorem 13

- We must show that that S is sufficient and that it is minimal.
- To show sufficiency, note that every $y \in \mathcal{Y}$ lies in an element of the partition \mathcal{C}_s generated by the possible values of S , and choose a representative dataset $y'_s \in \mathcal{C}_s$ for each s . For any y , $y'_{s(y)}$ is in the same equivalence set as y , so the ratio $f(y; \theta)/f(y'_{s(y)}; \theta)$ does not depend on θ , by the premise of the theorem. Hence

$$f(y; \theta) = f(y'_{s(y)}; \theta) \times \frac{f(y; \theta)}{f(y'_{s(y)}; \theta)} = g\{s(y); \theta\} \times h(y),$$

because $y'_{s(y)}$ is a function of $s(y)$. This factorisation shows that $S = s(Y)$ is sufficient.

- To show minimality, if $T = t(Y)$ is any other sufficient statistic the factorisation theorem gives

$$f(y; \theta) = g'\{t(y); \theta\} h'(y)$$

for some g' and h' . If two datasets y and z are such that $t(y) = t(z)$, then

$$\frac{f(y; \theta)}{f(z; \theta)} = \frac{g'\{t(y); \theta\} h'(y)}{g'\{t(z); \theta\} h'(z)} = \frac{h'(y)}{h'(z)}$$

does not depend on θ , and hence $s(y) = s(z)$. This implies that

$$\{z \in \mathcal{Y} : t(z) = t(y)\} \subset \{z \in \mathcal{Y} : s(z) = s(y)\},$$

i.e., the partition generated by the values of S is coarser than that generated by the values of T , and therefore it must be minimal.

stat.epfl.ch

Autumn 2022 – note 2 of slide 46

Exponential family models

- An elegant general theory puts many well-known distributions (Poisson, binomial, normal, ...) under the same roof.
- If $\theta \in \Theta \subset \mathbb{R}^d$, where $\dim \Theta = d$, and there exists a d -dimensional statistic $s = s(y)$ of data y and a parametrisation $\varphi \equiv \varphi(\theta)$, i.e., a 1-1 function of θ , such that

$$f(y; \theta) = m(y) \exp \{s^T \varphi - k(\varphi)\} = m(y) \exp [s^T \varphi(\theta) - k\{\varphi(\theta)\}], \quad \theta \in \Theta, y \in \mathcal{Y},$$

then this is a **regular (or full) (d, d) exponential family** of distributions, and

- the **canonical statistic** $S = s(Y)$ is minimal sufficient for θ ,
 - the **canonical parameter** is φ ,
 - the **mean parameter** $\eta = E(S; \varphi) = dk(\varphi)/d\varphi$ is obtained by differentiating the
 - the **cumulant-generating function** $k(\varphi + t) - k(\varphi)$, and
 - the **cumulant generator** $k(\cdot)$ is convex on the set $\mathcal{N} = \{\theta : \log \int e^{s^T \varphi} m(y) dy < \infty\}$.
- The cumulant-generating function is the log moment-generating function for S , i.e.,

$$K_S(t) = \log M_S(t) = k(\varphi + t) - k(\varphi), \quad t \in \mathcal{T} \subset \mathbb{R}^d,$$

where the open set \mathcal{T} contains $t = 0$. Check that $E(S) = \nabla k(\varphi)$, $\text{var}(S) = \nabla^2 k(\varphi)$.

Note on cumulant-generating functions

- The moment-generating function for the canonical statistic S of an exponential family is

$$M_S(t) = E \{ \exp(t^T S) \} = \int m(y) \exp \{s^T t + s^T \varphi - k(\varphi)\} dy,$$

and since this must equal unity when $t = 0$ we see that

$$\int m(y) \exp \{s^T \varphi\} dy = \exp \{k(\varphi)\},$$

and therefore that if it is defined,

$$M_S(t) = \int m(y) \exp \{s^T (t + \varphi) - k(\varphi)\} dy = \exp \{k(\varphi + t) - k(\varphi)\},$$

as required.

- To find the mean and variance we note that $M_S(0) = 1$,

$$\nabla M_S(t)|_{t=0} = E(S), \quad \nabla^2 M_S(t)|_{t=0} - E(S)E(S)^T = \text{var}(S),$$

and as $K_S(t) = \log M_S(t)$,

$$\nabla K_S(t) = M_S(t)^{-1} \nabla M_S(t), \quad \nabla^2 K_S(t) = M_S(t)^{-1} \nabla^2 M_S(t) - M_S(t)^{-2} \nabla M_S(t) \nabla^T M_S(t),$$

which reduce to the mean and covariance matrix when $t = 0$.

Examples

Example 14 (Poisson sample) *Are the models of Example 9 full exponential families?*

Example 15 (Satellite conjunction) *Show that the given model is an exponential family.*

stat.epfl.ch

Autumn 2022 – slide 48

Note to Example 14

- In the first model

$$f_y(y; \theta) = \prod_{j=1}^n \frac{\theta^{y_j}}{y_j!} e^{-\theta} = m(y) \exp(s \log \theta - n\theta),$$

where $m(y) = (\prod y_j)^{-1}$, $s = s(y) = \sum y_j$, $\varphi = \log \theta$, $k(\varphi) = n\theta = ne^\varphi$. This is a (1, 1) exponential family.

- In the second model we saw that

$$f(y, n; \theta) = m(y)(1-\theta)^{n-1}\theta^{s+1}e^{-n\theta} = m(y) \exp\{-n\theta + s \log \theta + n \log(1-\theta) + \log \theta - \log(1-\theta)\},$$

so we can take

$$s(y) = (n, s), \quad \varphi = (\log(1-\theta) - \theta, \log \theta), \quad k\{\varphi(\theta)\} = \log(1-\theta) - \log \theta,$$

with θ scalar but $s(y)$ of dimension 2. This is a (2, 1) curved exponential family (next slide).

stat.epfl.ch

Autumn 2022 – note 1 of slide 48

Note to Example 15

- The multivariate normal density is

$$\begin{aligned} f(y; \mu, \Omega) &= \frac{1}{(2\pi)^{n/2} |\Omega|^{1/2}} \exp\left\{-\frac{1}{2}(y - \mu)^T \Omega^{-1} (y - \mu)\right\}, \quad y \in \mathbb{R}^n \\ &= (2\pi)^{-n/2} \exp\left\{-\frac{1}{2}(y - \mu)^T \Omega^{-1} (y - \mu) - \frac{1}{2} \log |\Omega|\right\}, \end{aligned}$$

and the data y only appear in the exponent,

$$(y - \mu)^T \Omega^{-1} (y - \mu) = y^T \Omega^{-1} y - 2y^T \Omega^{-1} \mu + \mu^T \Omega^{-1} \mu.$$

- When Ω is known (as in the satellite case), the only unknown parameter is μ and we can set

$$\varphi = \Omega^{-1} \mu, \quad s(y) = y, \quad m(y) = (2\pi)^{-n/2} |\Omega|^{-1/2} e^{-y^T \Omega^{-1} y/2}, \quad k(\varphi) = \mu^T \Omega^{-1} \mu/2 = \varphi^T \Omega \varphi/2.$$

Note that $\nabla k(\varphi) = \Omega \varphi = \Omega \Omega^{-1} \mu = \mu = E(y)$ and $\nabla^2 k(\varphi) = \Omega = \text{var}(y)$, as expected.

- We could also have set $\varphi = \mu$ and $s(y) = \Omega^{-1} y$. Find $k(\varphi)$ in this case and find the corresponding mean and variance of the canonical statistic.
- If Ω is known, then y is a 1-1 transform of $s(y) = \Omega^{-1} y$, so y is sufficient also.
- In the satellite example $n = 2$ and $\Omega = D^{-1}$ is diagonal and known, so the expressions above simplify slightly.

stat.epfl.ch

Autumn 2022 – note 2 of slide 48

Exponential family models II

- When $\dim s = k > \dim \theta = d$ the model is called a **(k, d) curved exponential family**, and the $k \times 1$ vector $\varphi(\theta)$ gives a d -dimensional manifold in \mathbb{R}^k .
- The joint density of a random sample Y_1, \dots, Y_n from an exponential family is

$$\prod_{j=1}^n f(y_j; \theta) = \prod_{j=1}^n m(y_j) \exp \{s_j^T \varphi - k(\varphi)\} = m^*(y) \exp \{n\bar{s}^T \varphi - nk(\varphi)\},$$

say, where $n\bar{s} = \sum_j s(y_j)$, and this is also an exponential family, with canonical statistic $n\bar{s}$ and cumulant generator $nk(\varphi)$. Thus exponential families are **closed under sampling**.

- It is easy to check that the canonical statistic $S = \sum_{j=1}^n s(Y_j)$ of a random sample from an exponential family is minimal sufficient.
- The log likelihood $\ell(\theta) \equiv s^T \varphi - k(\varphi)$ in a full exponential family is concave as a function of φ , so the MLEs of φ and μ are given by

$$s = \nabla k(\varphi)|_{\varphi=\hat{\varphi}} = \hat{\mu},$$

and the observed and expected information quantities both equal $\nabla^2 k(\varphi)$.

stat.epfl.ch

Autumn 2022 – slide 49

Eliminating nuisance parameters

Sometimes the removal of nuisance parameters can be based on the following results.

Lemma 16 *In a statistical model $f(y; \psi, \lambda)$ let W_ψ be (minimal) sufficient for λ when ψ is regarded as fixed. Then the conditional density $f(y | w_\psi; \psi)$ depends only on ψ . This holds in particular if W_ψ does not depend on ψ .*

Lemma 17 *In a (d, d) exponential family in which $\varphi(\theta) = (\psi, \lambda)$ and $s = (t, w)$ is partitioned conformally with φ , the conditional density of T given $W = w^\circ$ is an exponential family that depends only on ψ .*

Example 18 (2×2 table) Apply Lemma 17 to the 2×2 table.

stat.epfl.ch

Autumn 2022 – slide 50

Note to Lemma 16

If ψ is regarded as fixed, then

$$f(y; \psi, \lambda) = f(w_\psi; \psi, \lambda) \times f(y | w_\psi; \psi),$$

where the rightmost term is free of λ , with logarithm

$$\log f(y; \psi, \lambda) - \log f(w_\psi; \psi, \lambda).$$

stat.epfl.ch

Autumn 2022 – note 1 of slide 50

Note to Lemma 17

In the discrete case, let \sum_o denote the sum over the set $\{y : w = w^o\}$ and note that

$$\begin{aligned} f(w^o; \psi, \lambda) &= \sum_o m^*(y) \exp \{t^T \psi + w^{oT} \lambda - k(\varphi)\} \\ &= \exp \{w^{oT} \lambda - k(\varphi)\} \sum_o m^*(y) \exp \{t^T \psi\} \end{aligned}$$

so

$$\begin{aligned} f(t \mid w^o; \psi) &= \frac{m^*(y) \exp \{t^T \psi + w^{oT} \lambda - k(\varphi)\}}{\exp \{w^{oT} \lambda - k(\varphi)\} \sum_o m^*(y) \exp(t^T \psi)} \\ &= m^*(y) \exp \left\{ t^T \psi - \log \sum_o m^*(y) \exp(t^T \psi) \right\} \\ &= m^*(y) \exp \{t^T \psi - k(\psi; w^o)\}, \end{aligned}$$

say, where the cumulant generator for the conditional density depends on w^o .

Note to Example 18

- A 2×2 table arises when m_1 individuals are allocated to a treatment and m_0 are allocated to a control. Responses from all individuals are independent and are binary with values 0/1, so the total number of successes for the control group $R_0 \sim B(m_0, \pi_0)$ is independent of those for the treatment group, $R_1 \sim B(m_1, \pi_1)$. If the parameter of interest is the difference in log odds of success. Here m_0 and m_1 are considered to be fixed, and R_0 and R_1 as random. If we write

$$\psi = \log\{\pi_1/(1 - \pi_1)\} - \log\{\pi_0/(1 - \pi_0)\} = \log\left\{\frac{\pi_1(1 - \pi_0)}{\pi_0(1 - \pi_1)}\right\}, \quad \lambda = \log\{\pi_0/(1 - \pi_0)\},$$

then we have

$$\pi_0 = \frac{e^\lambda}{1 + e^\lambda}, \quad \pi_1 = \frac{e^{\lambda+\psi}}{1 + e^{\lambda+\psi}}, \quad \psi, \lambda \in \mathbb{R}$$

and the joint density of the data reduces to

$$\binom{m_0}{r_0} \pi_0^{r_0} (1 - \pi_0)^{m_0 - r_0} \times \binom{m_1}{r_1} \pi_1^{r_1} (1 - \pi_1)^{m_1 - r_1} = \binom{m_0}{r_0} \binom{m_1}{r_1} \frac{e^{r_1\psi + (r_0 + r_1)\lambda}}{(1 + e^\lambda)^{m_0} (1 + e^{\lambda+\psi})^{m_1}},$$

which is a $(2, 2)$ exponential family with $\varphi = (\psi, \lambda)$, $s = (r_1, r_0 + r_1)$, and

$$m^*(y) = \binom{m_0}{r_0} \binom{m_1}{r_1}, \quad k(\varphi) = -m_0 \log(1 + e^\lambda) - m_1 \log(1 + e^{\lambda+\psi}).$$

- The result above implies that conditioning on $W = R_0 + R_1$ will eliminate λ , and

$$P(W = w) = \sum_{r=r_-}^{r_+} \binom{m_0}{w-r} \binom{m_1}{r} \frac{e^{r\psi + w\lambda}}{(1 + e^\lambda)^{m_0} (1 + e^{\lambda+\psi})^{m_1}},$$

where $r_- = \max(0, w - m_0)$, $r_+ = \min(w, m_1)$, and hence the conditional density of $T = R_1$ given $W = R_1 + R_0 = w$ is the **non-central hypergeometric density**

$$P(T = t \mid W = w; \psi) = \frac{\binom{m_0}{w-t} \binom{m_1}{t} e^{t\psi}}{\sum_{r=r_-}^{r_+} \binom{m_0}{w-r} \binom{m_1}{r} e^{r\psi}}, \quad t \in \{r_-, \dots, r_+\}.$$

Simple frequentist inference

- A basic frequentist recipe for inference on a parameter of interest:
- find the likelihood function for the data Y ;
 - find a sufficient statistic $S = s(Y)$ of the same dimension as θ ;
 - find a function T of S whose distribution depends only on ψ ;
 - invert the distribution of T to find confidence limits for ψ for arbitrary α ;
 - (use the conditional distribution of Y given S to assess model adequacy).

Example 19 (Exponential sample) Apply the recipe above to inference for the mean of an exponential random sample.

Example 20 (2×2 table) Apply the recipe above to the 2×2 table.

Note to Example 19

- ☐ Here $\theta = \psi$, so we can replace ψ by θ .
- ☐ The likelihood equals the joint density of y_1, \dots, y_n ,

$$L(\theta) = \prod_{j=1}^n \theta^{-1} \exp(-y_j/\theta) = \theta^{-n} \exp(-s/\theta), \quad y_1, \dots, y_n > 0, \theta > 0,$$

so $\ell(\theta) = -n \log \theta - s/\theta$ for $\theta > 0$.

- ☐ The scalar statistic $s = \sum_{j=1}^n y_j$ is clearly minimal sufficient for θ , and its distribution is gamma with parameters θ and n , i.e.,

$$f_S(s; \theta) = \frac{s^{n-1}}{\theta^n \Gamma(n)} \exp(-s/\theta), \quad s > 0, \theta > 0.$$

Note that $Q = S/\theta$ has density

$$f_Q(q; \theta) = f_S(s; \theta)|_{s=\theta q} \left| \frac{ds}{dw} \right| = \frac{q^{n-1}}{\Gamma(n)} \exp(-q), \quad q > 0,$$

so Q is a pivot, with quantiles q_p that satisfy $p = \int_0^{q_p} f_Q(u) du$, for $p \in (0, 1)$; these are just the quantiles of the gamma distribution with parameters 1 and n .

- ☐ For $\alpha \in (0, 1)$, the equation

$$1 - \alpha = P(q_{\alpha/2} < S/\theta \leq q_{1-\alpha/2}) = P(S/q_{1-\alpha/2} \leq \theta \leq S/q_{\alpha/2}),$$

yields exact $(1 - \alpha)$ confidence interval $(L, U) = (S/q_{1-\alpha/2}, S/q_{\alpha/2})$.

- ☐ The conditional distribution of Y_1, \dots, Y_n given $S = s$ is

$$\frac{f_Y(y; \theta)}{f_S(s; \theta)} = \frac{\theta^{-n} \exp(-s/\theta)}{\frac{s^{n-1}}{\theta^n \Gamma(n)} \exp(-s/\theta)} = \frac{\Gamma(n)}{s^{n-1}}, \quad (y_1, \dots, y_n)/s \in \mathcal{S}_n,$$

where \mathcal{S}_n denotes the simplex $\{(y_1, \dots, y_n) : y_1, \dots, y_n \geq 0, \sum_j y_j = 1\}$, and it follows that the joint distribution of $(Y_1/S, \dots, Y_n/S)$ is uniform on \mathcal{S}_n and does not depend on θ . Model assessment could be based on this.

Note to Example 20

- ☐ In this case

$$P(T \leq t \mid W = w; \psi) = \sum_{r=r_-}^t \frac{\binom{m_0}{w-r} \binom{m_1}{r} e^{r\psi}}{\sum_{r=r_-}^{r_+} \binom{m_0}{w-r} \binom{m_1}{r} e^{r\psi}}, \quad t \in \{r_-, \dots, r_+\},$$

and we can vary ψ to solve the equations

$$P(T \leq t \mid W = w; \psi_-) = \alpha/2, \quad P(T \leq t \mid W = w; \psi_+) = 1 - \alpha/2,$$

thus giving a $(1 - \alpha)$ confidence interval.

Ancillary statistics

Sometimes the dimension of the problem can be reduced by writing $S = (T, A)$ where $A = a(Y)$ is an **ancillary statistic**: a function of the minimal sufficient statistic whose distribution does not depend on the parameter. Then

$$f_Y(y; \theta) = f_{Y|S}(y | s) f_S(s; \theta) = f_{Y|S}(y | s) \times f_{T|A}(t | a; \theta) \times f_A(a),$$

and inference on θ is based on the second term only, with the other terms used for model-checking.

Example 21 (Location model) Show that writing

$$T = Y_{(1)}, \quad A = (0, Y_{(2)} - Y_{(1)}, \dots, Y_{(n)} - Y_{(1)}),$$

leads to inference based on the conditional density

$$f(t | a; \theta) = \frac{\prod_{j=1}^n g(t - \theta + a_j)}{\int \prod_{j=1}^n g(u + a_j) du}.$$

Note to Example 21

□ Write $y'_j = y_{(j)}$ for simplicity of notation, and note that

$$y'_1 = t, \quad y'_j = y'_1 + (y'_j - y'_1) = t + a_j, \quad j = 2, \dots, n,$$

so the Jacobian for the transformation is

$$\frac{\partial(y'_1, \dots, y'_n)}{\partial(t, a_2, \dots, a_n)} = \begin{vmatrix} 1 & 1 & 1 & \dots & 1 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{vmatrix} = 1,$$

and thus (setting $a_1 = 0$ for simplicity) the density of the **configuration** A is

$$f_A(a) = \int \prod_{j=1}^n g(t + a_j - \theta) dt = \int \prod_{j=1}^n g(u + a_j) du,$$

where we put $u = t - \theta$ in the second integral. We see that $Q = T - \theta$ is a pivot, because

$$P(Q \leq q | A = a) = P(T - \theta \leq q | A = a) = \frac{\int^q \prod_{j=1}^n g(u + a_j) du}{\int \prod_{j=1}^n g(u + a_j) du},$$

and using the quantiles $q_{\alpha/2}(a)$ and $q_{1-\alpha/2}(a)$ will give conditional confidence limits.

□ Assessment of model fit (i.e., of g) can be based on QQ plots of the values of a . We are familiar with this in regression problems.

Comments

- The essence of the recipe on slide 51 is to base an exact pivot $Q = q(Y; \psi)$ on a minimal sufficient statistic and use the **significance (or p-value) function**

$$P\{q(Y; \psi) \leq q_p\}, \quad p \in (0, 1)$$

to invert Q and thus make inference on ψ using the quantiles of Q .

- The difficulties are that:
 - finding the sufficient statistic and a function of it that depend exactly only on ψ are typically possible only in simple models;
 - finding the exact distribution of the pivot may be difficult; and
 - assessment of model fit using the conditional distribution is difficult in general.
- Nevertheless the recipe suggests how to proceed in more general settings, by basing **approximate pivots** on likelihood-based statistics, which will automatically depend on the minimal sufficient statistic.
- To finish this section we outline the two main likelihood-based approaches (to be justified later).

stat.epfl.ch

Autumn 2022 – slide 53

Maximum likelihood estimator

- In large samples from a **regular model** in which the true parameter is $\theta_{d \times 1}^0$, the maximum likelihood estimator $\hat{\theta}$ has an approximate normal distribution,

$$\hat{\theta} \sim \mathcal{N}_d \left\{ \theta^0, J(\hat{\theta})^{-1} \right\}.$$

- If we write v_{rr} for the r th diagonal element of the $d \times d$ matrix $J(\hat{\theta})^{-1}$, then the **Wald statistic (better, Wald pivot)**

$$\frac{\hat{\theta}_r - \theta_r^0}{v_{rr}^{1/2}} \sim \mathcal{N}(0, 1), \quad r = 1, \dots, d,$$

is an approximate pivot involving the r th component θ_r^0 of θ^0 , for which an approximate $(1 - 2\alpha)$ confidence interval is

$$\hat{\theta}_r \pm z_\alpha v_{rr}^{1/2}.$$

stat.epfl.ch

Autumn 2022 – slide 54

Profile log likelihood

- The Wald pivot is not invariant to interest-respecting transformations, but the **profile log likelihood** for ψ ,

$$\ell_p(\psi) = \max_{\lambda} \ell(\psi, \lambda) = \ell(\psi, \hat{\lambda}_{\psi}),$$

is invariant and provides a better (but computationally more demanding) approach to removing the nuisance parameter.

- If the model is regular and ψ is scalar then the **likelihood root**

$$r(\psi^0) = \text{sign}(\hat{\psi} - \psi^0) \left[2\{\ell_p(\hat{\psi}) - \ell_p(\psi^0)\} \right]^{1/2} \quad \dot{\sim} \quad \mathcal{N}(0, 1),$$

and therefore is an approximate pivot.

- A $(1 - \alpha)$ confidence set for the value ψ^0 generating the data is

$$\{\psi : z_{\alpha/2} \leq r(\psi) \leq z_{1-\alpha/2}\}.$$

- We will say what ‘regular’ means and justify these results later.

Discovery of the top quark (Abe et al., 1995, PRL)

Here are two extracts from the article announcing the discovery:

TABLE I. Number of lepton + jet events in the 67 pb^{-1} data sample along with the numbers of SVX tags observed and the estimated background. Based on the excess number of tags in events with ≥ 3 jets, we expect an additional 0.5 and 5 tags from $t\bar{t}$ decay in the 1- and 2-jet bins, respectively.

N_{jet}	Observed events	Observed SVX tags	Background tags expected
1	6578	40	50 ± 12
2	1026	34	21.2 ± 6.5
3	164	17	5.2 ± 1.7
≥ 4	39	10	1.5 ± 0.4

The numbers of SVX tags in the 1-jet and 2-jet samples are consistent with the expected background plus a small $t\bar{t}$ contribution (Table I and Fig. 1). However, for the $W + \geq 3$ -jet signal region, 27 tags are observed compared to a predicted background of 6.7 ± 2.1 tags [8]. The probability of the background fluctuating to ≥ 27 is calculated to be 2×10^{-5} (see Table II) using the procedure outlined in Ref. [1] (see [9]). The 27 tagged jets are in 21 events; the six events with two tagged jets can be compared with four expected for the top + background hypothesis and ≤ 1 for background alone. Figure 1 also shows the decay lifetime distribution

stat.epfl.ch

Autumn 2022 – slide 57

Performing a test

- ☐ There's a **null hypothesis** to be tested:

H_0 : the top quark does not exist.

This seems counter-intuitive, but as one cannot prove a hypothesis, we attempt to refute its opposite — '**proof by (stochastic) contradiction**'.

- ☐ We obtain data, $y_{\text{obs}} = 27$ events on the 3-jet, 4-jet, ... channels.
- ☐ We compare y_{obs} with its distribution P_0 supposing that H_0 is true.
- ☐ Here P_0 is $\text{Pois}(\lambda_0 = 6.7)$ and represents the baseline noise under H_0 .
- ☐ We compute the **P-value**

$$p_{\text{obs}} = P_0(Y \geq y_{\text{obs}}) = \sum_{y=y_{\text{obs}}}^{\infty} \frac{\lambda_0^y}{y!} e^{-\lambda_0} = 3 \times 10^{-9},$$

so

- either H_0 is true but a (very) rare event has occurred,
 - or H_0 is false and the top quark exists.
- ☐ Abe et al. announced a discovery, but if they had found $p_{\text{obs}} \approx 0.001$, maybe they would have decided that H_0 could not (yet) be rejected, and not published their work.

stat.epfl.ch

Autumn 2022 – slide 58

Industrial fraud?

DETAIL WEIGHT NOTE

No.	10	20	30	40	50	60	70	80	90	100	No.	Total
1	263	286	264	281	265	265	251				10	
2	266	266	266	284	265	261	265	264			20	
3	261	267	261	281	262	264	264	264			30	
4	261	261	261	262	263	263	261	265			40	
5	286	262	264	266	260	265	265	261	265		50	
6	287	262	264	262	261	264	265				60	
7	261	266	266	266	261	265	262				70	
8	265	264	263	261	262	262	262				80	
9	263	265	261	263	264	265	268				90	
10	265	266	266	266	262	262	265				100	
TOTAL	265	266	265	266	262	265	265	261	265			265
REDUCTIONS												
GROSS TOTAL												265

- ☐ $n = 92$ weighings of sacks upon the delivery of a commodity C:

261 289 291 265 281 291 285 283 280 261 263 281 291 289 280
 292 291 282 280 281 291 282 280 286 291 283 282 291 293 291
 300 302 285 281 289 281 282 261 282 291 291 282 280 261 283
 291 281 246 249 252 253 241 281 282 280 261 265 281 283 280
 242 260 281 261 281 282 280 241 249 251 281 273 281 261 281
 282 260 281 282 241 245 253 260 261 281 280 261 265 281 241
 260 241

- ☐ Their last digits are

0 1 2 3 4 5 6 7 8 9
 14 42 14 9 0 6 2 0 0 5

- ☐ How can we tell if fraud has taken place?

Benford's law

Definition 22 For $x \in \mathbb{R}$, let $d(x, j)$ denote the ' **j th significant digit function (base 10)**', so $d(31.4, 1) = 3$, $d(0.314, 2) = 1$ and $d(314, 3) = 4$.

Definition 23 If $x \in \mathbb{R}$ and $D_j = d(x, j)$, for $j = 1, 2, \dots$, then (discarding any leading zeros) the D_j follow **Benford's law** if

$$P(D_1 = d_1, D_2 = d_2, \dots, D_k = d_k) = \log_{10} \left\{ 1 + \left(\sum_{j=1}^k d_j \times 10^{k-j} \right)^{-1} \right\}, \quad d_j \in \{0, \dots, 9\}.$$

- ☐ For example, $P(D_1 = 3, D_2 = 1, D_3 = 4) = \log_{10}\{1 + (314)^{-1}\} \approx 0.0014$.
- ☐ This is an excellent model for the distribution of all sorts of digits.
- ☐ Frequencies (%) of the last digits D_3 for three-digit integers.

Digit	0	1	2	3	4	5	6	7	8	9
Uniform	10	10	10	10	10	10	10	10	10	10
Benford	10.178	10.137	10.097	10.057	10.017	9.978	9.940	9.901	9.864	9.826

- ☐ Apparently the last digits of the weighings should be approximately uniform.
- ☐ To detect forged deliveries we test the null hypothesis

H_0 : the last digits of the weighings are uniformly distributed on $0, \dots, 9$.

Pearson's statistic

Definition 24 If O_1, \dots, O_k are the numbers of observations from a random sample of size n falling in categories $1, \dots, k$, where

$$E(O_i) = E_i > 0, \quad i = 1, \dots, k, \quad \sum_{i=1}^k E_i = n,$$

then **Pearson's statistic (aka the ' χ^2 statistic')** is

$$T = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}.$$

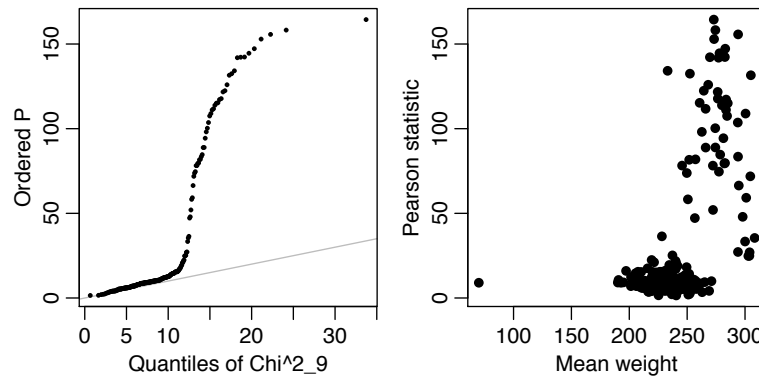
- ☐ If O_1, \dots, O_k are multinomially distributed with total n and probabilities $p_1 = E_1/n, \dots, p_k = E_k/n$, then $T \sim \chi_{k-1}^2$ (approximation OK if average $E_i \geq 5$).
- ☐ We use T to check whether data O_1, \dots, O_k agree with specified probabilities p_1, \dots, p_k .
- ☐ For the original dataset we found $t_{\text{obs}} = 158.2$ and hence

$$p_{\text{obs}} = P_0(T > t_{\text{obs}}) \doteq P(\chi_9^2 \geq 158.2) \doteq 0,$$

which is essentially impossible for uniformly distributed digits.

For 250 'deliveries' ...

- Left: Q-Q plot of t_{obs} for 250 different deliveries.
- Right: mean weights for deliveries and values of t_{obs} . Deliveries with the largest values of t_{obs} tend also to be heavier, also suggestive of fraud.



stat.epfl.ch

Autumn 2022 – slide 62

Elements of a test

- A **null hypothesis** H_0 to be tested.
- A **test statistic** T , large values of which will suggest that H_0 is false, and with observed value t_{obs} .
- A **P-value**

$$p_{\text{obs}} = P_0(T \geq t_{\text{obs}}),$$

where the **null distribution** $P_0(\cdot)$ denotes a probability computed under H_0 .

- The smaller p_{obs} is, the more we doubt that H_0 is true.
- If H_0 is true, then we can consider that p_{obs} is a realisation of a uniform random variable $P \sim U(0, 1)$, and then

$$P_0(P \leq p_{\text{obs}}) = p_{\text{obs}}.$$

- If I decide that H_0 is false, when in fact it is true, then I make an error whose probability under H_0 is exactly p_{obs} — so my uncertainty is quantified, because I know the probability of declaring a **“false positive”**.

stat.epfl.ch

Autumn 2022 – slide 63

Note: Why is a P-value uniform?

- Let T be a test statistic whose distribution is $F_0(t)$ when the null hypothesis is true. Then the corresponding P-value is

$$P_0(T \geq t_{\text{obs}}) = 1 - F_0(t_{\text{obs}}),$$

and if the value of t_{obs} is a realisation of T_{obs} (because the null hypothesis is true), then we can write the random value of p_{obs} seen in repetitions of the experiment as

$$P_{\text{obs}} = 1 - F_0(T_{\text{obs}}),$$

or equivalently $T_{\text{obs}} = F_0^{-1}(1 - P_{\text{obs}})$. Hence for $x \in [0, 1]$,

$$\begin{aligned} P_0(P_{\text{obs}} \leq x) &= P_0\{1 - F_0(T_{\text{obs}}) \leq x\} \\ &= P_0\{1 - x \leq F_0(T_{\text{obs}})\} \\ &= P_0\{T_{\text{obs}} \geq F_0^{-1}(1 - x)\} \\ &= 1 - F_0\{F_0^{-1}(1 - x)\} \\ &= x, \end{aligned}$$

which shows that $P_{\text{obs}} \sim U(0, 1)$.

- The above proof works for any continuous T_{obs} , but is only approximate if T_{obs} is discrete (e.g., has a Poisson distribution). In such cases P_{obs} can only take a finite or countable number of values known as the **achievable significance levels**.

stat.epfl.ch

Autumn 2022 – note 1 of slide 63

Some asides

- If we say that a hypothesis is **true**, we mean ‘it is reasonable to proceed as if the hypothesis was true’ — any model is an idealisation, so it cannot be exactly ‘true’.
- If we have a **discrete test statistic**, p_{obs} has at most a countable number of ‘achievable significance levels’. This is only problematic when comparing tests, though randomisation has (unfortunately) sometimes been proposed to overcome it.
- We may consider a **two-sided test**, with both unusually large and unusually small values of T of interest. We can then define

$$p_+ = P_0(T \geq t_{\text{obs}}), \quad p_- = P_0(T \leq t_{\text{obs}}), \quad p_{\text{obs}} = 2 \min(p_-, p_+),$$

so $p_- + p_+ = 1 + P_0(T = t_{\text{obs}})$, which equals 1 unless T is discrete;

- We sometimes avoid minor problems due to discreteness by computing ‘**continuity-corrected**’ P-values

$$p_+ = \sum_{t > t_{\text{obs}}} P_0(T = t) + \frac{1}{2} P_0(T = t_{\text{obs}}), \quad p_- = \sum_{t < t_{\text{obs}}} P_0(T = t) + \frac{1}{2} P_0(T = t_{\text{obs}}).$$

- The top quark and fraud examples illustrate **pure significance tests**, where the situation if H_0 is false is not explicitly considered. We look at the effect of alternatives now.

stat.epfl.ch

Autumn 2022 – slide 64

Testing as decision-making

Formulate testing as deciding between two hypotheses (Neyman–Pearson approach):

- the **null hypothesis** H_0 , which represents a baseline situation;
- the **alternative hypothesis** H_1 , which represents what happens if H_0 is false.
- We choose H_1 and ‘reject’ H_0 if p_{obs} is lower than some $\alpha \in (0, 1)$.
- For given α we partition the sample space \mathcal{Y} as

$$\mathcal{Y}_0 = \{y \in \mathcal{Y} : p_{\text{obs}}(y) > \alpha\}, \quad \mathcal{Y}_1 = \{y \in \mathcal{Y} : p_{\text{obs}}(y) \leq \alpha\},$$

where the notation $p_{\text{obs}}(y)$ indicates that the P-value depends on the data, or equivalently

$$\mathcal{Y}_0 = \{y \in \mathcal{Y} : t(y) < t_{1-\alpha}\}, \quad \mathcal{Y}_1 = \{y \in \mathcal{Y} : t(y) \geq t_{1-\alpha}\},$$

where t_p denotes the p quantile of the test statistic $T = t(Y)$ under H_0

- We call \mathcal{Y}_1 the **size α critical region** of the test, and we reject H_0 in favour of H_1 if $Y \in \mathcal{Y}_1$, or equivalently if the test statistic exceeds the **size α critical point** $t_{1-\alpha}$.
- Critical regions of different sizes for the same test should be nested, i.e., (in an obvious notation) if $\alpha' > \alpha$, then

$$\mathcal{Y}_1^\alpha \subset \mathcal{Y}_1^{\alpha'} \quad \text{and} \quad t_{1-\alpha} > t_{1-\alpha'}.$$

stat.epfl.ch

Autumn 2022 – slide 65

False positives and negatives

		Decision	
		Accept H_0	Reject H_0
State of Nature	H_0 true	Correct choice (True negative)	Type I Error (False positive)
	H_1 true	Type II Error (False negative)	Correct choice (True positive)

- We can make two sorts of wrong decisions:
 - Type I error (false positive)**: H_0 is true, but we wrongly reject it (and choose H_1);
 - Type II error (false negative)**: H_1 is true, but we wrongly choose H_0 .
- Notice that the consequences of bad decisions are not really taken into account in this framework.
- Statistics books and papers call
 - the **Type I error/false positive probability** the **size** $\alpha = P_0(Y \in \mathcal{Y}_1)$, and
 - the **true positive probability** the **power** $\beta = P_1(Y \in \mathcal{Y}_1)$.

Example 25 If $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, with σ^2 known, and $H_0 : \mu = \mu_0$ and $H_1 : \mu = \mu_1$, find the Type II error as a function of the Type I error.

stat.epfl.ch

Autumn 2022 – slide 66

Note to Example 31

- The minimal sufficient statistic for the normal model with both parameters unknown is (\bar{Y}, S^2) , and it is easy to check that if σ^2 is known the minimal sufficient statistic reduces to \bar{Y} , which has a $\mathcal{N}(\mu_0, \sigma^2/n)$ distribution under H_0 . Hence we take the test statistic T to be \bar{Y} , and $\mathcal{Y} = \mathbb{R}^n$.
- If $\mu_1 > \mu_0$, then clearly we will take

$$\mathcal{Y}_0 = \{y : \bar{y} < t_{1-\alpha}\}, \quad \mathcal{Y}_1 = \{y : \bar{y} \geq t_{1-\alpha}\};$$

this can be justified using the Neyman–Pearson lemma (below). Now

$$P_0(Y \in \mathcal{Y}_0) = P_0(\bar{Y} < t_{1-\alpha}) = P_0\{\sqrt{n}(\bar{Y} - \mu_0)/\sigma < \sqrt{n}(t_{1-\alpha} - \mu_0)/\sigma\} = \Phi\{\sqrt{n}(t_{1-\alpha} - \mu_0)/\sigma\},$$

because $Z = \sqrt{n}(\bar{Y} - \mu_0)/\sigma \sim \mathcal{N}(0, 1)$ under H_0 , and for this probability to equal $1 - \alpha$ we must take $t_{1-\alpha} = \mu_0 + \sigma n^{-1/2} z_{1-\alpha}$; this gives Type I error α .

- Note that although the form of \mathcal{Y}_0 above was determined by H_1 , the value of $t_{1-\alpha}$ is given by calculations under H_0 .
- $Z = \sqrt{n}(\bar{Y} - \mu_1)/\sigma \sim \mathcal{N}(0, 1)$ under H_1 , so the Type II error is

$$\begin{aligned} P_1(Y \in \mathcal{Y}_0) &= P_1(\bar{Y} < t_{1-\alpha}) \\ &= P_1(\bar{Y} < \mu_0 + \sigma n^{-1/2} z_{1-\alpha}) \\ &= P_1\{\sqrt{n}(\bar{Y} - \mu_1)/\sigma < \sqrt{n}(\mu_0 + \sigma n^{-1/2} z_{1-\alpha} - \mu_1)/\sigma\} \\ &= \Phi(z_{1-\alpha} - \delta), \end{aligned}$$

where $\delta = n^{1/2}(\mu_1 - \mu_0)/\sigma$. Hence the Type II error equals $1 - \alpha$ when $\mu_1 = \mu_0$ and decreases as a function of δ . We would expect this, because as μ_1 increases, the distribution of \bar{Y} under H_1 shifts to the right and we are less likely to make a false negative error.

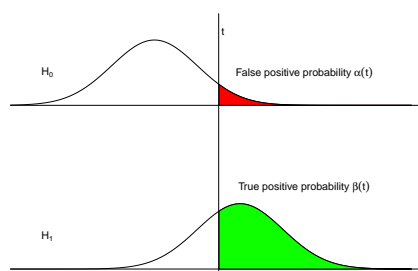
True and false positives: Example

- It is traditional to fix α and choose T (or equivalently \mathcal{Y}_1) to maximise β , but usually more informative to consider $P_0(T \geq t)$ and $P_1(T \geq t)$ as functions of t .
- In Example 31 we would
 - reject H_0 incorrectly (**false positive**) with probability

$$\alpha(t) = P_0(T \geq t) = 1 - \Phi\{n^{1/2}(t - \mu_0)/\sigma\},$$

- reject H_0 correctly (**true positive**) with probability

$$\beta(t) = P_1(T \geq t) = 1 - \Phi\{n^{1/2}(t - \mu_0)/\sigma - \delta\}.$$



ROC curve

Definition 26 The **receiver operating characteristic (ROC) curve** of a test plots $\beta(t)$ against $\alpha(t)$ as t varies, i.e., it shows $(P_0(T \geq t), P_1(T > t))$, when $t \in \mathbb{R}$.

- As μ increases, it becomes easier to detect when H_0 is false, because the densities under H_0 and H_1 become more separated, and the ROC curve moves 'further north-west'.
- When H_0 and H_1 are the same, i.e., $\mu = 0$, then the curve lies on the diagonal. Then the hypotheses cannot be distinguished.
- A common summary measure of the overall quality of a test is the **area under the curve**,

$$\text{AUC} = \int_0^1 \beta(\alpha) d\alpha,$$

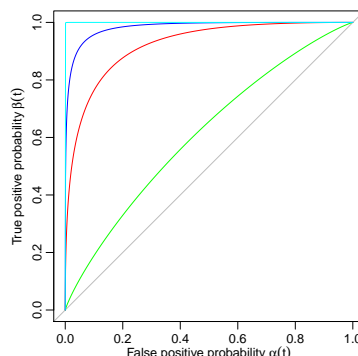
which ranges between 0.5 for a useless test and 1.0 for a perfect test.

Example

- In Example 31 $\alpha(t) = 1 - \Phi\{n^{1/2}(t - \mu_0)/\sigma\}$ and $\beta(t) = 1 - \Phi\{n^{1/2}(t - \mu_0)/\sigma - \delta\}$, so equivalently we graph

$$\beta(t) = 1 - \Phi(-z_{1-\alpha} - \delta) = \Phi(\delta + z_\alpha) \equiv \beta(\alpha) \text{ against } \alpha \in (0, 1).$$

- Here is the ROC curve with $\mu = 2$ (in red). Also shown are curves for $\mu = 0, 0.4, 3, 6$. Which is which?



Neyman–Pearson lemma

Definition 27 A **simple hypothesis** entirely fixes the distribution of the data Y , whereas a **composite hypothesis** does not fix the distribution of Y .

Definition 28 The **critical region** of a hypothesis test is the subset \mathcal{Y}_1 of the sample space \mathcal{Y} for which $Y \in \mathcal{Y}_1$ implies that the null hypothesis is rejected.

We aim to choose \mathcal{Y}_1 to maximise the power of the test for a given size, i.e., such that $P_1(Y \in \mathcal{Y}_1)$ is the largest possible such that $P_0(Y \in \mathcal{Y}_1) = \alpha$.

Lemma 29 (Neyman–Pearson) Let $f_0(y)$, $f_1(y)$ be the densities of Y under simple null and alternative hypotheses. Then if it exists, the set

$$\mathcal{Y}_1 = \{y \in \mathcal{Y} : f_1(y)/f_0(y) > t\}$$

such that $P_0(Y \in \mathcal{Y}_1) = \alpha$ maximises $P_1(Y \in \mathcal{Y}_1)$ amongst all \mathcal{Y}'_1 for which $P_0(Y \in \mathcal{Y}'_1) \leq \alpha$. Thus the test of size α with maximal power rejects H_0 when $Y \in \mathcal{Y}_1$.

Example 30 Construct an optimal test for testing $H_0 : \varphi = \varphi_0$ against $H_1 : \varphi = \varphi_1$ based on a random sample from a canonical exponential family.

Note to Lemma 29

Suppose that a region \mathcal{Y}_1 such that $P_0(Y \in \mathcal{Y}_1) = \alpha$ exists and let \mathcal{Y}'_1 be any other critical region of size α or less. Then for any density f ,

$$\int_{\mathcal{Y}_1} f(y) dy - \int_{\mathcal{Y}'_1} f(y) dy, \quad (2)$$

equals

$$\int_{\mathcal{Y}_1 \cap \mathcal{Y}'_1} f(y) dy + \int_{\mathcal{Y}_1 \cap \mathcal{Y}'_0} f(y) dy - \int_{\mathcal{Y}'_1 \cap \mathcal{Y}_1} f(y) dy - \int_{\mathcal{Y}'_1 \cap \mathcal{Y}'_0} f(y) dy,$$

where \mathcal{Y}_0 and \mathcal{Y}'_0 are the complements of \mathcal{Y}_1 and \mathcal{Y}'_1 in the sample space, and this is

$$\int_{\mathcal{Y}_1 \cap \mathcal{Y}'_0} f(y) dy - \int_{\mathcal{Y}'_1 \cap \mathcal{Y}'_0} f(y) dy. \quad (3)$$

If $f = f_0$, then (2) is non-negative, because \mathcal{Y}' has size at most that of \mathcal{Y}_1 , so (3) is also non-negative, giving

$$t \int_{\mathcal{Y}_1 \cap \mathcal{Y}'_0} f_0(y) dy \geq t \int_{\mathcal{Y}'_1 \cap \mathcal{Y}'_0} f_0(y) dy$$

for $t \geq 0$. But $f_1(y) > t f_0(y)$ for $y \in \mathcal{Y}_1$, and $t f_0(y) \geq f_1(y)$ for $y \in \mathcal{Y}_0$, so

$$\int_{\mathcal{Y}_1 \cap \mathcal{Y}'_0} f_1(y) dy \geq \int_{\mathcal{Y}'_1 \cap \mathcal{Y}'_0} f_1(y) dy.$$

On adding $\int_{\mathcal{Y}_1 \cap \mathcal{Y}'_1} f_1(y) dy$ to both sides we see that the power of \mathcal{Y}_1 is at least that of \mathcal{Y}'_1 , as required.

Note to Example 30

- The likelihood ratio is

$$\frac{f_1(y)}{f_0(y)} = \frac{m^*(y) \exp\{\varphi_1 s^* - nk(\varphi_1)\}}{m^*(y) \exp\{\varphi_0 s^* - nk(\varphi_0)\}} = \exp\{(\varphi_1 - \varphi_0)s^* + nk(\varphi_0) - nk(\varphi_1)\},$$

say, where $s^* = \sum_{j=1}^n s(y_j)$, so

$$\mathcal{Y}_1 = \{y : f_1(y)/f_0(y) > t\} = \{y : (\varphi_1 - \varphi_0)s^* + nk(\varphi_0) - nk(\varphi_1) > \log t\},$$

and if $\varphi_1 > \varphi_0$ then

$$\mathcal{Y}_1 = \{y : s^* > [\log t + nk(\varphi_1) - nk(\varphi_0)]/(\varphi_1 - \varphi_0)\},$$

This gives the form of \mathcal{Y}_1 and we should choose t so that $P_0(Y \in \mathcal{Y}_1) = \alpha$, or equivalently s_α so that (in the continuous case)

$$P_0(S^* > s_\alpha) = \int_{s_\alpha}^{\infty} f(s; \varphi_0) ds = \alpha.$$

We saw such a calculation already in Example 31 for normal data with known σ^2 and $\varphi_1 = \mu_1/\sigma^2 > \varphi_0 = \mu_0/\sigma^2$.

- If $\varphi_1 < \varphi_0$, then division by $\varphi_1 - \varphi_0 < 0$ leads to

$$\mathcal{Y}_1^* = \{y : s^* < [\log t + nk(\varphi_1) - nk(\varphi_0)]/(\varphi_1 - \varphi_0)\}.$$

- The Neyman–Pearson lemma tell us that \mathcal{Y}_1 gives a most powerful test, but as it does not depend on the value of φ , this test is **uniformly most powerful** for all $\varphi > \varphi_0$, and likewise \mathcal{Y}_1^* is **uniformly most powerful** for $\varphi_1 < \varphi_0$.

Exact and inexact tests

- Above we saw that $P \sim U(0, 1)$ under the null hypothesis, exactly in continuous cases and approximately in discrete cases.
- If the null distribution of the test statistic is estimated, we have $P \dot{\sim} U(0, 1)$ only.
- For example, if the true parameter is $\theta = (\psi_0, \lambda_0)$ and $H_0 : \psi = \psi_0$, then the P-value is

$$p_{\text{obs}} = P_0(T \geq t_{\text{obs}}) = P(T \geq t_{\text{obs}}; \psi_0, \lambda_0),$$

which we estimate by

$$\hat{p}_{\text{obs}} = P(T \geq t_{\text{obs}}; \psi_0, \hat{\lambda}_0),$$

where $\hat{\lambda}_0$ is the estimate of λ under H_0 .

- Exact tests, with $P \sim U(0, 1)$, can sometimes be obtained by using a pivot whose distribution is invariant to λ , or by removing λ by conditioning.

Example 31 If $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, show that the distribution of $T = (\bar{Y} - \mu)/\sqrt{S^2/n}$ is invariant to σ^2 .

Example 32 Find an exact test on a canonical parameter in a logistic regression model.

Note to Example 31

In this case the variables $\bar{Y} \sim \mathcal{N}(\mu, \sigma^2/n)$ and $(n-1)S^2/\sigma^2 \sim \chi_{n-1}^2$ are independent and we can write $\bar{Y} \stackrel{D}{=} \mu + \sigma n^{-1/2}Z$ and $S^2 \stackrel{D}{=} \sigma^2 V/(n-1)$, where $Z \sim \mathcal{N}(0,1)$ and $V \sim \chi_{n-1}^2$ are independent. Hence

$$T = \frac{\bar{Y} - \mu}{\sqrt{S^2/n}} \stackrel{D}{=} \frac{\mu + \sigma Z/n^{1/2} - \mu}{[\sigma^2 V/\{n(n-1)\}]^{1/2}} \stackrel{D}{=} \frac{Z}{\sqrt{V/(n-1)}} \sim t_{n-1},$$

is pivotal and thus allows tests on μ without reference to σ^2 .

Note to Example 32

- In a logistic regression model we have independent binary variables Y_1, \dots, Y_n each with density

$$P(Y_j = y_j; \beta) = \pi_j^{y_j} (1 - \pi_j)^{1-y_j} = \left(\frac{e^{x_j^T \beta}}{1 + e^{x_j^T \beta}} \right)^{y_j} \left(\frac{1}{1 + e^{x_j^T \beta}} \right)^{1-y_j} = \frac{e^{y_j x_j^T \beta}}{1 + e^{x_j^T \beta}},$$

for $y_j \in \{0, 1\}$, known covariate vectors $X_j \in \mathbb{R}^d$ and parameter $\beta \in \mathbb{R}^d$.

- The corresponding log likelihood is

$$\ell(\beta) = \sum_{j=1}^n \left\{ y_j x_j^T \beta - \log(1 + e^{x_j^T \beta}) \right\} = y^T X \beta - \sum_{j=1}^n \log(1 + e^{x_j^T \beta}), \quad \beta \in \mathbb{R}^d.$$

This is a (d, d) exponential family with canonical statistic $S = X^T y$, canonical parameter $\varphi = \beta$, and cumulant generator $k(\varphi) = \sum_{j=1}^n \log(1 + e^{x_j^T \varphi})$.

- Hence Lemma 17 implies that if $\varphi = (\psi, \lambda)$ and $S = (T, W) = (X_1^T y, X_2^T y)$, where X_1 is $n \times 1$ and X_2 is $n \times (d-1)$, an exact test on ψ is obtained from the conditional distribution

$$P(T = t \mid W = w^o; \psi) = \frac{e^{t\psi}}{\sum_{y' \in \mathcal{S}_{w^o}} e^{X_1^T y' \psi}},$$

where $\mathcal{S}_w = \{(y'_1, \dots, y'_n) : X_2^T y' = w^o\}$, with $w^o = X_2^T y^o$ and y^o respectively the observed data and the observed value of W .

- Calculation of this conditional density in applications may be awkward, but excellent approximations are available.

Interpretation of P-values

- Be careful about interpretation:
 - p_{obs} is a one-number summary of whether data are consistent with H_0 ;
 - it is NOT the probability that H_0 is true (require prior probabilities on H_0 and H_1);
 - even a tiny p_{obs} can support H_0 better than H_1 (consider $t_{\text{obs}} = 3$ when $T \sim \mathcal{N}(\mu, 1)$ with $\mu_0 = 0, \mu_1 = 10$);
 - the power depends on analogues of $\delta = n^{1/2}(\mu_1 - \mu_0)/\sigma$, where n is the **sample size**, $\mu_1 - \mu_0$ is the **effect size**, and σ is the **precision**, so
 - ▷ even a tiny (practically irrelevant) effect size can be detected with very large n ;
 - ▷ conversely a practically important effect might be undetectable if n is small;
 - ▷ i.e., 'statistical significance' \neq 'subject-matter importance'!
- A confidence interval, or estimate and its standard error, is often more informative.
- Hypothesis testing is often applied by rote — in some medical journals no statement is complete without an accompanying ' $P < 0.05$ ' — and is sometimes regarded as controversial, with certain journals now refusing to publish tests and P-values.
- The replication crisis is partly due to abuse of hypothesis testing, e.g., by not correcting for multiple tests, by formulating hypotheses in light of the data, ...

stat.epfl.ch

Autumn 2022 – slide 72

Contexts of testing

- It is unwise to be too categorical about testing, because of its different uses:
 - testing a clear hypothesis of scientific interest (e.g., top quark);
 - goodness of fit of a model (e.g., industrial fraud);
 - decision-making with a clearly-specific alternative (e.g., covid testing);
 - model simplification if null hypothesis true;
 - 'dividing hypothesis' used to split parameter space into different sets with sharply different interpretations;
 - as a technical device for generating confidence intervals;
 - to flag which of many null similar hypotheses might be false.

Example 33 *The generalized Pareto distribution, with survival function*

$$P(X > x) = \begin{cases} (1 + \xi x/\sigma)_+^{-1/\xi}, & \xi \neq 0, \\ \exp(-x/\sigma), & \xi = 0, \end{cases}$$

simplifies if $\xi = 0$, and has finite upper support point $x_+ = -\sigma/\xi$ when $\xi < 0$ but $x_+ = \infty$ when $\xi \geq 0$. Here $H_0 : \xi = 0$ is both a simplifying and a dividing hypothesis, of interest when the distribution is fitted to data on supercentenarians.

stat.epfl.ch

Autumn 2022 – slide 73

Generating confidence intervals

- Using a statistic T with observed value t_{obs} to test the null hypothesis $H_0 : \psi = \psi_0$ for a scalar parameter ψ gives P-value

$$p_{\text{obs}} = p(\psi_0) = P(T \geq t_{\text{obs}}; \psi_0),$$

and we regard ψ_0 as incompatible with the data if p_{obs} is too small.

- Recall that the corresponding random variable $P_{\text{obs}} \sim U(0, 1)$ under H_0 . Hence we can regard values of ψ for which the **P-value (or significance) function**

$$p(\psi) = P(T \geq t_{\text{obs}}; \psi)$$

is too extreme as incompatible with the data, leading to the (two-sided) $(1 - \alpha)$ confidence set

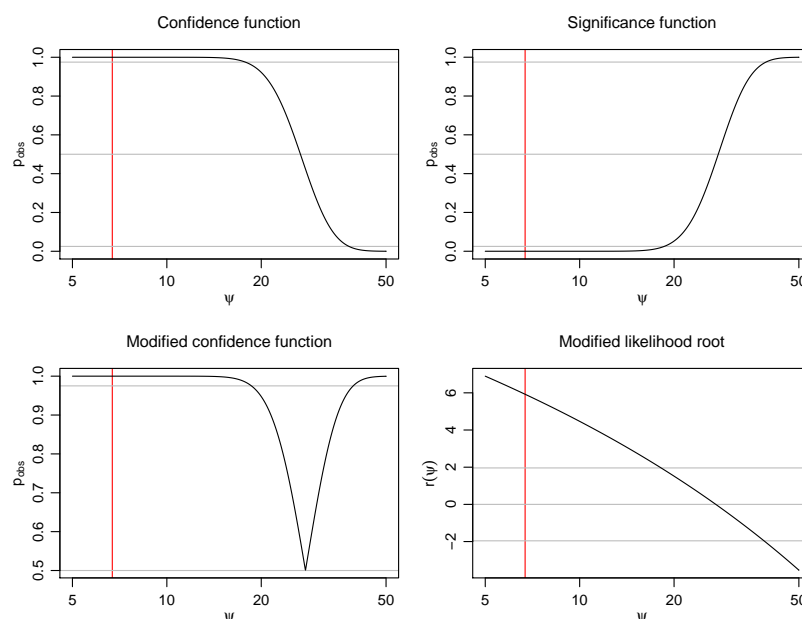
$$\{\psi : \alpha/2 \leq P(T \geq t_{\text{obs}}; \psi) \leq 1 - \alpha/2\}.$$

- Related functions include
 - the **confidence function** $1 - p(\psi)$;
 - the **modified confidence function** $\max\{p(\psi), 1 - p(\psi)\}$; and
 - a **pivot function** showing a (usually standard normal) pivot as a function of ψ .

stat.epfl.ch

Autumn 2022 – slide 74

Significance and related functions



stat.epfl.ch

Autumn 2022 – slide 75

Multiple testing

- Often require tests of several, even very many, hypotheses:
 - comparison of responses for several treatment groups with the same control group;
 - checking for a change in a series of observations;
 - screening genomic data for effects of many genes on a response.
- There are null hypotheses H_1, \dots, H_m , of which
 - m_0 are true, indexed by an unknown set \mathcal{I} ,
 - $m_1 = m - m_0$ are false, and
 - the **global null hypothesis** is $H_0 = H_1 \cap \dots \cap H_m$.
- We apply some testing procedure and declare R hypotheses to be significant, of which FP are false positives and TP are true positives. Only R and m are known.

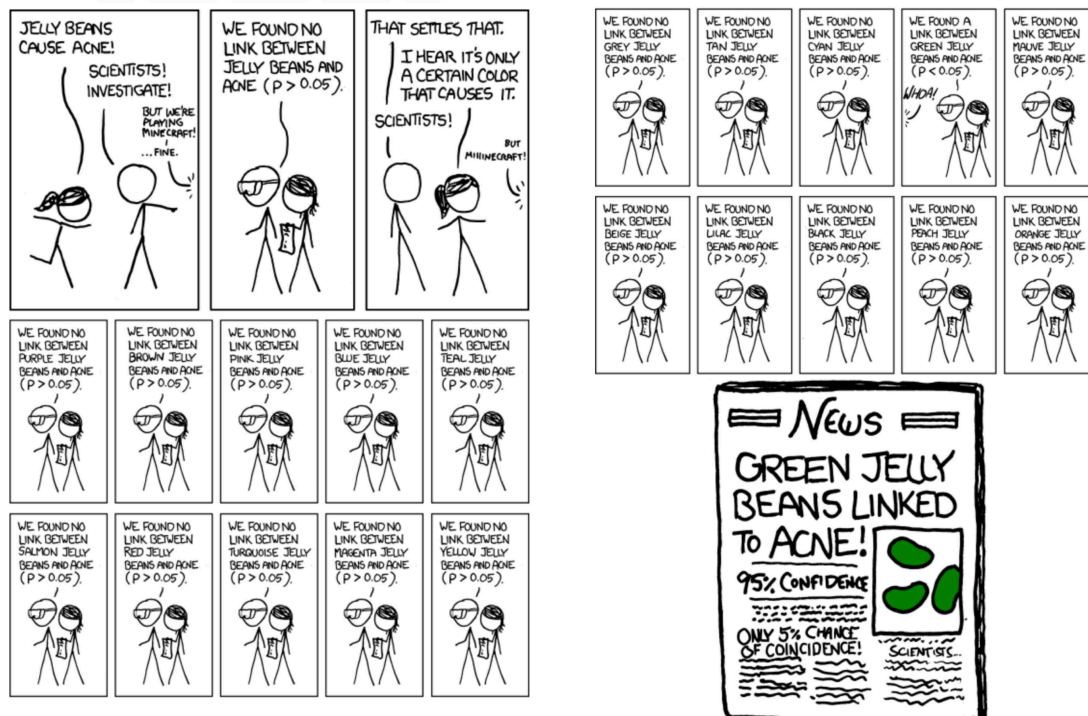
	Non-significant	Significant	
True nulls	TN	FP	m_0
False nulls	FN	TP	$m - m_0$
		R	m

- In the cartoon we have $m = 20$ hypotheses individually tested with $\alpha = 0.05$. We observe $R = 1$, but $E(\text{FP}) = m\alpha = 1$, so this is not a surprise.

stat.epfl.ch

Autumn 2022 – slide 76

The perils of multiple testing



stat.epfl.ch

Autumn 2022 – slide 77

Graphical approach

- Graphs can be helpful in suggesting which hypotheses are most suspect, and it is helpful to highlight the corresponding (i.e., smallest) P-values.
- $P \sim U(0, 1)$ implies $Z = -\log_{10} P \sim \exp(\lambda)$ with $\lambda = \ln 10$.
- With this transformation small P_j become large Z_j ; note that $Z_j > a$ iff $P_j < 10^{-a}$.
- If H_0 is true and the tests are independent, then $Z_1, \dots, Z_m \stackrel{\text{iid}}{\sim} \exp(\lambda)$ and the **Rényi representation**

$$Z_{(r)} \stackrel{D}{=} \lambda^{-1} \sum_{j=1}^r \frac{E_j}{m+1-j}, \quad r = 1, \dots, m, \quad E_1, \dots, E_m \stackrel{\text{iid}}{\sim} \exp(1),$$

applies to their order statistics. Then

- a plot of the ordered empirical Z_j against their expectations should be straight;
- outliers, very large Z_j (i.e., very small P_j), casting doubt on the corresponding H_j .
- For very small P_j (i.e., large Z_j) the uniformity may fail even under H_0 , because the null distributions give poor approximations in the extreme tail; then some form of model-fitting may be needed.
- Similar ideas apply to z statistics (e.g., in regression): use a normal QQ-plot (excluding the intercept etc.) as a basis for discussion of significant effects.

GWAS, I

- A **genome-wide association study (GWAS)** tests the association between SNPs ('single nucleotide polymorphisms') and a phenotype such as the expression of a protein. The null hypotheses are

$$H_{0,j} : \text{no association between the expression of the protein and SNP}_j, \quad j = 1, \dots, m.$$

- In a simple model we construct statistics Y_j such that $Y_j \sim \mathcal{N}(\theta_j, 1)$, where $\theta_j = 0$ under $H_{0,j}$, and we take $T_j = |Y_j|$, which is likely to be far from zero if $\theta_j \gg 0$ or $\theta_j \ll 0$.
- If $t_{\text{obs},j}$ denotes the observed value of T_j , then the P-value for association j is

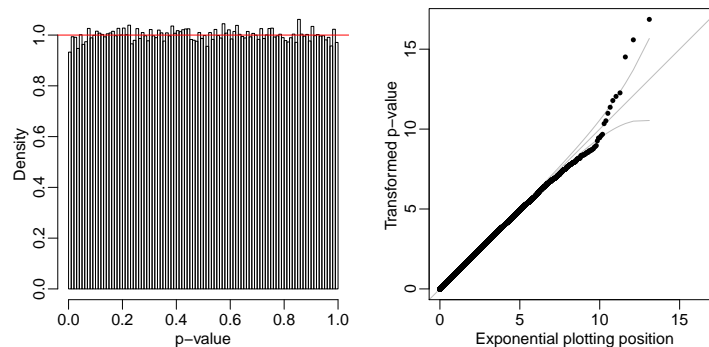
$$p_{\text{obs},j} = P_0(T_j > t_{\text{obs},j}) = 1 - P_0(-t_{\text{obs},j} \leq Y_j \leq t_{\text{obs},j}) \doteq 2\Phi(-t_{\text{obs},j}),$$

where the approximation comes from the fact that $Y_j \sim \mathcal{N}(0, 1)$ under $H_{0,j}$.

- Here it is reasonable to expect that the effects are **sparse**, i.e., most of the $\theta_j = 0$, and we seek a needle in a haystack.
- With many tests it is essential to ensure that the true positives are not drowned in the mass of false positives.

GWAS, II

- ☐ Left: a histogram of the P-values for tests of the association between $m = 275297$ SNPs and the expression of the protein CFAB.
- ☐ The P-values for SNPs not associated with CFAB are uniformly distributed. Is there an excess of small P-values?
- ☐ Right: exponential Q-Q plot of the $Z_j = -\log P_j$. What do you make of it?



Control

- ☐ With several tests Type I error generalises to the **familywise error rate (FWER)**, i.e., the probability of at least one false positive when the individual hypotheses are tested,

$$\text{FWER} = P(\text{FP} \geq 1) = 1 - P(\text{accept all } H_j, j \in \mathcal{I}),$$

and we aim to control this by ensuring that $\text{FWER} \leq \alpha$.

- ☐ There are different notions of control:
 - **weak control** guarantees $\text{FWER} \leq \alpha$ only under H_0 , i.e., $m_0 = m$;
 - **strong control** guarantees $\text{FWER} \leq \alpha$ for any configuration of null and alternative hypotheses.
- ☐ If all the tests are independent and we use individual levels α , then

$$\text{FWER} = 1 - P(\text{FP} = 0) = 1 - (1 - \alpha)^{m_0} \rightarrow 1, \quad m_0 \rightarrow \infty.$$

- ☐ If conversely we fix FWER and the tests are independent we need

$$\alpha = 1 - (1 - \text{FWER})^{1/m_0},$$

so with $m_0 = 20$ and $\text{FWER} = 0.05$ we need $\alpha \doteq 0.0026$, so the power for individual tests will be tiny (recall ROC curves).

Bonferroni methods

- If P_j is the P-value for the j th test and we reject H_j if $P_j < \alpha/m$, then **Boole's inequality** (the first of the **Bonferroni inequalities**) gives

$$\text{FWER} = P(\text{FP} \geq 1) = P\left(\bigcup_{j=1}^{m_0} \left\{P_j \leq \frac{\alpha}{m}\right\}\right) \leq \sum_{j=1}^{m_0} P\left(P_j \leq \frac{\alpha}{m}\right) = m_0 \frac{\alpha}{m} \leq \alpha,$$

so we have strong control of FWER, even if the tests are dependent.

- Note that we could replace α/m for test j by α_j such that $\sum_{j=1}^m \alpha_j \leq \alpha$.
- The resulting **Bonferroni procedure** lacks power when m is large (because α/m is very small), but its assumptions are very weak.
- An improvement is the **Holm–Bonferroni procedure**: for given α ,
 - order the P-values as $P_{(1)} \leq \dots \leq P_{(m)}$ and the hypotheses as $H_{(1)}, \dots, H_{(m)}$, then
 - reject $H_{(1)}, \dots, H_{(S-1)}$, where

$$S = \min \left\{ s : P_{(s)} > \frac{\alpha}{m+1-s} \right\}.$$

This gives strong control and is more powerful than the basic Bonferroni procedure.

stat.epfl.ch

Autumn 2022 – slide 82

Note: Bonferroni–Holm procedure

If $\text{FP} \geq 1$, then we must have wrongly rejected some H_j for which $j \in \mathcal{I}$. If $H_{(s)}$ is the first such hypothesis rejected in the sequential procedure, then the $s-1$ hypotheses rejected before it must have been among the $m-m_0$ false null hypotheses, so $s-1 \leq m-m_0$, i.e., $m_0 \leq m+1-s$. As $H_{(s)}$ was rejected, the corresponding P-value satisfies

$$P_{(s)} \leq \frac{\alpha}{m+1-s} \leq \frac{\alpha}{m_0}.$$

This implies that if $\text{FP} \geq 1$ then the P-value for at least one of the true null hypotheses satisfies $P_j \leq \alpha/m_0$, and so Boole's inequality gives

$$\text{FWER} = P(\text{FP} \geq 1) \leq P\left(\bigcup_{j \in \mathcal{I}} \{P_j \leq \alpha/m_0\}\right) \leq \alpha.$$

As the only assumption needed for the above argument was that the null P-values are $U(0,1)$, this procedure strongly controls the FWER at level α .

stat.epfl.ch

Autumn 2022 – note 1 of slide 82

False discovery rate

- When m is large and the goal is exploratory, Bonferroni procedures are unreasonably stringent, and it seems preferable to try and control the **false discovery proportion**

$$I(R > 0)FP/R,$$

where R is the number of rejected null hypotheses. The intention is to bound the proportion of false positives among the rejections.

- Control of $I(R > 0)FP/R$ is impossible because \mathcal{I} is unknown, so instead we try and control the **false discovery rate (FDR)**

$$FDR = E\{I(R > 0)FP/R\}.$$

- Strong control is achieved by the **Benjamini–Hochberg procedure**: specify α , then
 - order the P-values as $P_{(1)} \leq \dots \leq P_{(m)}$ and the hypotheses as $H_{(1)}, \dots, H_{(m)}$,
 - reject $H_{(1)}, \dots, H_{(R)}$, where

$$R = \max \left\{ r : P_{(r)} < \frac{\alpha r}{m} \right\}.$$

This guarantees that $FDR \leq \alpha$, but does not bound the actual proportion of false positives, just its expectation. Often $\alpha = 0.1, 0.2, \dots$

Note: Derivation of the Benjamini–Hochberg procedure

- Let the P-values for the false null hypotheses be P'_1, \dots, P'_{m_1} , say, independent of the true null P-values $P_1, \dots, P_{m_0} \stackrel{\text{iid}}{\sim} U(0, 1)$. Then

$$\{R = r\} \cap \{P_1 \leq r\alpha/m\} = \{P_1 \leq r\alpha/m\} \cap \{R_{-1} = r - 1\},$$

where $\{R_{-1} = r - 1\}$ is the event that there are exactly $r - 1$ rejections among H_2, \dots, H_m . Then the false discovery proportion is

$$\sum_{r=1}^m \frac{\text{FP}}{r} I(R = r) = \sum_{r=1}^m \frac{I(R = r)}{r} \sum_{j=1}^{m_0} I(P_j \leq r\alpha/m),$$

and by symmetry of the P_j this has the same expectation as

$$m_0 \sum_{r=1}^m \frac{I(R = r)}{r} I(P_1 \leq r\alpha/m) = m_0 \sum_{r=1}^m \frac{I(R_{-1} = r - 1)}{r} I(P_1 \leq r\alpha/m).$$

Thus the false discovery rate is

$$\begin{aligned} \text{FDR} &= m_0 \sum_{r=1}^m \frac{1}{r} P(R_{-1} = r - 1, P_1 \leq r\alpha/m) \\ &= m_0 \sum_{r=1}^m \frac{1}{r} P(R_{-1} = r - 1 \mid P_1 \leq r\alpha/m) P(P_1 \leq r\alpha/m) \\ &= m_0 \sum_{r=1}^m \frac{1}{r} P(R_{-1} = r - 1) \frac{r\alpha}{m} \\ &= \frac{m_0\alpha}{m} \sum_{r=0}^{m-1} P(R_{-1} = r) \\ &= \frac{m_0\alpha}{m} \leq \alpha. \end{aligned}$$

The main steps above successively use the definition of conditional probability, the facts that P_1 and R_{-1} are independent and $P_1 \sim U(0, 1)$, and the fact that $R_{-1} \in \{0, 1, \dots, m - 1\}$.

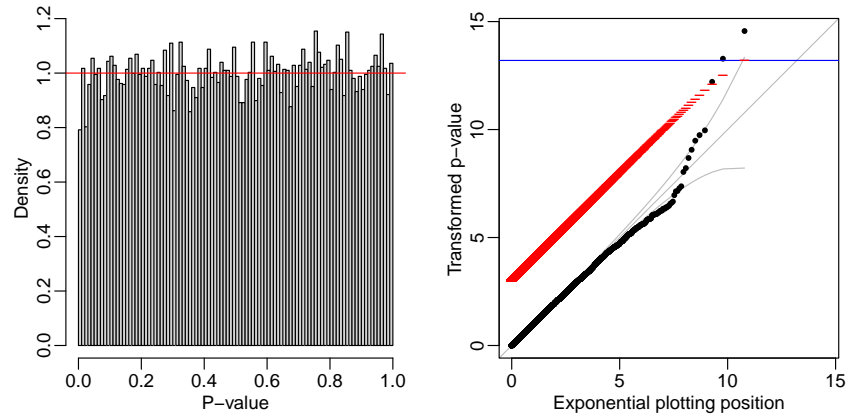
- We see that the Benjamini–Hochberg procedure strongly controls the FDR under the conditions above.
- Note that
- if $m_0 \ll m$, then the last inequality may be very unequal, so the FDR may in fact be much lower than α .
 - if the P-values are dependent in such a way that

$$P(R_{-1} = r - 1 \mid P_1 \leq r\alpha/m) \leq P(R_{-1} = r - 1),$$

then the result also holds.

GWAS, II

- Left: a histogram of $Q_j = 10P_j$ (when $P_j < 0.1$) for tests of the association between $m = 27530$ SNPs and the expression of the protein CFAB.
- Right: exponential Q-Q plot of $Z_j = -\log Q_j$, with Bonferroni cutoff (blue) and Benjamini–Hochberg cutoffs, both with $\alpha = 0.05$.



stat.epfl.ch

Autumn 2022 – slide 84

Comments

- The Bonferroni–Holm procedure compares $P_{(1)}, P_{(2)}, \dots$ to $\alpha/m, \alpha/(m-1), \dots$, whereas the ordinary Bonferroni procedure compares all the P_j to α/m .
- The **Simes procedure** (exercises) has exact FWER α for independent tests and then is preferable to the Bonferroni–Holm procedure.
- The Benjamini–Hochberg procedure strongly controls the false discovery rate, comparing the ordered P-values to $\alpha/m, 2\alpha/m, \dots, \alpha$.
- The first two also give strong control when the P-values are dependent. So does the third, using the comparison

$$P_{(j)} \leq \frac{j\alpha}{mc(m)},$$

with $c(m) = 1$ when the tests are independent or positively dependent, and $c(m) = \sum_{j=1}^m 1/j$ under arbitrary dependence.

- Many variants exist, but these versions are simple and widely used.
- Other classical procedures for multiple testing in regression settings are named after
 - Tukey — bounds the maximum of t statistics for different tests;
 - Scheffé — simultaneously bounds all possible linear combinations of estimates $\hat{\beta}$;
 - Dunnett — compares different treatments with the same control.

stat.epfl.ch

Autumn 2022 – slide 85

Overview

- In theoretical discussion we glibly write something like

$$\text{“Let } Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} f(y; \theta) \dots \text{”}$$

but in applications this cannot be taken for granted.

- Ideally we can ensure random sampling and full measurement of observations from a well-specified population, but if not, possible complications include:
 - selection of observations based on their values, especially truncation;
 - censoring;
 - dependence;
 - missing data.
- We now briefly discuss these ...

stat.epfl.ch

Autumn 2022 – slide 87

Selection

- If the available data were selected from a population using a mechanism expressible in probabilistic terms, then the likelihood is

$$P(Y = y \mid \mathcal{S}; \theta),$$

where \mathcal{S} is the selection event. If \mathcal{S} is unknown or not probabilistic, only sensitivity analysis is possible (at best).

- A common example is **truncation** of independent data, where $\mathcal{S}_j = \{Y_j \in \mathcal{I}_j\}$ for some set \mathcal{I}_j , giving likelihood

$$\prod_{j=1}^n f(y_j \mid y_j \in \mathcal{I}_j; \theta).$$

Example 34 In certain demographic databases on very old persons, an individual born on calendar date x is included only if they die aged $u_0 + t$, where u_0 is a high threshold (e.g., 100 years) and $t \geq 0$, between two calendar dates c_1 and c_2 . The likelihood contribution for this person is then of form

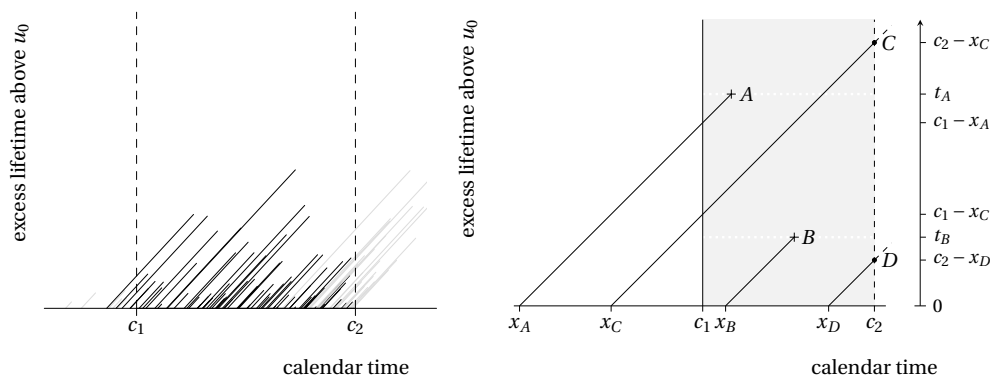
$$\frac{f(t)}{\mathcal{F}(a) - \mathcal{F}(b)}, \quad a < t < b, \quad [a, b] = [\max(0, c_1 - x), c_2 - x],$$

where x is the calendar date at which they reach age u_0 . See the next page.

stat.epfl.ch

Autumn 2022 – slide 88

Selection in a Lexis diagram



Lexis diagrams showing age on the vertical axis and calendar time on the horizontal axis. Only ages over u_0 are shown.

Left: only the individuals with solid lines appear in the sample.

Right: explanation of the intervals for which different individuals are observed.

Censoring

- ☐ Truncation determines which observations appear in a sample, whereas censoring reduces the information available in the sample.
- ☐ **Censoring** is very common in studies on lifetime data and leads to the precise values of certain observations being unknown:
 - **right-censoring** results in $(T = \min(Y, b), D = I(Y \leq b))$ for some b ;
 - **left-censoring** results in $(T = \max(Y, a), D = I(Y > a))$ for some a ;
 - **interval-censoring** results in $(Y, I(a < Y \leq b))$, $(a, I(Y \leq a))$ or $(b, I(Y > b))$, or it is known only which of the disjoint intervals $\mathcal{I}_1, \dots, \mathcal{I}_K$ contains Y .
- ☐ In each case we lose information when Y lies within some (possibly random) interval \mathcal{I} , often with the assumption that $Y \perp\!\!\!\perp \mathcal{I}$.
- ☐ **Rounding** is a form of interval censoring, and we have already seen (exercises) that little information is lost if the rounding is not too coarse.
- ☐ Likelihood contributions based on right- and left-censored observations are

$$f_Y(t)^d \{1 - F_Y(t)\}^{1-d}, \quad f_Y(t)^d \{F_Y(t)\}^{1-d}.$$

- ☐ Truncation and censoring can arise in the same study; see the Lexis diagram.

Dependent data

- If the joint density of $Y = (Y_1, \dots, Y_n)$ is known, we can write

$$f(y; \theta) = f(y_1, \dots, y_n; \theta) = f(y_1; \theta) \prod_{j=2}^n f(y_j | y_1, \dots, y_{j-1}; \theta),$$

the so-called **prediction decomposition**.

- This is most useful if the data arise in time order and satisfy the **Markov property**, that given the 'present' Y_{j-1} , the 'future', Y_j, Y_{j+1}, \dots , is independent of the 'past', \dots, Y_{j-3}, Y_{j-2} , so

$$f(y_j | y_1, \dots, y_{j-1}; \theta) = f(y_j | y_{j-1}; \theta)$$

and the product above simplifies to

$$f(y; \theta) = f(y_1; \theta) \prod_{j=2}^n f(y_j | y_{j-1}; \theta).$$

- Many variants of this are possible.

Example 35 (Poisson birth process) Find the likelihood when $Y_0 \sim \text{Pois}(\theta)$ and Y_0, \dots, Y_n are such that $Y_{j+1} | Y_0 = y_0, \dots, Y_j = y_j \sim \text{Pois}(\theta y_j)$.

stat.epfl.ch

Autumn 2022 – slide 91

Note to Example 35

Here

$$f(y_{j+1} | y_j; \theta) = \frac{(\theta y_j)^{y_{j+1}}}{y_{j+1}!} \exp(-\theta y_j), \quad y_{j+1} = 0, 1, \dots, \quad \theta > 0.$$

If Y_0 is Poisson with mean θ , the joint density of data y_0, \dots, y_n is

$$f(y_0; \theta) \prod_{j=1}^n f(y_j | y_{j-1}; \theta) = \frac{\theta^{y_0}}{y_0!} \exp(-\theta) \prod_{j=0}^{n-1} \frac{(\theta y_j)^{y_{j+1}}}{y_{j+1}!} \exp(-\theta y_j),$$

so the likelihood is

$$L(\theta) = \left(\prod_{j=0}^n y_j! \right)^{-1} \exp(s_0 \log \theta - s_1 \theta), \quad \theta > 0,$$

where $s_0 = \sum_{j=0}^n y_j$ and $s_1 = 1 + \sum_{j=0}^{n-1} y_j$. This is a (2,1) exponential family.

stat.epfl.ch

Autumn 2022 – note 1 of slide 91

Missing data

- ☐ Missing data are widespread in applications, especially those involving living subjects.
- ☐ Central problems are:
 - uncertainty increases due to missingness;
 - assumptions about missingness cannot be checked directly, so inferences are fragile.
- ☐ Suppose ideal is inference on θ based on n independent pairs (X, Y) , but some Y are missing, indicated by a variable I , so we observe either $(x, y, 1)$ or $(x, ?, 0)$.
- ☐ The likelihood contributions from individuals with complete data and with y missing are respectively

$$P(I = 1 \mid x, y)f(y \mid x; \theta)f(x; \theta), \quad \int P(I = 0 \mid x, y)f(y \mid x; \theta)f(x; \theta) dy,$$

and there are three possibilities:

- data are **missing completely at random**, $P(I = 0 \mid x, y) = P(I = 0)$;
- data are **missing at random**, $P(I = 0 \mid x, y) = P(I = 0 \mid x)$; and
- **non-ignorable non-response**, $P(I = 0 \mid x, y)$ depends on y and maybe on x .

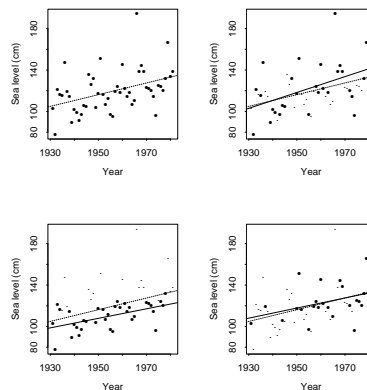
The first two are sometimes called **ignorable non-response**, as then I has no information about θ and can (mostly) be ignored.

stat.epfl.ch

Autumn 2022 – slide 92

Example

Missing data in straight-line regression for Venice sea-level data. Clockwise from top left: original data, data with values missing completely at random, data with values missing at random — missingness depends on x but not on y , and data with non-ignorable non-response — missingness depends on both x and y . Missing values are represented by a small dot. The dotted line is the fit from the full data, the solid lines those from the non-missing data.



stat.epfl.ch

Autumn 2022 – slide 93

Example

	Truth	Average estimate (average standard error)			
		Full	MCAR	MAR	NIN
β_0	120	120 (2.79)	120 (4.02)	120 (4.73)	132 (3.67)
β_1	0.50	0.49 (0.19)	0.48 (0.28)	0.50 (0.32)	0.20 (0.25)

- Average estimates and standard errors for missing value simulation based on Venice data, for full dataset, with data missing completely at random (MCAR), missing at random (MAR) and with non-ignorable non-response (NIN) and non-response mechanisms

$$P(I = 0 \mid x, y) = \begin{cases} 0.5, \\ \Phi \{0.05(x - \bar{x})\}, \\ \Phi [0.05(x - \bar{x}) + \{y - \beta_0 - \beta_1(x - \bar{x})\} / \sigma]; \end{cases}$$

In each case roughly one-half of the observations are missing.

- Data loss increases the variability of the estimates but their means are unaffected when the non-response is ignorable; otherwise they become entirely unreliable.
- Standard errors for the averages for $\hat{\beta}_0$ and $\hat{\beta}_1$ are at most 0.16 and 0.01; those for their standard errors are at most 0.03 and 0.002.

stat.epfl.ch

Autumn 2022 – slide 94

Discussion

- Truncation, censoring and other forms of **data coarsening** are widely observed in time-to-event data and there is a huge literature on dealing with them, especially in terms of non- and semi-parametric estimation.
- Selection (especially self-selection!) can totally undermine analyses if ignored or if it can't be modelled appropriately.
- The Markov property plays a key simplifying role in inference based on time series, and generalisations are important in spatial and other types of complex data.
- Missingness is usually the most annoying of the complications above:
 - it is quite common in applications, often for ill-specified reasons;
 - when there is NIN and a non-negligible proportion of the data is missing, correct inference requires us to specify the missingness mechanism correctly;
 - in practice it is hard to tell whether missingness is ignorable, so fully reliable inference is largely out of reach;
 - sensitivity analysis and or bounds to assess how heavily the conclusions depend on plausible mechanisms for non-response is then useful.

stat.epfl.ch

Autumn 2022 – slide 95

Motivation

□ Likelihood

- provides a general paradigm for inference on parametric models, with many generalisations and variants;
- is a central concept in both frequentist and Bayesian statistics;
- has a simple, general and widely-applicable 'large-sample' theory; but
- is not a panacea!

□ Plan below:

- recall some basics on convergence;
- give (fairly) general setup for parameter;
- prove main results for scalar parameter;
- discussion of inference;
- vector parameter, nuisance parameters, ...

stat.epfl.ch

Autumn 2022 – slide 97

Reminders I

Definition 36 Let X, X_1, X_2, \dots be random variables with cumulative distribution functions F, F_1, F_2, \dots . Then

- (a) X_n converges to X **in probability**, $X_n \xrightarrow{P} X$, if $\lim_{n \rightarrow \infty} P(|X_n - X| > \varepsilon) = 0$ for all $\varepsilon > 0$;
- (b) X_n converges to X **in distribution**, $X_n \xrightarrow{D} X$, if $\lim_{n \rightarrow \infty} F_n(x) = F(x)$ at each point x where $F(x)$ is continuous.
- (c) A sequence X_1, X_2, \dots of estimators of a parameter θ is **(weakly) consistent** if $X_n \xrightarrow{P} \theta$.

Theorem 37 Let x_0, y_0 be constants, $X, Y, \{X_n\}, \{Y_n\}$ random variables and $g(\cdot)$ and $h(\cdot, \cdot)$ continuous functions. Then

$$\begin{aligned} X_n \xrightarrow{P} X &\Rightarrow X_n \xrightarrow{D} X, \\ X_n \xrightarrow{D} x_0 &\Rightarrow X_n \xrightarrow{P} x_0, \\ X_n \xrightarrow{P} X &\Rightarrow g(X_n) \xrightarrow{P} g(X), \\ X_n \xrightarrow{D} X \text{ and } Y_n \xrightarrow{D} y_0 &\Rightarrow h(X_n, Y_n) \xrightarrow{D} h(X, y_0). \end{aligned}$$

The last two lines are known as the **continuous mapping theorem** and **Slutsky's theorem**.

stat.epfl.ch

Autumn 2022 – slide 98

Reminders II

Theorem 38 (Weak law of large numbers) If X, X_1, X_2, \dots are independent identically distributed random variables and $E(X)$ is finite, then $\bar{X} = n^{-1}(X_1 + \dots + X_n) \xrightarrow{P} E(X)$.

Theorem 39 (Central limit theorem, CLT) If $X_1, X_2, \dots \stackrel{\text{iid}}{\sim} (\mu, \sigma^2)$ and $0 < \sigma^2 < \infty$, then

$$Z_n = \frac{n^{1/2}(\bar{X} - \mu)}{\sigma} \xrightarrow{D} Z \sim \mathcal{N}(0, 1), \quad n \rightarrow \infty.$$

Theorem 40 ('Delta method') If $a_n(X_n - \mu) \xrightarrow{D} Z$, where $a_n, \mu \in \mathbb{R}$ for all n , $a_n \rightarrow \infty$ as $n \rightarrow \infty$, and g is continuously differentiable at μ , then $a_n\{g(X_n) - g(\mu)\} \xrightarrow{D} g'(\mu)Z$.

- Many more general laws of large numbers and versions of the CLT exist.
- The delta method also applies with $X_n, Z \in \mathbb{R}^p$, $g(x) : \mathbb{R}^p \rightarrow \mathbb{R}^q$ continuously differentiable and $g'(\mu)$ replaced by $J_g(\mu) = \partial g(\mu)/\partial \mu^T$.

Theorem 41 If X is a random variable, $a > 0$ a constant, h a non-negative function and g a convex function, then

$$\begin{aligned} P\{h(X) \geq a\} &\leq E\{h(X)\}/a, & (\text{basic inequality}), \\ g\{E(X)\} &\leq E\{g(X)\}, & (\text{Jensen's inequality}). \end{aligned}$$

Basic setup

- Let $Y, Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} g$, and define the **Kullback–Leibler divergence** from g to a density f ,

$$\text{KL}(g, f) = E_g\{\log g(Y) - \log f(Y)\} = E_g\left[-\log \left\{\frac{f(Y)}{g(Y)}\right\}\right] \geq 0,$$

where the inequality follows because $-\log x$ is convex and is strict unless $f \equiv g$.

- g is the unknown model from which the Y_j are drawn, and f is a candidate model.
- In a parametric setting there is a family of models, $f \in \mathcal{F} = \{f_\theta : \theta \in \Theta\}$, so minimising $\text{KL}(g, f)$ over f is equivalent to maximising $E_g \log f(Y; \theta)$, which is estimated by

$$\bar{\ell}(\theta) = n^{-1} \sum_{j=1}^n \log f(Y_j; \theta) \xrightarrow{P} E_g \log f(Y; \theta), \quad n \rightarrow \infty.$$

- $\theta_g = \arg \max_{\theta} E_g \log f(Y; \theta)$ gives the best fit of f_θ to g .
- In an ideal case $g = f_{\theta_g}$, i.e., $g \in \mathcal{F}$, but the theory does not require this (yet).
- $\hat{\theta} = \arg \max_{\theta} \bar{\ell}(\theta)$ is the natural estimator of θ_g , and we hope that $\hat{\theta} \xrightarrow{P} \theta_g$ as $n \rightarrow \infty$.
- Unfortunately this requires conditions to restrict the variation of $\bar{\ell}$ as it converges.

Regular models

- Recall the notation $\nabla g(\theta) = \partial g(\theta)/\partial \theta$ and $\nabla^2 g(\theta) = \nabla \nabla^T g(\theta) = \partial^2 g(\theta)/\partial \theta \partial \theta^T$.
- The following conditions ensure consistency and asymptotic normality of the MLE:
 - (C1) θ_g is interior to $\Theta \subset \mathbb{R}^d$ for some finite d , and Θ is compact;
 - (C2) the densities f_θ defined by any two different values of $\theta \in \Theta$ are distinct;
 - (C3) there is a neighbourhood \mathcal{N} of θ_g within which the first three derivatives of the log likelihood with respect to θ exist almost surely, and for $r, s, t = 1, \dots, d$ satisfy $|\partial^3 \log f(Y; \theta)/\partial \theta_r \partial \theta_s \partial \theta_t| < m(Y)$ with $E_g\{m(Y)\} < \infty$; and
 - (C4) within \mathcal{N} , the $d \times d$ matrices

$$I_1(\theta) = E_g \{-\nabla^2 \log f(Y; \theta)\}, \quad K_1(\theta) = E_g \{\nabla \log f(Y; \theta) \nabla^T \log f(Y; \theta)\},$$

are finite and positive definite. When $g = f_{\theta_g}$ we shall see that $K_1(\theta_g) = I_1(\theta_g)$.

- Comments:
 - (C1) ensures that $\hat{\theta}$ can be 'on all sides' of θ_g in the limit;
 - (C2) is essential for consistency, otherwise $\hat{\theta}$ might not converge;
 - (C3) is a technical condition needed to bound terms of a Taylor series; and
 - (C4) ensures that the asymptotic variance of $\hat{\theta}$ is positive definite.

Consistency of the MLE

Lemma 42 *If θ is scalar, then a sequence of maximum likelihood estimators $\hat{\theta}$ exists such that $\hat{\theta} \xrightarrow{P} \theta_g$.*

This result:

- does not require f_θ to be smooth, so it is quite general;
- guarantees that a consistent sequence exists, but not that we can find it;
- but if the log likelihood is convex (as in exponential families, for example), then there is (at most) one maximum for any n , and if it exists this must converge to θ_g ;
- can be generalized to vector θ , but the argument is more delicate;
- van der Vaart (1998, *Asymptotic Statistics*, Chapter 5) gives more general proofs.

Note to Lemma 42

- As the θ s correspond to different densities, there is precisely one θ_g that minimises $\text{KL}(g, f_\theta)$.
- Take any $\varepsilon > 0$ and let $\theta_+, \theta_- = \theta_g \pm \varepsilon$, write $D_n(\theta) = \bar{\ell}(\theta_g) - \bar{\ell}(\theta)$, so $D_n(\theta_g) = 0$, and note that as $n \rightarrow \infty$,

$$D_n(\theta_+) \xrightarrow{P} \text{KL}(g, f_{\theta_+}) - \text{KL}(g, f_{\theta_g}) = a_+ > 0, \quad D_n(\theta_-) \xrightarrow{P} \text{KL}(g, f_{\theta_-}) - \text{KL}(g, f_{\theta_g}) = a_- > 0.$$

- If A_n and B_n denote the events $D_n(\theta_+) > 0$ and $D_n(\theta_-) > 0$, Boole's inequality gives

$$P(A_n \cap B_n) = 1 - P(A_n^c \cup B_n^c) \geq 1 - P(A_n^c) - P(B_n^c).$$

Now

$$P(A_n^c) = P\{D_n(\theta_+) \leq 0\} = P\{a_+ - D_n(\theta_+) \geq a_+\} \leq P\{|D_n(\theta_+) - a_+| \geq a_+\} \rightarrow 0, \quad n \rightarrow \infty,$$

and likewise $P(B_n^c) \rightarrow 0$. Hence $P(A_n \cap B_n) \rightarrow 1$.

- Hence there is a local minimum of $D_n(\theta)$, or equivalently a local maximum of $\bar{\ell}(\theta)$, inside the interval $(\theta_g - \varepsilon, \theta_g + \varepsilon)$ with probability one as $n \rightarrow \infty$, and as this is true for arbitrary ε , the corresponding sequence of maximisers $\hat{\theta}$ satisfies $P(|\hat{\theta} - \theta_g| > \varepsilon) \rightarrow 0$ and therefore is consistent.

Asymptotic normality of the MLE

Theorem 43 *If θ is scalar and the regularity conditions hold, then the sequence of consistent maximum likelihood estimators $\hat{\theta}$ satisfies*

$$n^{1/2}(\hat{\theta} - \theta_g) \xrightarrow{D} \mathcal{N}_d\{0, I_1^{-1}(\theta_g)K_1(\theta_g)I_1^{-1}(\theta_g)\},$$

where $d = 1$ and for a single observation Y we define

$$I_1(\theta) = \mathbb{E}_g \{-\nabla^2 \log f(Y; \theta)\}, \quad K_1(\theta) = \mathbb{E}_g \{\nabla \log f(Y; \theta) \nabla^T \log f(Y; \theta)\}.$$

- In the vector case, $d > 1$, the above **sandwich variance matrix** expression also applies.
- This implies that for large n ,

$$\hat{\theta} \dot{\sim} \mathcal{N}_d\{\theta_g, I^{-1}(\theta_g)K(\theta_g)I^{-1}(\theta_g)\},$$

where $I(\theta) = nI_1(\theta)$, $K(\theta) = nK_1(\theta)$ correspond to a sample of size n .

- This provides tests and confidence intervals based on the approximate pivots

$$v_{rr}^{-1/2}(\hat{\theta}_r - \theta_{g,r}) \dot{\sim} \mathcal{N}(0, 1), \quad r = 1, \dots, d,$$

where v_{rr} are the diagonal elements of an estimate of $I^{-1}(\theta_g)K(\theta_g)I^{-1}(\theta_g)$.

- When $g = f_{\theta_g}$, $I(\theta_g) = K(\theta_g)$ and the variance (matrix) becomes $I(\theta_g)^{-1}$.

Note to Theorem 43

- We first note that under the given conditions, θ_g gives a stationary point of $\text{KL}(g, f_\theta)$, and therefore

$$0 = \nabla \text{KL}(g, f_\theta)|_{\theta=\theta_g} = - \nabla \int \log f(y; \theta) g(y) dy \Big|_{\theta=\theta_g} = - \int \nabla \log f(y; \theta) \Big|_{\theta=\theta_g} g(y) dy,$$

so $E_g\{\nabla \log f(Y; \theta)\} = 0$.

- As $\hat{\theta}$ gives a local maximum of the differentiable function $\bar{\ell}(\theta) = n^{-1} \sum_{j=1}^n \log f(Y_j; \theta)$,

$$0 = \nabla \bar{\ell}(\hat{\theta}) = n^{-1} \sum_{j=1}^n \nabla \log f(Y_j; \hat{\theta}),$$

and (supposing now that θ is scalar, to simplify the expressions), Taylor series expansion gives

$$0 = \nabla \bar{\ell}(\theta_g) + (\hat{\theta} - \theta_g) \nabla^2 \bar{\ell}(\theta_g) + \frac{1}{2} (\hat{\theta} - \theta_g)^2 \nabla^3 \bar{\ell}(\theta^*),$$

where θ^* lies between θ_g and $\hat{\theta}$ (so $\theta^* \xrightarrow{P} \theta_g$), and hence we can write

$$n^{1/2}(\hat{\theta} - \theta_g) = \frac{n^{1/2} \nabla \bar{\ell}(\theta_g)}{-\nabla^2 \bar{\ell}(\theta_g) - R_n/2}, \quad R_n = (\hat{\theta} - \theta_g) \nabla^3 \bar{\ell}(\theta^*). \quad (4)$$

- Now

$$n^{1/2} \nabla \bar{\ell}(\theta_g) = n^{-1/2} \sum_{j=1}^n \nabla \log f(Y_j; \theta_g)$$

has mean (vector) zero and variance (matrix)

$$\text{var} \left\{ n^{-1/2} \sum_{j=1}^n \nabla \log f(Y_j; \theta_g) \right\} = n^{-1} \sum_{j=1}^n E_g \{ \nabla \log f(Y_j; \theta_g) \nabla^T \log f(Y_j; \theta_g) \} = K_1(\theta_g).$$

so the numerator of (4) converges in distribution to $\mathcal{N}\{0, K_1(\theta_g)\}$, using the CLT.

- Moreover the weak law of large numbers gives

$$-\nabla^2 \bar{\ell}(\theta_g) = -\frac{1}{n} \sum_{j=1}^n \nabla^2 \log f(Y_j; \theta_g) \xrightarrow{P} I_1(\theta_g).$$

- Lemma 44 shows that $R_n \xrightarrow{P} 0$, so the denominator of (4) tends in probability to $I_1(\theta_g)$.

- Putting the pieces together, we find that

$$n^{1/2}(\hat{\theta} - \theta_g) \xrightarrow{D} \mathcal{N}_d\{0, I_1(\theta_g)^{-1} K_1(\theta_g) I_1(\theta_g)^{-1}\}, \quad n \rightarrow \infty,$$

where the variance formula is also valid when I_1 and K_1 are $d \times d$ matrices.

- The information quantities based on a random sample of size n are $I(\theta_g) = n I_1(\theta_g)$ and $K(\theta_g) = n K_1(\theta_g)$, giving

$$\hat{\theta} \sim \mathcal{N}_d(\theta_g, I(\theta_g)^{-1} K(\theta_g) I(\theta_g)^{-1}),$$

in which the variance is of the usual order $1/n$.

Note: Lemma 44

Lemma 44 Under the conditions of Theorem 43, $R_n = (\hat{\theta} - \theta_g) \nabla^3 \bar{\ell}(\theta^*) \xrightarrow{P} 0$ as $n \rightarrow \infty$.

- For $\varepsilon > 0$, $B_n = \{|R_n| > \varepsilon\}$, $A_n = \{|\hat{\theta} - \theta_g| > \delta\}$ and $\delta > 0$ small enough that \mathcal{N} contains a ball of radius δ around θ_g , we have

$$P(|R_n| > \varepsilon) = P(B_n \cap A_n) + P(B_n \cap A_n^c) \leq P(A_n) + P(B_n \cap A_n^c),$$

where the first term tends to zero because the sequence $\hat{\theta}$ is consistent.

- If $|\hat{\theta} - \theta_g| < \delta$, then (C3) implies that

$$|R_n| \leq \delta n^{-1} \sum_{j=1}^n |\partial^3 \log f(Y_j; \theta^*) / \partial \theta^3| \leq \delta n^{-1} \sum_{j=1}^n m(Y_j) = \delta \bar{M}_n,$$

say, and clearly $\bar{M}_n \xrightarrow{P} M$, say. Therefore

$$P(B_n \cap A_n^c) = P(B_n \cap |\hat{\theta} - \theta_g| > \delta) \leq P(B_n \cap |R_n| \leq \delta \bar{M}_n)$$

and for $\eta > 0$ this equals

$$P(B_n \cap |R_n| \leq \delta \bar{M}_n \cap \bar{M}_n \leq M + \eta) + P(B_n \cap |R_n| \leq \delta \bar{M}_n \cap \bar{M}_n > M + \eta),$$

which is bounded by

$$P\{|R_n| > \varepsilon \cap |R_n| \leq \delta(M + \eta)\} + P(|\bar{M}_n - M| > \eta).$$

The last term here tends to zero, because $\bar{M}_n \xrightarrow{P} M$, and the first can be made equal to zero by choosing δ such that $\delta(M + \eta) < \varepsilon$. This proves the lemma.

Classical asymptotics

- The true model is supposed to lie in the candidate family, i.e., $g \in \mathcal{F}$, so $\theta_g \in \Theta$.
- We can differentiate under the integral sign and get the **Bartlett identities**:

$$1 = \int f(y; \theta) dy,$$

$$0 = \int \nabla \log f(y; \theta) \times f(y; \theta) dy,$$

$$0 = \int \nabla^2 \log f(y; \theta) \times f(y; \theta) dy + \int \nabla \log f(y; \theta) \nabla^T \log f(y; \theta) \times f(y; \theta) dy,$$

$$0 = \dots$$

giving the moments of the $d \times 1$ **score vector** $U(\theta) = \nabla \ell(\theta)$, viz

$$E\{U(\theta)\} = 0, \quad \text{var}\{U(\theta)\} = E\{\nabla \ell(\theta) \nabla^T \ell(\theta)\} = E\{-\nabla^2 \ell(\theta)\}, \quad \dots$$

- Hence $I(\theta) = K(\theta)$, and $I(\theta) = nI_1(\theta) = nK_1(\theta)$ when $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} g$.
- The assumption that $g \in \mathcal{F}$ is technically always false, but this is irrelevant if model-checking suggests that \mathcal{F} is 'close enough' to g .
- Crucially, the interest parameter ψ should have a stable interpretation for candidates likely to be close to g (i.e., within $n^{-1/2}$), so \mathcal{F} is 'robustly specified'.

Note to the Bartlett identities

- The first is true for any θ , and provided we can exchange the order of integration and differentiation we have
- $$0 = \nabla \int f(y; \theta) dy = \int \nabla f(y; \theta) dy = \int \nabla f(y; \theta) \frac{f(y; \theta)}{f(y; \theta)} dy = \int \nabla \log f(y; \theta) f(y; \theta) dy.$$
- The second stems from a second differentiation and applying the chain rule to the terms in the final integral here; likewise for the third and higher-order ones, which give higher-order moments of $U(\theta)$.
 - For independent data Y_1, \dots, Y_n we have $U(\theta) = \sum_{j=1}^n U_j(\theta)$, where the $U_j = \nabla \log f(Y_j; \theta)$ are independent, so using the Bartlett identities for the individual densities $f_j(y_j; \theta)$ we have

$$\text{var}\{U(\theta)\} = \sum_{j=1}^n \text{var}\{U_j(\theta)\} = \sum_{j=1}^n E\{U_j(\theta) U_j^T(\theta)\} = \sum_{j=1}^n -E\{\nabla^T U_j(\theta)\} = -E\{\nabla^T U(\theta)\}$$

and this equals $E\{-\nabla^2 \ell(\theta)\} = I(\theta)$, and this in turn equals $nI_1(\theta)$ if $Y_j \stackrel{\text{iid}}{\sim} f_{\theta_g}$.

In practice . . .

- We usually assume classical asymptotics and replace the sandwich matrix

$$I(\theta_g)^{-1}K(\theta_g)I(\theta_g)^{-1} \quad \text{by the \textbf{observed information matrix} } \hat{J} = -\nabla^2 \ell(\hat{\theta}),$$

which

- can be computed numerically without (possibly awkward) expectations,
- will (helpfully!) misbehave if the maximisation is questionable,
- has been found to give generally good results in applications,
- has the heuristic justification that $(\hat{\theta}, \hat{J})$ are approximately sufficient for θ_g , as

$$\ell(\theta_g) \doteq \ell(\hat{\theta}) - \frac{1}{2}(\hat{\theta} - \theta_g)^T \hat{J}(\hat{\theta} - \theta_g).$$

- Standard errors for $\hat{\theta}$ are the square roots of the diagonal elements of \hat{J}^{-1} .
- To make the sandwich we can replace $I(\theta_g)$ by \hat{J} and $K(\theta_g)$ by (some version of)

$$\hat{K} = \sum_{j=1}^n \nabla \log f(Y_j; \hat{\theta}) \nabla^T \log f(Y_j; \hat{\theta}),$$

though $\hat{J}^{-1} \hat{K} \hat{J}^{-1}$ can be unstable because of numerical problems with \hat{K} .

Related statistics

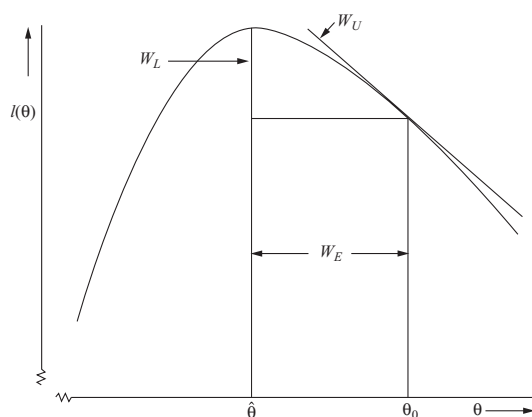


Figure 6.2. Three asymptotically equivalent ways, all based on the log likelihood function of testing null hypothesis $\theta = \theta_0$: W_E , horizontal distance; W_L vertical distance; W_U slope at null point.

From Cox (2006, *Principles of Statistical Inference*)

Related statistics

- When θ is scalar the asymptotic arguments support inference based on any of the pivots

$$T = t(\theta_g) = \hat{j}^{1/2}(\hat{\theta} - \theta_g) \sim \mathcal{N}(0, 1), \quad \text{Wald statistic,}$$

$$S = s(\theta_g) = \hat{j}^{-1/2}U(\theta_g) \sim \mathcal{N}(0, 1), \quad \text{score statistic,}$$

$$W = w(\theta_g) = 2\{\ell(\hat{\theta}) - \ell(\theta_g)\} \sim \chi_1^2, \quad \text{likelihood ratio statistic,}$$

$$R = r(\theta_g) = \text{sign}(\hat{\theta} - \theta_g)w(\theta_g)^{1/2} \sim \mathcal{N}(0, 1), \quad \text{likelihood root.}$$

The likelihood root has other names (e.g., directed likelihood ratio statistic).

- The distribution of W follows from the expansion on the previous slide.
 □ If $\hat{\theta}^o$ and $j(\hat{\theta}^o)$ have been obtained for observed data y^o , then the approximation

$$P_g\{T(\theta_g) \leq t^o(\theta_g)\} \doteq \Phi\{t^o(\theta_g)\}$$

leads to $(1 - \alpha)$ confidence interval $\hat{\theta}^o \pm j(\hat{\theta}^o)^{-1/2}z_{1-\alpha/2}$ based on T , while that based on W is

$$\{\theta : W^o(\theta) \leq \chi_1^2(1 - \alpha)\} = \{\theta : \ell^o(\theta) \geq \ell^o(\hat{\theta}^o) - \frac{1}{2}\chi_1^2(1 - \alpha)\},$$

where z_p and $\chi_\nu^2(p)$ are respectively the p quantiles of the $N(0, 1)$ and χ_ν^2 distributions.

Comments

- Comparative comments:
- confidence intervals based on T are symmetric, but those based on W or R take the shape of ℓ into account and are parametrisation-invariant;
 - in small samples the distributional approximations for W and R are better than that for T , and that for W can be improved by **Bartlett correction**, using $W_B = W/(1 + b/n)$;
 - confidence sets based on W may not be connected (and if so those based on T or R are unreliable);
 - the main use of S is for testing in situations where maximisation of ℓ is awkward, and then \hat{j} is often replaced by $I(\theta_g)$;
 - a variant of R , the **modified likelihood root**

$$R^* = r^*(\theta_g) = r(\theta_g) + \frac{1}{r(\theta_g)} \log \frac{q(\theta_g)}{r(\theta_g)},$$

often gives almost perfect inferences even in small samples (more later ...).

Example 45 Compute the above statistics when $y_1, \dots, y_n \stackrel{\text{iid}}{\sim} \exp(\theta)$ and compare the resulting inferences with those from an exact pivot.

Note to Example 45

□ The log likelihood is $\ell(\theta) = n(\log \theta - \theta \bar{y})$, for $\theta > 0$, which is clearly unimodal with $\hat{\theta} = 1/\bar{y}$ and $j(\theta) = n/\theta^2$.

□ Hence

$$t(\theta) = n^{1/2}(1 - \theta \bar{y}),$$

$$s(\theta) = n^{1/2}\{1/(\theta \bar{y}) - 1\},$$

$$w(\theta) = 2n \{\theta \bar{y} - \log(\theta \bar{y}) - 1\},$$

$$r(\theta) = \text{sign}(1 - \theta \bar{y}) [2n \{\theta \bar{y} - \log(\theta \bar{y}) - 1\}]^{1/2}.$$

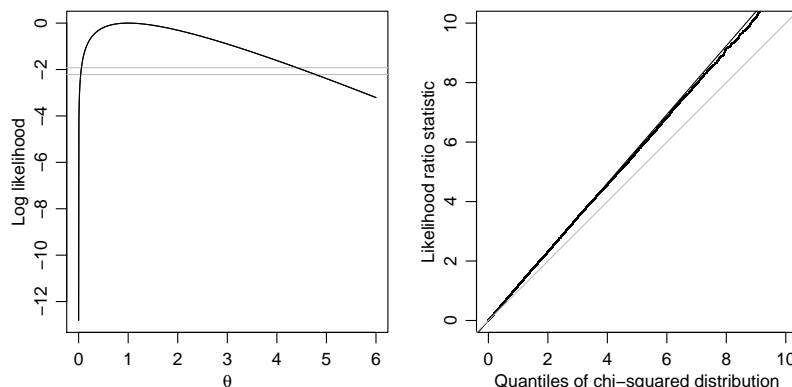
□ The exact pivot is $\theta \sum Y_j$ whose distribution is gamma with unit scale and shape parameter n .

□ Consider an exponential sample with $n = 1$ and $\bar{y} = 1$; then $\hat{\theta} = 1$. The log likelihood $\ell(\theta)$, shown in the left-hand panel of the figure, is unimodal but strikingly asymmetric, suggesting that confidence intervals based on an approximating normal distribution for $\hat{\theta}$ will be poor. The right-hand panel is a chi-squared probability plot in which the ordered values of simulated $w(\theta)$ are graphed against quantiles of the χ_1^2 distribution—if the simulations lay along the diagonal line $x = y$, then this distribution would be a perfect fit. The simulations do follow a straight line rather closely, but with slope $(1 + b/n)\chi_1^2$, where $b = 0.1544$. This indicates that the distribution of the Bartlett-adjusted likelihood ratio statistic $w(\theta)/(1 + b/n)$ would be essentially χ_1^2 . The 95% confidence intervals for θ based on the unadjusted and adjusted likelihood ratio statistics are (0.058, 4.403) and (0.042, 4.782) respectively.

stat.epfl.ch

Autumn 2022 – note 1 of slide 108

Exponential example

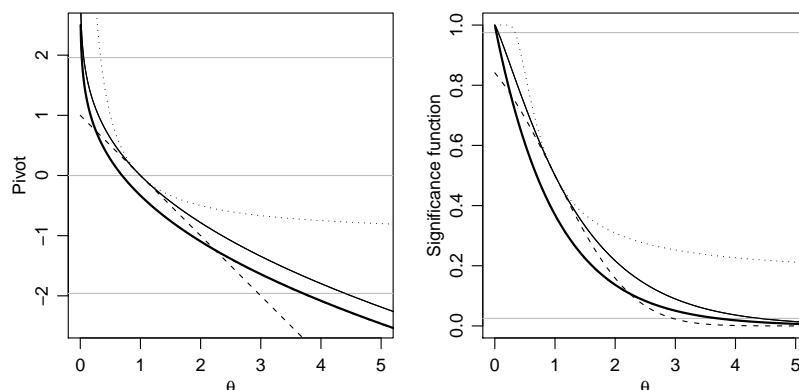


Likelihood inference for exponential sample of size $n = 1$. Left: log likelihood $\ell(\theta)$. Intersection of the function with the two horizontal lines gives two 95% confidence intervals for θ : the upper line is based on the χ_1^2 approximation to the distribution of $w(\theta)$, and the lower line is based on the Bartlett-corrected statistic. Right: comparison of simulated values of likelihood ratio statistic $w(\theta)$ with χ_1^2 quantiles. The χ_1^2 approximation is shown by the line of unit slope, while the $(1 + b/n)\chi_1^2$ approximation is shown by the upper straight line.

stat.epfl.ch

Autumn 2022 – slide 109

Exponential example



Approximate pivots and P-values based on an exponential sample of size $n = 1$. Left: likelihood root $r(\theta)$ (solid), score pivot $s(\theta)$ (dots), Wald pivot $t(\theta)$ (dashes), modified likelihood root $r^*(\theta)$ (heavy), and exact pivot $\theta \sum y_j$ (dot-dash). The modified likelihood root is indistinguishable from the exact pivot. The horizontal lines are at $0, \pm 1.96$. Right: corresponding significance functions, with horizontal lines at 0.025 and 0.975.

stat.epfl.ch

Autumn 2022 – slide 110

Vector case

- When θ is a vector and under classical asymptotics we base inference on the approximations

$$\hat{\theta} \sim \mathcal{N}_d(\theta_g, \hat{J}^{-1}), \quad w(\theta_g) = 2 \left\{ \ell(\hat{\theta}) - \ell(\theta_g) \right\} \sim \chi_d^2, \quad s(\theta_g) = \hat{J}^{-1/2} U(\theta_g) \sim \mathcal{N}_d(0, I_d),$$

with

- the first very commonly used for inferences on parameters;
- the second used to test whether $\theta = \theta_g$;
- the third much less used than the others, generally in the form $s(\theta_g)^T s(\theta_g) \sim \chi_d^2$.

- If θ divides into a $p \times 1$ **interest parameter** ψ and a $q \times 1$ **nuisance parameter** λ , then

$$\hat{\theta} = \begin{pmatrix} \hat{\psi} \\ \hat{\lambda} \end{pmatrix} \sim \mathcal{N}_{p+q} \left\{ \begin{pmatrix} \psi_g \\ \lambda_g \end{pmatrix}, \begin{pmatrix} \hat{J}_{\psi\psi} & \hat{J}_{\psi\lambda} \\ \hat{J}_{\lambda\psi} & \hat{J}_{\lambda\lambda} \end{pmatrix}^{-1} \right\},$$

where for brevity we now write $\hat{\lambda}_\psi = \max_\lambda \ell(\psi, \lambda)$, $\tilde{\theta} = \hat{\theta}_\psi = (\psi, \hat{\theta}_\psi)$,

$$\ell_\psi = \frac{\partial \ell(\theta)}{\partial \psi} \Big|_{\theta=\theta_g}, \quad \hat{J}_{\psi\psi} = -\hat{\ell}_{\psi\psi} = -\frac{\partial^2 \ell(\theta)}{\partial \psi \partial \psi^T} \Big|_{\theta=\hat{\theta}}, \quad \tilde{\ell}_{\psi\psi} = \frac{\partial^2 \ell(\theta)}{\partial \psi \partial \psi^T} \Big|_{\theta=\tilde{\theta}}, \quad \text{etc.}$$

stat.epfl.ch

Autumn 2022 – slide 111

Inference on ψ

- Under classical asymptotics and setting $\hat{J}^{\psi\psi} = (\hat{J}_{\psi\psi} - \hat{J}_{\psi\lambda} \hat{J}_{\lambda\lambda}^{-1} \hat{J}_{\lambda\psi})^{-1}$ we have

$$\begin{aligned}\hat{\psi} &\dot{\sim} \mathcal{N}_p(\psi_g, \hat{J}^{\psi\psi}) && \text{maximum likelihood estimator,} \\ w_p(\psi_g) &= 2 \left\{ \ell_p(\hat{\psi}) - \ell_p(\psi_g) \right\} \dot{\sim} \chi_p^2 && \text{(generalized) likelihood ratio statistic,} \\ s(\psi_g) &= \tilde{\ell}_{\psi}^T \hat{J}^{\psi\psi} \tilde{\ell}_{\psi} \dot{\sim} \chi_p^2 && \text{score statistic,}\end{aligned}$$

where we defined w_p using the **profile log likelihood** $\ell_p(\psi) = \ell(\psi, \hat{\lambda}_{\psi}) = \max_{\lambda} \ell(\psi, \lambda)$.

- If ψ is scalar ($p = 1$, the usual situation), the **likelihood root**

$$r(\psi_g) = \text{sign}(\hat{\psi} - \psi_g) \sqrt{w(\psi_g)} \dot{\sim} \mathcal{N}(0, 1).$$

- Properties:

- inferences using $w(\psi_g)$ and $r(\psi_g)$ are invariant to interest-respecting reparametrisation, so are preferable but more computationally burdensome;
- $s(\psi_g)$ is mainly used for tests, since only λ must be estimated (as $\psi = \psi_g$ is known).

- A $(1 - \alpha)$ confidence set based on $w_p(\psi_g)$ (or equivalently on $\ell_p(\psi)$) is

$$\{\psi : w_p(\psi) \leq \chi_p^2(1 - \alpha)\} = \left\{ \psi : \ell(\psi, \hat{\lambda}_{\psi}) \geq \ell(\hat{\psi}, \hat{\lambda}) - \frac{1}{2} \chi_p^2(1 - \alpha) \right\}.$$

Note: Large-sample distribution of the likelihood ratio statistic $w_p(\psi_g)$

□ We write

$$w_p(\psi_g) = 2\{\ell(\hat{\theta}) - \ell(\hat{\theta}_\psi)\} = 2\{\ell(\hat{\theta}) - \ell(\theta_g)\} - 2\{\ell(\hat{\theta}_\psi) - \ell(\theta_g)\}$$

and shall use Taylor series to approximate both terms by quadratic forms in $\hat{\theta} - \theta_g$ and $\hat{\lambda}_\psi - \lambda_g$.

□ We shall need to express ℓ_θ , ℓ_λ and $\hat{\lambda}_\psi - \lambda_g$ in terms of $\hat{\theta} - \theta_g$. Taylor expansion gives

$$0 = \tilde{\ell}_\theta = \ell_\theta + \ell_{\theta\theta}(\hat{\theta} - \theta_g) + \dots = \ell_\theta - \imath_{\theta\theta}(\hat{\theta} - \theta_g) + \dots,$$

where $\imath_{\theta\theta}$ denotes the expected information matrix evaluated at θ_g and \dots denotes terms of smaller order containing third derivatives. Likewise

$$0 = \tilde{\ell}_\lambda = \ell_\lambda + \ell_{\lambda\lambda}(\hat{\lambda}_\psi - \lambda_g) + \dots = \ell_\lambda - \imath_{\lambda\lambda}(\hat{\lambda}_\psi - \lambda_g) + \dots.$$

This implies that

$$\ell_\lambda \doteq \imath_{\lambda\psi}(\hat{\psi} - \psi_g) + \imath_{\lambda\lambda}(\hat{\lambda} - \lambda_g) = \imath_{\lambda\lambda}(\hat{\lambda}_\psi - \lambda_g),$$

so the necessary approximations are

$$\ell_\theta \doteq \imath_{\theta\theta}(\hat{\theta} - \theta_g), \quad \ell_\lambda \doteq \imath_{\lambda\lambda}(\hat{\lambda}_\psi - \lambda_g), \quad \hat{\lambda}_\psi - \lambda_g \doteq \hat{\lambda} - \lambda_g + \imath_{\lambda\lambda}^{-1} \imath_{\lambda\psi}(\hat{\psi} - \psi_g).$$

□ To obtain the quadratic forms we write

$$\begin{aligned} \ell(\hat{\theta}) &= \ell(\theta_g) + (\hat{\theta} - \theta_g)^T \ell_\theta + \frac{1}{2}(\hat{\theta} - \theta_g)^T \ell_{\theta\theta}(\hat{\theta} - \theta_g) + \dots \\ &\doteq \ell(\theta_g) + (\hat{\theta} - \theta_g)^T \imath_{\theta\theta}(\hat{\theta} - \theta_g) - \frac{1}{2}(\hat{\theta} - \theta_g)^T \imath_{\theta\theta}(\hat{\theta} - \theta_g), \end{aligned}$$

resulting in

$$2\{\ell(\hat{\theta}) - \ell(\theta_g)\} \doteq (\hat{\theta} - \theta_g)^T \imath_{\theta\theta}(\hat{\theta} - \theta_g),$$

and with a similar expression for $2\{\ell(\hat{\theta}_\psi) - \ell(\theta_g)\}$ we obtain

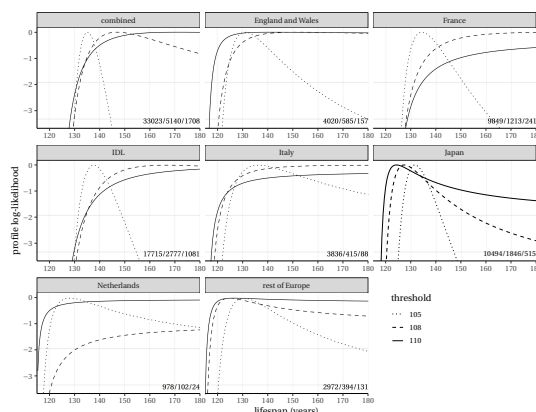
$$\begin{aligned} w_p(\psi_g) &\doteq (\hat{\theta} - \theta_g)^T \imath_{\theta\theta}(\hat{\theta} - \theta_g) - (\hat{\lambda}_\psi - \lambda_g)^T \imath_{\lambda\lambda}(\hat{\lambda}_\psi - \lambda_g) \\ &\doteq (\hat{\psi} - \psi_g)^T \imath_{\psi\psi}(\hat{\psi} - \psi_g) + 2(\hat{\psi} - \psi_g)^T \imath_{\psi\lambda}(\hat{\lambda} - \lambda_g) + (\hat{\lambda} - \lambda_g)^T \imath_{\lambda\lambda}(\hat{\lambda} - \lambda_g) \\ &\quad - \left\{ (\hat{\lambda} - \lambda_g) + \imath_{\lambda\lambda}^{-1} \imath_{\lambda\psi}(\hat{\psi} - \psi_g) \right\}^T \imath_{\lambda\lambda} \left\{ (\hat{\lambda} - \lambda_g) + \imath_{\lambda\lambda}^{-1} \imath_{\lambda\psi}(\hat{\psi} - \psi_g) \right\} \\ &= (\hat{\psi} - \psi_g)^T (\imath_{\psi\psi} - \imath_{\psi\lambda} \imath_{\lambda\lambda}^{-1} \imath_{\lambda\psi}) (\hat{\psi} - \psi_g), \end{aligned}$$

and as $\hat{\psi} \sim \mathcal{N}\{\psi_g, (\imath_{\psi\psi} - \imath_{\psi\lambda} \imath_{\lambda\lambda}^{-1} \imath_{\lambda\psi})^{-1}\}$, we see that $w_p(\psi_g) \sim \chi_p^2$, as claimed.

□ Arguments along the lines of Lemma 44 show that the terms dropped above all tend in probability to zero, and thus do not affect the approximation.

Example: Human lifespan

Example 46 The figure below shows profile log likelihoods for the endpoint ψ of a generalized Pareto distribution fitted to data on lifetimes of semi-supercentenarians from different databases, with thresholds at 105, 108, 110 years. Here λ is scalar, so $p = q = 1$, and the horizontal line at $-\frac{1}{2}\chi_1^2(0.95) = -1.92$ indicates 95% confidence regions.



From Belzile et al. (2022, *Annual Reviews of Statistics and its Application*).

Model selection

- The fact that the Kullback–Leibler divergence

$$\text{KL}(g, f) = \mathbb{E}_g\{\log g(Y) - \log f(Y)\} = \mathbb{E}_g \left[-\log \left\{ \frac{f(Y)}{g(Y)} \right\} \right] \geq 0,$$

is minimised when $f = g$ suggested that we compare competing models $\mathcal{F}_1, \dots, \mathcal{F}_M$ in terms of their maximised log likelihoods $\log f_m(y; \hat{\theta}_m) = \hat{\ell}_m$.

- But $\hat{\ell}_m$ should be penalized, because
 - $\hat{\ell}_m \geq \log f_m(y; \theta_m)$ even if \mathcal{F}_m is the true model class, and
 - enlarging θ_m will increase $\hat{\ell}_m$ even if further parameters are unnecessary.
- Akaike proposed minimising $2\mathbb{E}_g \mathbb{E}_g^+ \left[-\log \{f(Y^+; \hat{\theta})/g(Y^+)\} \right]$, where $Y^+, Y \stackrel{\text{iid}}{\sim} g$ are independent datasets. The idea is that if $\hat{\theta} = \hat{\theta}(Y)$ is estimated separately from Y^+ , there will be a penalty due to ‘missing θ_g ’ which will grow with $\dim(\theta)$ (picture ...)
- This leads to choosing m to minimise the **Akaike** or the **network** information criteria

$$\text{AIC}_m = 2(d_m - \hat{\ell}_m), \quad \text{NIC}_m = 2 \left\{ \text{tr}(\hat{K}_m \hat{J}_m^{-1}) - \hat{\ell}_m \right\},$$

where the first takes $\text{tr}(\hat{K}_m \hat{J}_m^{-1}) \approx d_m = \dim(\theta_m)$.

Note: Derivation of AIC/NIC

□ Now

$$2E_g E_g^+ \left[-\log \{f(Y^+; \hat{\theta})/g(Y^+)\} \right] = 2E_g^+ \{ \log g(Y^+) \} - 2E_g E_g^+ \{ \log f(Y^+; \hat{\theta}) \},$$

so we can ignore the first term in the minimisation over f . An unbiased estimator of the second term would be $2\ell^+(\hat{\theta})$, where ℓ^+ is the log likelihood based on Y^+ and $\hat{\theta}$ is based on Y , but the estimator we have available is $2\ell(\hat{\theta})$, in which the log likelihood and $\hat{\theta}$ are both based on Y . Clearly $\ell(\hat{\theta})$ is upwardly biased, but by how much?

□ To find out we consider the expectation over Y^+ and Y of

$$2 \{ \ell(\hat{\theta}) - \ell^+(\hat{\theta}) \} = 2 \{ \ell(\hat{\theta}) - \ell(\theta_g) \} + 2 \{ \ell(\theta_g) - \ell^+(\theta_g) \} + 2 \{ \ell^+(\theta_g) - \ell^+(\hat{\theta}) \}, \quad (5)$$

where as before θ_g is the best candidate parameter value under f .

□ As $\hat{\theta}$ maximises the log likelihood, $\ell_{\theta}(\hat{\theta}) = 0$, so the first term on the right-hand side of (5) is

$$\begin{aligned} 2 \{ \ell(\hat{\theta}) - \ell(\theta_g) \} &\doteq 2 \left\{ \ell(\hat{\theta}) - \ell(\hat{\theta}) - \ell_{\theta}(\hat{\theta})(\theta_g - \hat{\theta}) - \frac{1}{2}(\theta_g - \hat{\theta})^T \ell_{\theta\theta}(\hat{\theta})(\theta_g - \hat{\theta}) \right\} \\ &\doteq (\hat{\theta} - \theta_g)^T \iota_{\theta\theta}(\theta_g)(\hat{\theta} - \theta_g), \end{aligned}$$

where we have neglected terms that are $o_p(1)$. The expectation of this scalar equals that of its trace, and the large-sample normal distribution of $\hat{\theta}$ gives

$$\begin{aligned} E_g \left[\text{tr} \left\{ (\hat{\theta} - \theta_g)^T \iota_{\theta\theta}(\theta_g)(\hat{\theta} - \theta_g) \right\} \right] &= E_g \left[\text{tr} \left\{ (\hat{\theta} - \theta_g)(\hat{\theta} - \theta_g)^T \iota_{\theta\theta}(\theta_g) \right\} \right] \\ &\doteq \text{tr} \left\{ \iota_{\theta\theta}^{-1}(\theta_g) K(\theta_g) \iota_{\theta\theta}^{-1}(\theta_g) \iota_{\theta\theta}(\theta_g) \right\} \\ &= \text{tr} \left\{ K(\theta_g) \iota_{\theta\theta}^{-1}(\theta_g) \right\}. \end{aligned}$$

□ The second term on the right-hand side of (5) has expectation zero.

□ The third term on the right-hand side of (5) can be written as

$$2 \{ \ell^+(\theta_g) - \ell^+(\hat{\theta}) \} \doteq 2 \left\{ \ell^+(\theta_g) - \ell^+(\theta_g) - \ell_{\theta}^+(\theta_g)(\hat{\theta} - \theta_g) - \frac{1}{2}(\hat{\theta} - \theta_g)^T \ell_{\theta\theta}^+(\theta_g)(\hat{\theta} - \theta_g) \right\},$$

plus $o_p(1)$ terms. Now $E_g^+ \{ \ell_{\theta}^+(\theta_g) \} = 0$ and $E_g^+ \{ \ell_{\theta\theta}^+(\theta_g) \} = -\iota_{\theta\theta}(\theta_g)$, so

$$2E_g E_g^+ \{ \ell^+(\theta_g) - \ell^+(\hat{\theta}) \} \doteq E_g \left\{ (\hat{\theta} - \theta_g)^T \iota_{\theta\theta}(\theta_g)(\hat{\theta} - \theta_g) \right\} \doteq \text{tr} \left\{ K(\theta_g) \iota_{\theta\theta}^{-1}(\theta_g) \right\}.$$

□ Hence

$$2E_g E_g^+ \left[-\log f(Y^+; \hat{\theta}) \right] \doteq 2E_g E_g^+ \left[-\log f(Y; \hat{\theta}) \right] + 2 \text{tr} \left\{ K(\theta_g) \iota_{\theta\theta}^{-1}(\theta_g) \right\}.$$

If $K(\theta_g) \doteq \iota_{\theta\theta}(\theta_g)$, then this final expression can be estimated by $\text{AIC} = 2\{d - \ell(\hat{\theta})\}$, where $d = \dim(\theta)$, or by the *network information criterion* $\text{NIC} = 2\{\text{tr}(\hat{K}\hat{J}^{-1}) - \ell(\hat{\theta})\}$, though neither gives consistent estimation of the true model, which would require the penalty to grow with n . The calculations above rely on generic large-sample likelihood results, and could be improved in specific cases (e.g., with normal errors).

Dealing with nuisance parameters

- Profiling removes nuisance parameters, but the bias of $\hat{\psi}$ is $O(d^3/n)$ in general, and then we require $d = o(n^{1/3})$ for consistency. Hence accuracy may be low if $\dim(\lambda)$ is high.
- Other approaches to dealing with λ include:
 - basing inference on a **marginal likelihood** or a **conditional likelihood**,

$$f(y; \psi, \lambda) = f(w; \psi) \times f(y | w; \psi, \lambda) = f(y | w_\psi; \psi) \times f(w_\psi; \psi, \lambda),$$

where w_ψ may not depend on ψ (recall Lemmas 16 and 17) — OK for any configuration of λ s, but may lose information on ψ ;

- constructing a **partial likelihood** (like the above, but harder to build);
 - **higher-order inference** such as using a **modified profile likelihood**, which can approximate both conditional and marginal likelihoods;
 - taking $\lambda \sim h(\cdot)$ and using the **integrated likelihood** $\int f(y; \psi, \lambda) h(\lambda) d\lambda$ — depends on h , like Bayesian inference;
 - using **orthogonal parameters**, i.e., mapping $\lambda \mapsto \zeta(\lambda, \psi)$ which is orthogonal to ψ ; or
 - using a **composite likelihood** in which λ does not appear.
- Below we sketch some of these.

Modified profile likelihood

- Replace profile likelihood $\exp\{\ell_p(\psi)\}$ by the **modified profile likelihood**

$$L_{\text{mp}}(\psi) = \exp\{\ell_{\text{mp}}(\psi)\} = M(\psi)L_p(\psi),$$

with $M(\psi)$ chosen to mimic properties of marginal or conditional likelihood.

- Taking

$$M(\psi) = \left| J_{\lambda\lambda}(\psi, \hat{\lambda}_\psi) \right|^{-1/2} \left| \frac{\partial \hat{\lambda}}{\partial \hat{\lambda}_\psi^\top} \right|$$

does this in some generality.

- The
 - first term can be obtained numerically if need be, but
 - the second term is hard to compute in general.
- Simpler to base a likelihood on the normal distribution of the modified likelihood root $r^*(\psi)$ (next).

Higher-order inference . . .

- Classical theory gives first-order accuracy. With ψ scalar using the likelihood root gives

$$P \{r(\psi_g) \leq r^o(\psi_g)\} = \Phi\{r^o(\psi)\} + O(n^{-1/2}),$$

so tests and confidence sets based on data y^o have error $n^{-1/2}$.

- If we use the **modified likelihood root**,

$$r^*(\psi) = r(\psi) + \frac{1}{r(\psi)} \log \left\{ \frac{q(\psi)}{r(\psi)} \right\},$$

where $q(\psi)$ depends on the model, then the error drops to $O(n^{-3/2})$ for continuous responses and to $O(n^{-1})$ for discrete responses, so

$$P \{r^*(\psi_g) \leq r^{*o}(\psi_g)\} = \Phi\{r^{*o}(\psi_g)\} + O(n^{-3/2}),$$

for continuous data (often almost exact even for tiny n ; see Example 45).

- Highly accurate even into the distribution tails, because the relative error is bounded.
- A $1 - 2\alpha$ confidence interval,

$$\{\psi : z_\alpha \leq r^{*o}(\psi) \leq 1 - \alpha\},$$

has error of order $n^{-3/2}$ (often effectively perfect).

. . . with nuisance parameters

- With nuisance parameters, $r(\psi) = \text{sign}(\hat{\psi} - \psi) \sqrt{w_p(\psi)}$, and

$$q(\psi) = \frac{|\varphi(\hat{\theta}) - \varphi(\hat{\theta}_\psi)|}{|\varphi(\hat{\theta})|} \left\{ \frac{|\hat{J}|}{|J_{\lambda\lambda}(\hat{\theta}_\psi)|} \right\}^{1/2}$$

where φ is the $d \times 1$ canonical parameter of a local exponential family approximation to the model at the observed data y^o , with $\varphi_\theta(\theta) = \partial \varphi(\theta) / \partial \theta^T$, etc.

- In a general exponential family $\varphi(\theta)$ is the canonical parameter, and in a linear exponential family,

$$q(\psi) = (\hat{\psi} - \psi) \left\{ \frac{|\hat{J}|}{|J_{\lambda\lambda}(\hat{\theta}_\psi)|} \right\}^{1/2}.$$

- In general for independent continuous observations we write

$$\varphi(\theta)_{d \times 1} = V_{d \times n}^T \frac{\partial \ell(\theta; y)}{\partial y} \Big|_{y=y^o} = \sum_{j=1}^n V_j^T \frac{\partial \log f(y_j; \psi, \lambda)}{\partial y_j} \Big|_{y=y^o},$$

where the $1 \times d$ $V_j = \partial y_j / \partial \theta^T$ are evaluated at y^o and $\hat{\theta}^o$.

Properties of higher order approximations

- ☐ Invariant to interest-respecting reparameterization.
- ☐ Computation almost as easy as first order versions.
- ☐ Error $O(n^{-3/2})$ in continuous response models, $O(n^{-1})$ in discrete response models.
- ☐ Relative (not absolute) error, so highly accurate in tails.
- ☐ Bayesian version is also available (and easier to derive).

Example 47 (Location-scale model) Compute $\varphi(\theta)$ for a location-scale model, in which independent observations Y_j have density $\tau^{-1}h\{(y - \eta)/\tau\}$. What about the normal density?

stat.epfl.ch

Autumn 2022 – slide 119

Note to Example 47

- ☐ In this case the overall log likelihood is

$$\ell(\eta, \tau) = -n \log \tau + \sum_{j=1}^n \log h\{(y_j - \eta)/\tau\},$$

so the vector $\partial \ell(\eta, \tau)/\partial y$ has components $\tau^{-1}(\log h)' \{(y_j - \eta)/\tau\}$, evaluated at the maximum likelihood estimates $\hat{\eta}^o$ and $\hat{\tau}^o$ and observed data vector y_1^o, \dots, y_n^o .

- ☐ To compute the V_j we use the structural expression $y = \eta + \tau \varepsilon$, where $\varepsilon \sim h$. This represents y as a function of $\theta^T = (\eta, \tau)$, and yields $\partial y_j / \partial \theta^T = (1, \varepsilon_j)$. This has to be evaluated at the observed data point y^o , and at that point the parameters are replaced by their maximum likelihood estimates, giving $V_j^T = (1, (y_j^o - \hat{\eta}^o)/\hat{\tau}^o)$.
- ☐ This yields

$$\varphi(\theta) = \sum_{j=1}^n \tau^{-1} (\log h)' \{(y_j^o - \eta)/\tau\} (1, e_j)^T,$$

where we have set $e_j = (y_j^o - \hat{\eta}^o)/\hat{\tau}^o$.

- ☐ If h is normal, then $\log h(u) \equiv -u^2/2$, so $(\log h)' \{(y_j^o - \eta)/\tau\} = -(y_j^o - \eta)/\tau^2$, leading to

$$\varphi(\theta)^T = \left(\sum_{j=1}^n (\eta - y_j^o)/\tau^2, \sum_{j=1}^n (\eta - y_j^o)/\tau^2 \times e_j \right) \equiv (\eta/\tau^2, 1/\tau^2),$$

because it turns out that inferences are invariant under non-singular affine transformations of $\varphi(\theta)$ (exercise).

stat.epfl.ch

Autumn 2022 – note 1 of slide 119

Orthogonal parameters

- If the expected information matrix is block diagonal, with $i_{\psi,\lambda}(\theta) = 0$ for all θ , then $\hat{\psi}$ is asymptotically independent of $\hat{\lambda}$, and we can hope that the effect on $\hat{\psi}$ of estimating λ will be limited. If so, we say that ψ and λ are **orthogonal**.
- This suggests mapping (ψ, γ) to (ψ, λ) , where $\lambda = \lambda(\psi, \gamma)$ is orthogonal to ψ .
- Writing $\gamma = \gamma(\psi, \lambda)$ gives

$$\ell(\psi, \lambda) = \ell^* \{ \psi, \gamma(\psi, \lambda) \},$$

and differentiation with respect to ψ and λ leads to

$$\frac{\partial^2 \ell}{\partial \lambda \partial \psi} = \frac{\partial \gamma^T}{\partial \lambda} \frac{\partial^2 \ell^*}{\partial \gamma \partial \psi} + \frac{\partial \gamma^T}{\partial \lambda} \frac{\partial^2 \ell^*}{\partial \gamma \partial \gamma^T} \frac{\partial \gamma}{\partial \psi} + \frac{\partial^2 \gamma^T}{\partial \lambda \partial \psi} \frac{\partial \ell^*}{\partial \gamma}.$$

- For orthogonality this must have expectation zero, so

$$0 = \frac{\partial \gamma^T}{\partial \lambda} i_{\gamma\psi}^* + \frac{\partial \gamma^T}{\partial \lambda} i_{\gamma\gamma}^* \frac{\partial \gamma}{\partial \psi},$$

where $i_{\gamma\psi}^*$ and $i_{\gamma\gamma}^*$ are components of the expected information matrix in the non-orthogonal parametrization, so λ solves the system of q PDEs

$$\frac{\partial \gamma}{\partial \psi} = -i_{\gamma\gamma}^{*-1}(\psi, \gamma) i_{\gamma\psi}^*(\psi, \gamma).$$

Orthogonal parameters II

- A (possibly numerical) solution always exists when $\dim(\psi) = 1$, but need not exist when ψ is vector, because then we must simultaneously solve

$$\frac{\partial \gamma}{\partial \psi_1} = -i_{\gamma\gamma}^{*-1}(\psi, \gamma) i_{\gamma\psi_1}^*(\psi, \gamma), \quad \frac{\partial \gamma}{\partial \psi_2} = -i_{\gamma\gamma}^{*-1}(\psi, \gamma) i_{\gamma\psi_2}^*(\psi, \gamma),$$

for all γ , ψ_1 and ψ_2 , but the compatibility condition

$$\frac{\partial^2 \gamma}{\partial \psi_1 \partial \psi_2} = \frac{\partial^2 \gamma}{\partial \psi_2 \partial \psi_1}$$

may fail.

Example 48 (Linear exponential family) What parameter is orthogonal to ψ in the linear exponential family with log likelihood

$$\ell^*(\psi, \gamma) \equiv s_1^T \psi + s_2^T \gamma - k(\psi, \gamma)?$$

Consider normal and Poisson likelihoods in particular.

Note to Example 48

- The parameters $\lambda = \lambda(\psi, \gamma)$ orthogonal to ψ are determined by

$$\frac{\partial \gamma}{\partial \psi^T} = -k_{\gamma\gamma}^{-1}(\psi, \gamma)k_{\gamma\psi}(\psi, \gamma). \quad (6)$$

If we reparametrize in terms of ψ and $\lambda = k_{\gamma}(\psi, \gamma) = \partial k(\psi, \gamma)/\partial \gamma$, then in this new parametrization, γ is a function of ψ and λ , and

$$0 = \frac{\partial \lambda^T}{\partial \psi} = \frac{\partial \gamma^T}{\partial \psi} k_{\gamma\gamma}(\psi, \gamma) + k_{\psi\gamma}(\psi, \gamma),$$

so $\lambda = k_{\gamma}(\psi, \gamma)$ is a solution to (6). That is, the parameter orthogonal to ψ is the so-called complementary mean parameter $\lambda(\psi, \gamma) = E(S_2; \psi, \gamma)$. By symmetry, $E(S_1; \psi, \gamma)$ is orthogonal to γ .

- The normal distribution with mean μ and variance σ^2 has canonical parameter $(\mu/\sigma^2, -1/(2\sigma^2))$. The canonical statistic (Y, Y^2) has expectation $(\mu, \mu^2 + \sigma^2)$, so μ is orthogonal to $-1/(2\sigma^2)$, and hence to σ^2 , while μ/σ^2 is orthogonal to $\mu^2 + \sigma^2$.
- Independent Poisson variables Y_1 and Y_2 with means $\exp(\gamma)$ and $\exp(\gamma + \psi)$ have log likelihood

$$\ell^*(\psi, \gamma) \equiv (y_1 + y_2)\gamma + y_2\psi - e^{\gamma} - e^{\gamma+\psi}.$$

The discussion above suggests that

$$\lambda = E(Y_1 + Y_2) = \exp(\gamma) + \exp(\gamma + \psi) = e^{\gamma}(1 + e^{\psi})$$

is orthogonal to ψ , so $\gamma = \log \lambda - \log(1 + e^{\psi})$ and

$$\ell(\psi, \lambda) \equiv y_2\psi - (y_1 + y_2)\log(1 + e^{\psi}) + (y_1 + y_2)\log \lambda - \lambda.$$

The separation of ψ and λ implies that the profile and modified profile likelihoods for ψ are proportional. They correspond to the conditional likelihood obtained from the density of Y_2 given $Y_1 + Y_2$.

Composite likelihood

- Used when full likelihood can't be computed but densities for distinct subsets of the observations, y_{S_1}, \dots, y_{S_C} , are available, can use a **composite (log) likelihood**

$$\ell_C(\theta) = \sum_{c=1}^C \log f(y_{S_c}; \theta).$$

- The choice of subsets S_1, \dots, S_C determines what parameters can be estimated.
- Special cases:
 - **independence likelihood** takes $S_j = \{y_j\}$ and treats (possibly dependent) y_j as independent;
 - **pairwise likelihood** uses subsets of distinct pairs $\{y_j, y_{j'}\}$.
- May be useful with spatial data, and then contributions from distant pairs may be downweighted or dropped entirely.
- $\ell_C(\theta)$ satisfies the first Bartlett identity, so can give consistent estimators $\tilde{\theta}$, but requires a sandwich variance matrix (or some other approach) to estimate $\text{var}(\tilde{\theta})$.
- Model comparisons use the **composite likelihood information criterion**

$$\text{CLIC} = 2 \left[\text{tr}\{K(\tilde{\theta})J(\tilde{\theta})^{-1}\} - \ell_C(\tilde{\theta}) \right].$$

stat.epfl.ch

Autumn 2022 – slide 122

Empirical likelihood

- Empirical likelihood allows nonparametric estimation of constrained distributions.
- If $y_1, \dots, y_n \stackrel{\text{iid}}{\sim} G$, then the \hat{G} that maximises the 'nonparametric likelihood'

$$L(G) = \prod_{j=1}^n G(dy_j) = \prod_{j=1}^n p_j, \quad \text{subject to } p_j \geq 0, \quad \sum_{j=1}^n p_j \leq 1,$$

sets $G(dy_j) = \hat{p}_j \equiv n^{-1}$: \hat{G} is the empirical distribution function of y_1, \dots, y_n .

- Adding the constraint $E\{c(Y; \theta)\} = 0$ leads to maximising

$$\sum_{j=1}^n \log p_j \quad \text{subject to } p_j \geq 0, \quad \sum_{j=1}^n p_j \leq 1, \quad \sum_{j=1}^n p_j c(y_j; \theta) = 0,$$

and a use of Lagrange multipliers shows that we must find $a = a_\theta$ to solve

$$\sum_{j=1}^n \frac{c(y_j; \theta)}{n\{1 + ac(y_j; \theta)\}} = 0$$

giving $\hat{a} = \hat{a}_\theta$, $\hat{p}_j(\theta) = n^{-1}/\{1 + \hat{a}c(y_j; \theta)\}$ and empirical likelihood ratio statistic $w(\theta) = 2 \sum \log\{1 + \hat{a}c(y_j; \theta)\}$

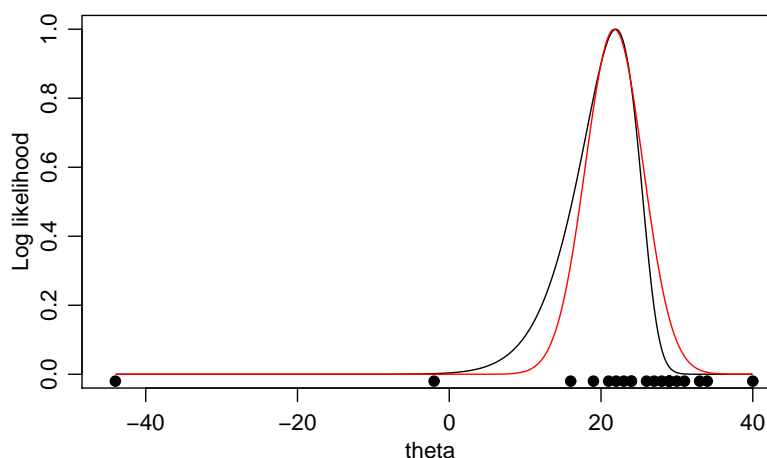
- The usual χ^2 result applies to $w(\theta)$. Can be widely generalised ...

stat.epfl.ch

Autumn 2022 – slide 123

Example: Newcomb data

$n = 20$ observations on the speed of light (made by Simon Newcomb), empirical likelihood (black) and normal likelihood (red) for the mean θ . Note how the empirical likelihood adapts to the 'shape' of the data.



stat.epfl.ch

Autumn 2022 – slide 124

Non-regular models

- The regularity conditions (C1)–(C4) apply in many settings met in practice, but not universally. The most common failures arise when
 - some of the parameters are discrete (e.g., change point problems),
 - the model is not identifiable (distinct θ values give the same model),
 - θ_g is on the boundary of the parameter space (e.g., testing for a zero variance),
 - $d = \dim(\theta)$ grows (too fast) with n , or
 - the support of $f(y; \theta)$ depends on θ (so the Bartlett identities fail).
- Even when the conditions are satisfied there can be datasets for which maximum likelihood estimation fails, e.g.,
 - there is no unique maximum to the likelihood, or
 - the maximum is on the edge of the parameter space,and then penalisation (equivalent to using a prior) is often used.

Example 49 If $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} U(0, \theta)$, find the likelihood and $\hat{\theta}$, and show that the limit distribution of $n(\theta - \hat{\theta})/\theta$ when $n \rightarrow \infty$ is $\exp(1)$. Discuss.

stat.epfl.ch

Autumn 2022 – slide 125

Note to Example 49

Owing to the independence,

$$L(\theta) = \prod_{j=1}^n f_Y(y_j; \theta) = \prod_{j=1}^n \{\theta^{-1} I(0 < y_j < \theta)\} = \theta^{-n} I(\max y_j < \theta), \quad \theta > 0,$$

and therefore $\hat{\theta} = M = \max Y_j$, whose distribution is

$$P(M \leq x) = (x/\theta)^n, \quad 0 < x < \theta.$$

Now

$$P\left\{n(\theta - \hat{\theta})/\theta \leq x\right\} = P(\hat{\theta} \geq \theta - x\theta/n) = 1 - \{(\theta - x\theta/n)/\theta\}^n \rightarrow 1 - \exp(-x),$$

as required. Note that:

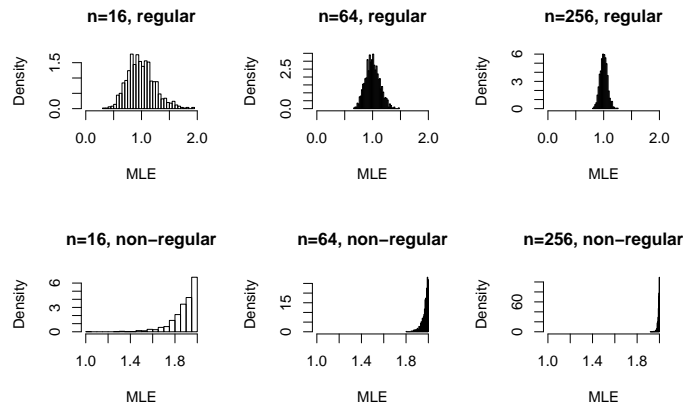
- ☐ the scaling needed to get a limiting distribution is much faster here than in the regular case (we have to multiply by n to get a non-degenerate limit);
- ☐ the limit is not normal.

stat.epfl.ch

Autumn 2022 – note 1 of slide 125

Uniform example

Comparison of the distributions of $\hat{\theta}$ in a regular case (panels above, with standard deviation $\propto n^{-1/2}$) and in a nonregular case (Example 49, panels below, with standard deviation $\propto n^{-1}$). In other nonregular cases it might happen that the distribution is nasty (unlike here) and/or that the convergence is slower than in regular cases.



stat.epfl.ch

Autumn 2022 – slide 126

Comments

- Other likelihoods and/or likelihood-like functions are widely used, especially
 - **partial likelihood**, used to eliminate nuisance functions for inference (survival data),
 - **quasi-likelihood**, used to model over-dispersion in exponential family models,
 - **pseudo-likelihood**, treats data as Gaussian even when they are not (econometrics).
- Strengths of likelihood approach:
 - heuristic as plausibility of a model as explanation of data;
 - we 'just' have to write down the density of the observed data;
 - invariance to data and parameter transformations;
 - simple and (fairly) general approximate theory for inference under regularity conditions, also easily implemented numerically;
 - large-sample optimality properties in regular cases;
 - close links to Bayesian inference (next).
- Weaknesses of likelihood approach:
 - requires 'parametric' model for data;
 - can have trouble in high-dimensional settings;
 - not all models are regular.

Parameters and functionals

- ☐ Parametric models are determined by a finite vector $\theta \in \Theta$. Does this generalise?
- ☐ If $Y \sim G$, then we can define a parameter in terms of a **statistical functional**, e.g.,

$$\mu = t_1(G) = \int y \, dG(y), \quad \sigma^2 = t_2(G) = \int y^2 \, dG(y) - \left\{ \int y \, dG(y) \right\}^2.$$

- ☐ Below we always assume that such functionals are well-defined.
- ☐ We apply the '**plug-in principle**' and replace G by an estimator \hat{G} , giving

$$\hat{\mu} = t_1(\hat{G}) = \int y \, d\hat{G}(y), \quad \hat{\sigma}^2 = t_2(\hat{G}) = \int y^2 \, d\hat{G}(y) - \left\{ \int y \, d\hat{G}(y) \right\}^2.$$

- ☐ With a parametric model we can write $G \equiv G_\theta$ and $\hat{G} \equiv G_{\hat{\theta}}$, but a general estimator of G based on $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} G$ is the **empirical distribution function (EDF)**

$$\hat{G}(y) = \frac{1}{n} \sum_{j=1}^n H(y - Y_j), \quad H(x) = \begin{cases} 0, & x < 0, \\ 1, & x \geq 0, \end{cases}$$

where $H(\cdot)$ is the **Heaviside function**.

stat.epfl.ch

Autumn 2022 – slide 129

Algorithmic approach

Example 50 Give general definitions of the median and the parameter obtained from a maximum likelihood fit of a density $f(y; \theta)$. What are the corresponding estimators (a) under a fitted exponential model, and (b) a nonparametric model?

- ☐ This approach is essentially algorithmic: $t(\cdot)$ is an algorithm that
 - when applied to the distribution G gives the parameter $t(G)$;
 - when applied to an estimator \hat{G} based on data Y_1, \dots, Y_n gives the estimator $t(\hat{G})$.
- ☐ The algorithm $t(\cdot)$ can be (almost) arbitrarily complex.
- ☐ This point of view suggests a sampling approach to frequentist inference:
 - if we knew G , we could assess the properties of $t(\hat{G})$ by generating many samples $\hat{G} \equiv \{Y_1, \dots, Y_n\}$ from G and looking at the corresponding values of $t(\hat{G})$;
 - since G is unknown, we replace it by \hat{G} , generate samples $\hat{G}^* \equiv \{Y_1^*, \dots, Y_n^*\}$ from \hat{G} , and use the corresponding values of $t(\hat{G}^*)$ to estimate the distribution of $t(\hat{G})$.
- ☐ The samples $\hat{G}^* \equiv \{Y_1^*, \dots, Y_n^*\}$ are known as **bootstrap samples**, and the overall procedure is known as a **bootstrap**, one of many possible **resampling** procedures.

stat.epfl.ch

Autumn 2022 – slide 130

Example 50

- The usual definition of the p quantile is

$$t_1(G) = \inf\{x : G(x) \geq p\},$$

for $p \in (0, 1)$. For the median we set $p = 1/2$.

- The maximum likelihood estimator is defined as

$$t_2(G) = \arg \max_{\theta} E_G\{\log f(Y; \theta)\} = \arg \max_{\theta} \int \log f(y; \theta) d\hat{G}(y),$$

which we earlier called θ_g .

- Under an exponential model

$$t_1(G) = \inf\{x : 1 - \exp(-\lambda x) \geq p\} = -\lambda^{-1} \log(1 - p) = \lambda^{-1} \log 2,$$

so if the fitted model has parameter $\hat{\lambda}$, then $t_1(\hat{G}) = \hat{\lambda}^{-1} \log 2$.

Likewise θ_g is estimated by

$$\arg \max_{\theta} \int \log f(y; \theta) \hat{\lambda} e^{-\hat{\lambda} y} dy;$$

note that f is not necessarily exponential.

- Under the general model and with order statistics $Y_{(1)} \leq Y_{(2)} \leq \dots \leq Y_{(n)}$,

$$t_1(\hat{G}) = \inf\{x : \hat{G}(x) \geq p\} = Y_{(m)},$$

where $m = \lfloor (n+1)/2 \rfloor$, and as $dH(u)$ puts a unit mass at $u = 0$,

$$\begin{aligned} t_2(\hat{G}) &= \arg \max_{\theta} \int \log f(y; \theta) d\hat{G}(y) \\ &= \arg \max_{\theta} \int \log f(y; \theta) d \left\{ n^{-1} \sum_{j=1}^n H(y - Y_j) \right\} \\ &= \arg \max_{\theta} n^{-1} \sum_{j=1}^n \int \log f(y; \theta) dH(y - Y_j) \\ &= \arg \max_{\theta} n^{-1} \sum_{j=1}^n \log f(Y_j; \theta), \end{aligned}$$

i.e., the maximum likelihood estimator of θ based on the sample.

Example: Handedness data

Table 1: Data from a study of handedness; *hand* is an integer measure of handedness, and *dnan* a genetic measure. Data due to Dr Gordon Claridge, University of Oxford.

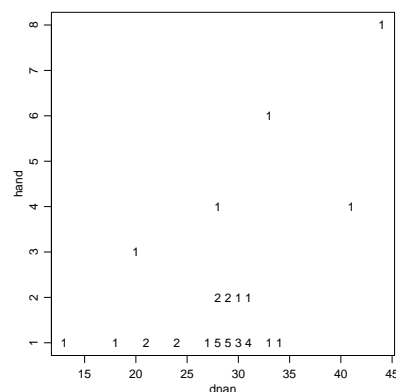
	dnan	hand	dnan	hand	dnan	hand	dnan	hand
1	13	1	11	28	1	21	29	2
2	18	1	12	28	2	22	29	1
3	20	3	13	28	1	23	29	1
4	21	1	14	28	4	24	30	1
5	21	1	15	28	1	25	30	1
6	24	1	16	28	1	26	30	2
7	24	1	17	29	1	27	30	1
8	27	1	18	29	1	28	31	1
9	28	1	19	29	1	29	31	1
10	28	2	20	29	2	30	31	1

stat.epfl.ch

Autumn 2022 – slide 131

Example: Handedness data

Scatter plot of handedness data. The numbers show the multiplicities of the observations.



stat.epfl.ch

Autumn 2022 – slide 132

Example: Handedness data

- How do we quantify dependence between `dnan` and `hand` for these $n = 37$ individuals?
- A standard measure is the **product-moment (Pearson) correlation** for $G(u, v)$, i.e.,

$$\theta = t(G) = \frac{\int \{u - \int u dG(u, v)\} \{v - \int v dG(u, v)\} dG(u, v)}{\left[\int \{u - \int u dG(u, v)\}^2 dG(u, v) \int \{v - \int v dG(u, v)\}^2 dG(u, v) \right]^{1/2}}.$$

- With $(u, v) = (\text{dnan}, \text{hand})$, the sample version is

$$\begin{aligned} \hat{\theta} = t(\hat{G}) &= \frac{\sum_{j=1}^n (\text{dnan}_j - \overline{\text{dnan}})(\text{hand}_j - \overline{\text{hand}})}{\left\{ \sum_{j=1}^n (\text{dnan}_j - \overline{\text{dnan}})^2 \sum_{j=1}^n (\text{hand}_j - \overline{\text{hand}})^2 \right\}^{1/2}} \\ &= 0.509. \end{aligned}$$

- Standard (bivariate normal) 95% confidence interval is (0.221, 0.715), but this is obviously inappropriate (the data look highly non-normal).
- Try simulation approach ...

stat.epfl.ch

Autumn 2022 – slide 133

Bootstrap simulation

- Whether \hat{G} is parametric or non-parametric, we simulate as follows:

- For $r = 1, \dots, R$:
 - ▷ generate a bootstrap sample $y_1^*, \dots, y_n^* \stackrel{\text{iid}}{\sim} \hat{G}$,
 - ▷ compute $\hat{\theta}_r^*$ using y_1^*, \dots, y_n^* ,
- so the output is a set of **bootstrap replicates**,

$$\hat{\theta}_1^*, \dots, \hat{\theta}_R^*.$$

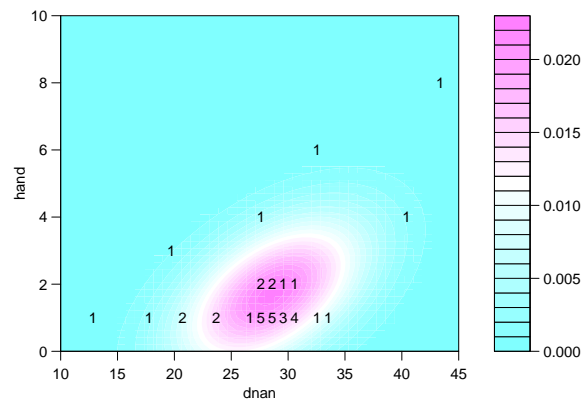
- We then use $\hat{\theta}_1^*, \dots, \hat{\theta}_R^*$ to estimate properties of $\hat{\theta}$ (histogram, ...).
- If $R \rightarrow \infty$, then get perfect match to theoretical calculation based on \hat{G} (if this is available): Monte Carlo error disappears completely.
- In practice R is finite, so some Monte Carlo error remains.
- If \hat{G} is the EDF, then $y_1^*, \dots, y_n^* \stackrel{\text{iid}}{\sim} \hat{G}$ are sampled with replacement and equal probabilities from y_1, \dots, y_n , so if $f_i^* = \#\{y_j^* = y_i\}$, then (f_1^*, \dots, f_n^*) has the multinomial distribution with denominator n and probability vector (n^{-1}, \dots, n^{-1}) .
- Although $E^*(f_j^*) = 1$, y_j can appear 0, 1, ..., n times in the bootstrap sample.

stat.epfl.ch

Autumn 2022 – slide 134

Handedness data: Fitted bivariate normal model

Contours of bivariate normal distribution fitted to handedness data; parameter estimates are $\hat{\mu}_1 = 28.5$, $\hat{\mu}_2 = 1.7$, $\hat{\sigma}_1 = 5.4$, $\hat{\sigma}_2 = 1.5$, $\hat{\rho} = 0.509$. The data are also shown.

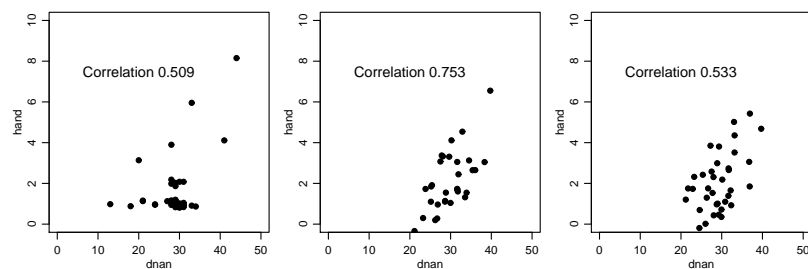


stat.epfl.ch

Autumn 2022 – slide 135

Handedness data: Parametric bootstrap samples

Left: original data, with jittered vertical values. Centre and right: two samples generated from the fitted bivariate normal distribution.

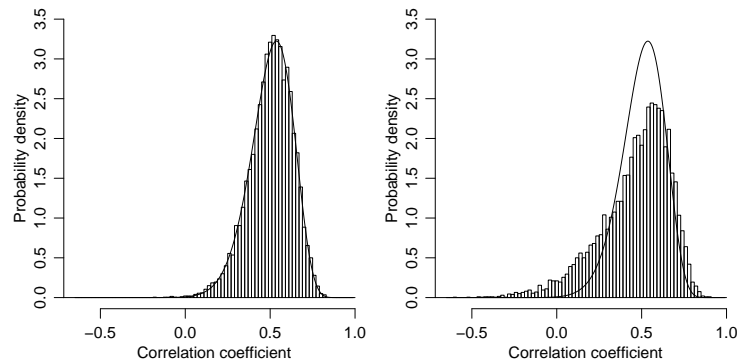


stat.epfl.ch

Autumn 2022 – slide 136

Handedness data: Correlation coefficient

Bootstrap distributions with $R = 10000$. Left: simulation from fitted bivariate normal distribution. Right: nonparametric sampling from the EDF. The lines show the theoretical probability density function of the correlation coefficient under sampling from a fitted bivariate normal distribution.

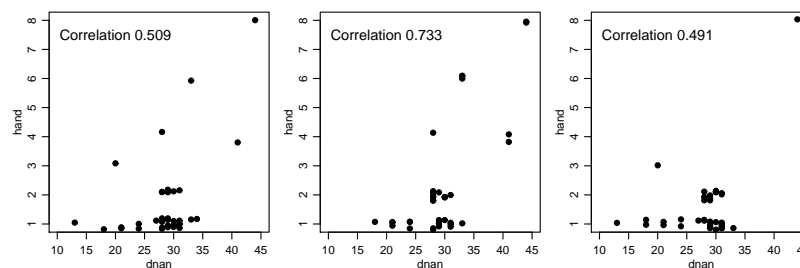


stat.epfl.ch

Autumn 2022 – slide 137

Handedness data: Bootstrap samples

Left: original data, with jittered vertical values. Centre and right: two bootstrap samples, with jittered vertical values.



stat.epfl.ch

Autumn 2022 – slide 138

Using the $\hat{\theta}^*$

- The **bias** and **variance** of $\hat{\theta}$ as an estimator of $\theta = t(G)$,

$$\beta(G) = E(\hat{\theta} \mid y_1, \dots, y_n \stackrel{\text{iid}}{\sim} G) - t(G), \quad \nu(G) = \text{var}(\hat{\theta} \mid G),$$

are estimated by replacing the unknown G by its known estimate \hat{G} :

$$\beta(\hat{G}) = E(\hat{\theta} \mid y_1, \dots, y_n \stackrel{\text{iid}}{\sim} \hat{G}) - t(\hat{G}), \quad \nu(\hat{G}) = \text{var}(\hat{\theta} \mid y_1, \dots, y_n \stackrel{\text{iid}}{\sim} \hat{G}).$$

- The Monte Carlo approximations to $\beta(\hat{G})$ and $\nu(\hat{G})$ are

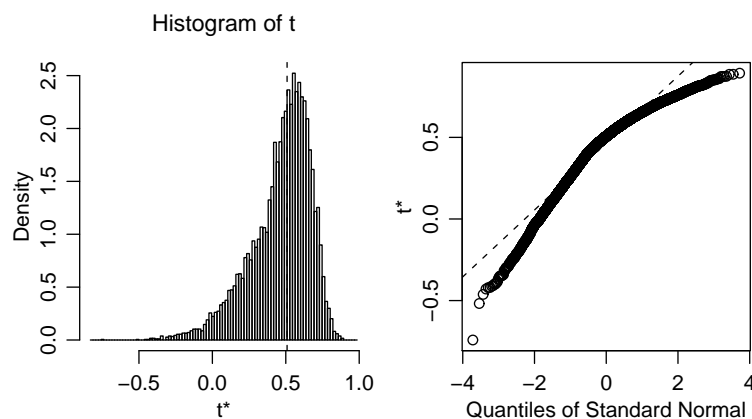
$$b = \bar{\hat{\theta}^*} - \hat{\theta} = R^{-1} \sum_{r=1}^R \hat{\theta}_r^* - \hat{\theta}, \quad v = \frac{1}{R-1} \sum_{r=1}^R (\hat{\theta}_r^* - \bar{\hat{\theta}^*})^2.$$

For the handedness data, $R = 10^4$ and $b = -0.046$, $v = 0.043 = 0.205^2$.

- We estimate the **p quantile** of $\hat{\theta}$ using the p quantile of $\hat{\theta}_1^*, \dots, \hat{\theta}_R^*$, i.e., $\hat{\theta}_{((R+1)p)}^*$.

Handedness data

Summaries of the $\hat{\theta}^*$. Left: histogram, with vertical line showing $\hat{\theta}$. Right: normal Q–Q plot of $\hat{\theta}^*$.



Common questions

- ☐ **How big should n be?** — depends on the context
- ☐ **What if the sample is unrepresentative?** — this is always a potential problem in statistics, not specific to resampling methods.
- ☐ **How big should R be?** — at least 1000 for most purposes
- ☐ **Why take resamples of size n ?**
 - We usually want to mimic the sampling properties of samples like the original one, so take resamples of size n ,
 - but sometimes we take resamples of size $m \ll n$ in order to achieve validity of the bootstrap—e.g., for extreme quantiles.
- ☐ **Why resample from the EDF?**
 - The EDF is the nonparametric MLE of G , so is a natural choice, but
 - sometimes (e.g., testing) we resample from a constrained version of \hat{G} ,
 - sometimes it may be useful to smooth \hat{G} ;
 - sometimes it may be useful to simulate from (several) parametric fits.

stat.epfl.ch

Autumn 2022 – slide 141

How big should n be?

- ☐ For the **average** $\hat{\theta} = \bar{y}$, the number of distinct samples is

$$m_n = \binom{2n-1}{n},$$

the most probable of which has probability $p_n = n!/n^n$.

For $n > 12$, we have $m_n > 10^6$ and $p_n < 6 \times 10^{-5}$.

- ☐ Bootstrapping of smooth statistics like the average will often work OK provided $n > 20$.
- ☐ For the **median** of a sample of size $n = 2m + 1$, the possible distinct values of $\hat{\theta}^*$ are $y_{(1)} < \dots < y_{(n)}$, and

$$P^*(\hat{\theta}^* > y_{(l)}) = \sum_{r=0}^m \binom{n}{r} \left(\frac{l}{n}\right)^r \left(1 - \frac{l}{n}\right)^{n-r},$$

so exact calculations of the variance etc. are possible.

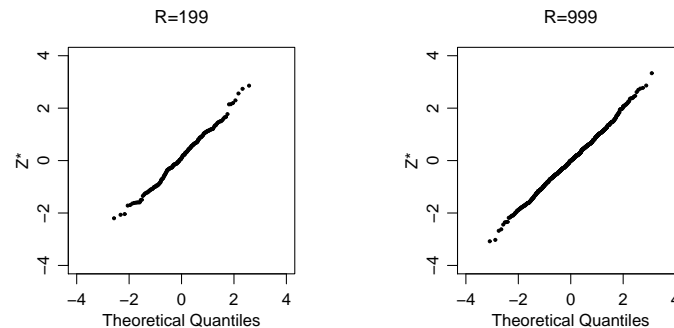
- ☐ However the median is very vulnerable to bad sample values, so for the median (and other 'non-smooth' statistics) much larger n is needed for reliable inference.

stat.epfl.ch

Autumn 2022 – slide 142

How many bootstraps?

- ☐ Must estimate moments and quantiles of $\hat{\theta}$ and derived quantities. Often feasible to take $R \gg 1000$
- ☐ Need $R \geq 200$ to estimate bias, variance, etc.
- ☐ Need $R \gg 100$, preferably $R \geq 2500$ to estimate quantiles needed for 95% confidence intervals

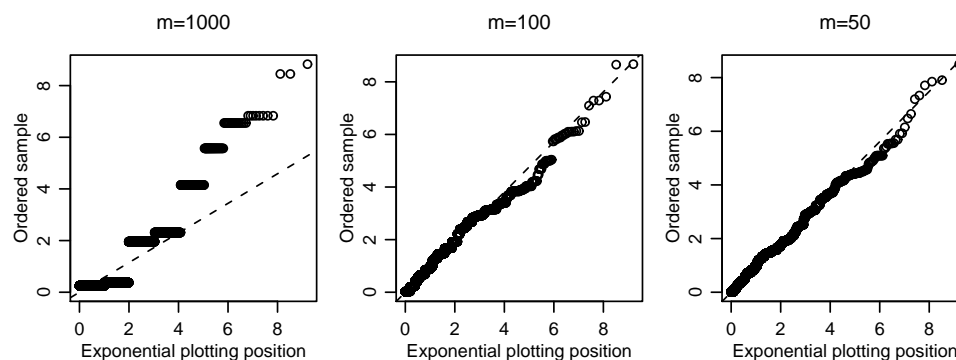


stat.epfl.ch

Autumn 2022 – slide 143

Resamples of size n ?

- ☐ Exponential sample of size $n = 1000$
- ☐ Distribution of $n \min(Y_1, \dots, Y_n)$ is $\exp(1)$
- ☐ Resampling distribution $m \min(Y_1^*, \dots, Y_m^*)$ using resamples of size $m = 1000, 100, 50$
- ☐ To avoid discreteness must choose $m \ll n$, but how?



stat.epfl.ch

Autumn 2022 – slide 144

Variants of \hat{G} ?

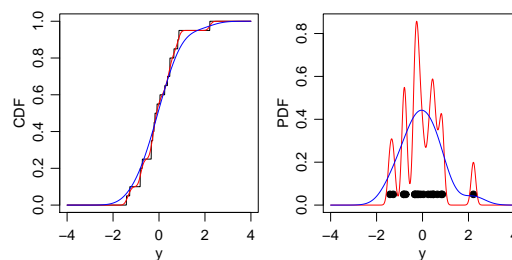
- Can be useful to simulate from a smoothed EDF, given by

$$Y^* = y_{j^*} + h\varepsilon^*, \quad \varepsilon^* \sim \mathcal{N}(0, 1) \perp\!\!\!\perp j^* \sim U\{1, \dots, n\},$$

equivalent to simulating from a kernel density estimate. Below, with $h = 0.1$ (red) and $h = 0.5$ (blue).

- Since $\text{var}^*(Y^*) = \hat{\sigma}^2 + h^2$, may prefer a shrunk smoothed estimate, given by

$$Y^* = \bar{y} + \frac{(y_{j^*} - \bar{y}) + h\varepsilon^*}{(1 + h^2/\hat{\sigma}^2)^{1/2}}.$$



When does the bootstrap work?

- ‘Work’ might mean the bootstrap gives
 - **reliable** answers when used in practice, or
 - **mathematically correct** answers under ‘suitable’ regularity conditions.
- For the second of these, suppose we seek to estimate properties of a standardized quantity $Q = q(Y_1, \dots, Y_n; G)$, maybe $Q = n^{1/2}(\bar{Y} - \theta)$. Let $n \rightarrow \infty$ to get limiting results for the distribution function

$$H_{G,n}(q) = P_G \{Q(Y_1, \dots, Y_n; G) \leq q\},$$

where subscript G indicates that Y_1, \dots, Y_n is a random sample from G .

- Bootstrap estimate of this is

$$H_{\hat{G},n}(q) = P_{\hat{G}} \{Q(Y_1^*, \dots, Y_n^*; \hat{G}) \leq q\}$$

where $Q(Y_1^*, \dots, Y_n^*; \hat{G}) = n^{1/2}(\bar{Y}^* - \bar{y})$.

- We need conditions under which $H_{\hat{G},n} \xrightarrow{D} H_{G,n}$ as $n \rightarrow \infty$.

Regularity conditions

- The true distribution G is surrounded by a neighbourhood \mathcal{N} in a suitable space of distributions, and as $n \rightarrow \infty$, \hat{G} eventually falls into \mathcal{N} with probability one. Also:
 1. for any $F \in \mathcal{N}$, $H_{F,n}$ converges weakly to a limit $H_{F,\infty}$;
 2. this convergence must be uniform on \mathcal{N} ; and
 3. the function mapping F to $H_{F,\infty}$ must be continuous.
- Weak convergence of $H_{F,n}$ to $H_{F,\infty}$ means that for all integrable $b(\cdot)$,

$$\int b(u) dH_{F,n}(u) \rightarrow \int b(u) dH_{F,\infty}(u), \quad n \rightarrow \infty.$$

- Under these conditions the bootstrap is **consistent**: for any q and $\varepsilon > 0$,

$$P\{|H_{\hat{G},n}(q) - H_{G,\infty}(q)| > \varepsilon\} \rightarrow 0, \quad n \rightarrow \infty.$$

- The first condition ensures that there is a limit for $H_{G,n}$ to converge to.
- As n increases, \hat{G} changes, so the second and third conditions are needed to ensure that $H_{\hat{G},n}$ approaches $H_{G,\infty}$ along every possible sequence of \hat{G} s.
- If any one of these conditions fails, the bootstrap can fail. For the minimum (for example) the convergence is not uniform on suitable neighbourhoods of G .

stat.epfl.ch

Autumn 2022 – slide 147

Summary

- **Estimator is algorithm:**
 - applied to original data y_1, \dots, y_n gives original $\hat{\theta}$;
 - applied to simulated data y_1^*, \dots, y_n^* gives $\hat{\theta}^*$;
 - $\hat{\theta}$ can be of (almost) any complexity; but
 - for more sophisticated ideas to work, $\hat{\theta}$ must often be smooth function of data.
- **Sample is used to estimate G :**
 - $\hat{G} \approx G$ — heroic assumption
- **Simulation replaces theoretical calculation:**
 - removes need for mathematical skill;
 - does not remove need for thought; and in particular,
 - check code **very** carefully — garbage in, garbage out!
- **Two sources of error:**
 - statistical ($\hat{G} \neq G$) — reduce by thought; and
 - simulation ($R \neq \infty$) — reduce by taking R large (enough).

stat.epfl.ch

Autumn 2022 – slide 148

Bootstrap confidence Intervals: Desiderata

- A $(1 - \alpha)$ **upper confidence limit** for a scalar parameter θ based on data Y is a random variable $\theta_\alpha = \theta_\alpha(Y)$ for which

$$P(\theta \leq \theta_\alpha) = \alpha, \quad 0 < \alpha < 1, \theta \in \Theta. \quad (7)$$

- We may seek invariance to monotone transformations $\psi = \psi(\theta)$, that is

$$P\{\psi(\theta) \leq \psi_\alpha\} = \alpha, \quad 0 < \alpha < 1, \theta \in \Theta.$$

- In practice exact intervals are rarely available, and we seek intervals such that (7) is satisfied as closely as possible. If $Y \equiv Y_1, \dots, Y_n$, then we typically have

$$P(\theta \leq \theta_\alpha) = \alpha + \mathcal{O}(n^{-1/2}), \quad 0 < \alpha < 1, \theta \in \Theta,$$

and the corresponding two-sided interval satisfies

$$P(\theta_\alpha < \theta \leq \theta_{1-\alpha}) = (1 - 2\alpha) + \mathcal{O}(n^{-1}), \quad 0 < \alpha < 1/2, \theta \in \Theta.$$

Normal confidence intervals

- If $\hat{\theta} \sim \mathcal{N}(\theta + \beta, \nu)$ with known bias $\beta = \beta(G)$ and variance $\nu = \nu(G)$, then a $(1 - 2\alpha)$ confidence interval is based on the equation

$$P\left(z_\alpha < \frac{\hat{\theta} - \theta - \beta}{\nu^{1/2}} \leq z_{1-\alpha}\right) = 1 - 2\alpha,$$

and has limits $\hat{\theta} - \beta \pm z_\alpha \nu^{1/2}$, where $\Phi(z_\alpha) = \alpha$.

- We replace β, ν by the bootstrap estimates

$$\beta(G) \doteq \beta(\hat{G}) \doteq b = \overline{\hat{\theta}^*} - \hat{\theta},$$

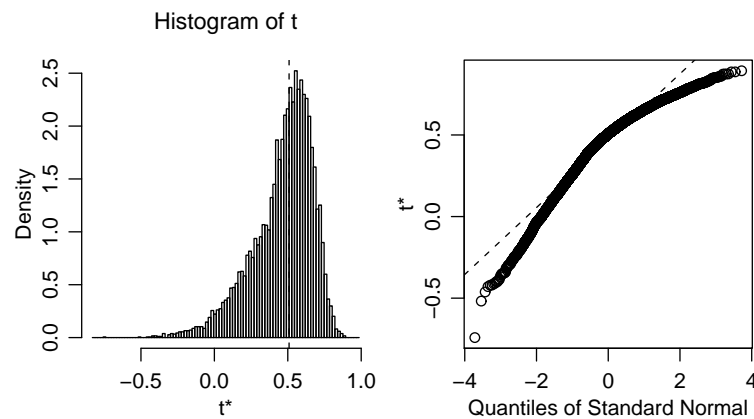
$$\nu(G) \doteq \nu(\hat{G}) \doteq v = (R - 1)^{-1} \sum_r (\hat{\theta}_r^* - \overline{\hat{\theta}^*})^2,$$

to get the $(1 - 2\alpha)$ interval with limits $\hat{\theta} - b \pm z_\alpha v^{1/2}$.

- For the handedness data we have $R = 10,000$, $b = -0.046$, $v = 0.205^2$, $\alpha = 0.025$, $z_\alpha = -1.96$, so 95% CI is (0.147, 0.963)
- We can use the $\hat{\theta}_1^*, \dots, \hat{\theta}_R^*$ to check the quality of the normal approximation, and perhaps to suggest transformations.

Handedness data

Summaries of the $\hat{\theta}^*$. Left: histogram, with vertical line showing $\hat{\theta}$. Right: normal Q-Q plot of $\hat{\theta}^*$.

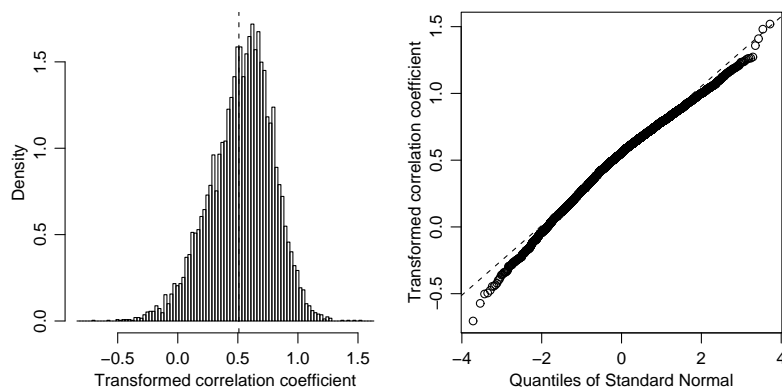


stat.epfl.ch

Autumn 2022 – slide 151

Handedness data: Transformed scale?

Plots for $\hat{\psi}^* = \frac{1}{2} \log\{(1 + \hat{\theta}^*)/(1 - \hat{\theta}^*)\}$:



stat.epfl.ch

Autumn 2022 – slide 152

Normal confidence intervals

- Correlation coefficient: try Fisher's z transformation:

$$\hat{\psi}^* = \psi(\hat{\theta}^*) = \frac{1}{2} \log\{(1 + \hat{\theta}^*)/(1 - \hat{\theta}^*)\}$$

with bias and variance estimates

$$b_\psi = R^{-1} \sum_{r=1}^R \hat{\psi}_r^* - \hat{\psi}, \quad v_\psi = \frac{1}{R-1} \sum_{r=1}^R (\hat{\psi}_r^* - \hat{\psi})^2,$$

- Then the $(1 - 2\alpha)$ confidence interval for θ is

$$\psi^{-1} \left\{ \hat{\psi} - b_\psi - z_{1-\alpha} v_\psi^{1/2} \right\}, \quad \psi^{-1} \left\{ \hat{\psi} - b_\psi - z_\alpha v_\psi^{1/2} \right\}$$

- For handedness data, get (0.074, 0.804) ... but how do we choose a transformation in general?

stat.epfl.ch

Autumn 2022 – slide 153

Pivots

- Assume properties of $\hat{\theta}_1^*, \dots, \hat{\theta}_R^*$ mimic effect of sampling from original model (plug-in principle) — false in general, but more nearly true for pivots.
- **Pivot** is combination of data and parameter whose distribution is independent of underlying model, such as t statistic

$$Z = \frac{\bar{Y} - \mu}{(S^2/n)^{1/2}} \sim t_{n-1},$$

when $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$.

- Exact pivot generally unavailable in nonparametric case, but if we can estimate the variance of $\hat{\theta}^*$ using V , we use

$$Z = \frac{\hat{\theta} - \theta}{V^{1/2}}$$

- If the quantiles z_α of Z known, then

$$P(z_\alpha \leq Z \leq z_{1-\alpha}) = P\left(z_\alpha \leq \frac{\hat{\theta} - \theta}{V^{1/2}} \leq z_{1-\alpha}\right) = 1 - 2\alpha$$

(z_α no longer denotes a normal quantile!) gives $(1 - 2\alpha)$ CI $(\hat{\theta} - V^{1/2} z_{1-\alpha}, \hat{\theta} - V^{1/2} z_\alpha)$

stat.epfl.ch

Autumn 2022 – slide 154

Studentized statistic

- Bootstrap sample gives $(\hat{\theta}^*, V^*)$ and hence

$$Z^* = \frac{\hat{\theta}^* - \hat{\theta}}{V^{*1/2}}.$$

- We bootstrap to get R copies of $(\hat{\theta}, V)$, i.e.,

$$(\hat{\theta}_1^*, V_1^*), (\hat{\theta}_2^*, V_2^*), \dots, (\hat{\theta}_R^*, V_R^*),$$

and the corresponding

$$z_1^* = \frac{\hat{\theta}_1^* - \hat{\theta}}{V_1^{*1/2}}, \quad z_2^* = \frac{\hat{\theta}_2^* - \hat{\theta}}{V_2^{*1/2}}, \quad \dots, \quad z_R^* = \frac{\hat{\theta}_R^* - \hat{\theta}}{V_R^{*1/2}},$$

then order these to estimate quantiles of Z , with z_p estimated by $z_{(p(R+1))}^*$.

- Get $(1 - 2\alpha)$ **Studentized bootstrap confidence interval**

$$\hat{\theta} - V^{1/2} z_{((1-\alpha)(R+1))}^*, \quad \hat{\theta} - V^{1/2} z_{(\alpha(R+1))}^*.$$

- This is not invariant to transformation and needs an estimated variance V_r^* for each $\hat{\theta}_r^*$.

stat.epfl.ch

Autumn 2022 – slide 155

Why Studentize?

- If we Studentize, then $Z \xrightarrow{D} N(0, 1)$ as $n \rightarrow \infty$, and we can use Edgeworth series to write

$$P_G(Z \leq z) = \Phi(z) + n^{-1/2} a(z) \phi(z) + O(n^{-1}),$$

where $a(\cdot)$ is an even quadratic polynomial.

- For example, if we use $\hat{\theta} = \bar{Y}$ and $V = n^{-1} S^2$ to compute Z for data with skewness γ , then $a(x) = \gamma(2x^2 + 1)/6$ and (next slide) $a'(x) = -\gamma(x^2 - 1)/6$.
- The corresponding expansion for Z^* is

$$P_{\hat{G}}(Z^* \leq z) = \Phi(z) + n^{-1/2} \hat{a}(z) \phi(z) + O_p(n^{-1}).$$

- Typically $\hat{a}(z) = a(z) + O_p(n^{-1/2})$, so

$$P_{\hat{G}}(Z^* \leq z) - P_G(Z \leq z) = O_p(n^{-1}),$$

so the order of error is n^{-1} .

stat.epfl.ch

Autumn 2022 – slide 156

Why Studentize? II

- Without Studentization, $Z = n^{1/2}(\hat{\theta} - \theta) \xrightarrow{D} N(0, \nu')$, and then

$$P_G(Z \leq z) = \Phi\left(\frac{z}{\nu'^{1/2}}\right) + n^{-1/2}a'\left(\frac{z}{\nu'^{1/2}}\right)\phi\left(\frac{z}{\nu'^{1/2}}\right) + O(n^{-1})$$

and

$$P_{\hat{G}}(Z^* \leq z) = \Phi\left(\frac{z}{\hat{\nu}'^{1/2}}\right) + n^{-1/2}\hat{a}'\left(\frac{z}{\hat{\nu}'^{1/2}}\right)\phi\left(\frac{z}{\hat{\nu}'^{1/2}}\right) + O_p(n^{-1}).$$

- Typically $\hat{\nu}' = \nu' + O_p(n^{-1/2})$, giving

$$P_{\hat{G}}(Z^* \leq z) - P_G(Z \leq z) = O_p(n^{-1/2}),$$

and the difference in the leading terms means that the overall error is of order $n^{-1/2}$.

- Thus Studentizing reduces error from $O_p(n^{-1/2})$ to $O_p(n^{-1})$: better than using large-sample asymptotics, for which error is usually $O_p(n^{-1/2})$.

Other confidence intervals

- Simpler approaches:

- **Basic bootstrap** interval: treat $\hat{\theta} - \theta$ as pivot, get

$$\hat{\theta} - (\hat{\theta}_{((R+1)(1-\alpha))}^* - \hat{\theta}), \quad \hat{\theta} - (\hat{\theta}_{((R+1)\alpha)}^* - \hat{\theta}).$$

- **Percentile interval**: use empirical quantiles of $\hat{\theta}_1^*, \dots, \hat{\theta}_R^*$:

$$\hat{\theta}_{((R+1)\alpha)}^*, \quad \hat{\theta}_{((R+1)(1-\alpha))}^*.$$

- The percentile interval is transformation-invariant, not the basic bootstrap interval.

- **Bias-corrected and accelerated (BC_a)** intervals replace percentile interval with $(\hat{\theta}_{((R+1)\alpha')}^*, \hat{\theta}_{((R+1)(1-\alpha'))}^*)$, where

$$\alpha' = \Phi\left\{w + \frac{w + z_\alpha}{1 - a(w + z_\alpha)}\right\}, \quad w = \Phi^{-1}\left\{\hat{G}^*(\hat{\theta})\right\}, \quad a = \frac{\frac{1}{6} \sum_{j=1}^n l_j^3}{\left(\sum_{j=1}^n l_j^2\right)^{3/2}},$$

with \hat{G}^* the EDF of the $\hat{\theta}_1^*, \dots, \hat{\theta}_R^*$, and l_1, \dots, l_n the empirical influence values (soon).

- If the **bias** $w = 0$, then $\hat{G}^*(\hat{\theta}) = \frac{1}{2}$, so $\hat{\theta}$ is at the median of the EDF of $\hat{\theta}^*$
- If the **acceleration** $a = 0$, then the effect of the data y_1, \dots, y_n on $\hat{\theta}$ is symmetric.

Comparisons

Table 2: Empirical error rates (%) for nonparametric bootstrap confidence limits in ratio estimation: rates for sample sizes $n_1 = n_2 = 10$ are given above those for sample sizes $n_1 = n_2 = 25$. $R = 999$ for all bootstrap methods. 10,000 data sets generated from Gamma distributions.

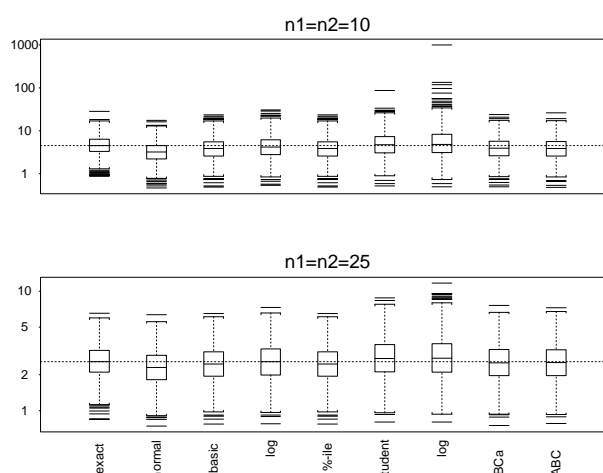
Method	Nominal error rate							
	Lower limit				Upper limit			
	1	2.5	5	10	10	5	2.5	1
Exact	1.0	2.8	5.5	10.5	9.8	4.8	2.6	1.0
	1.0	2.3	4.8	9.9	10.2	4.9	2.5	1.1
Normal approximation	0.1	0.5	1.7	6.3	20.6	15.7	12.5	9.6
	0.1	0.5	2.1	6.4	16.3	11.5	8.2	5.5
Basic bootstrap	0.0	0.0	0.2	1.8	24.4	21.0	18.6	16.4
	0.0	0.1	0.4	3.0	19.2	15.0	12.5	10.3
Basic bootstrap, log scale	2.6	4.9	8.1	12.9	13.1	7.5	4.8	2.5
	1.6	3.2	6.0	11.4	11.5	6.3	3.3	1.7
Studentized bootstrap	0.6	2.1	4.6	9.9	11.9	6.7	4.0	2.0
	0.8	2.3	4.6	9.9	10.9	5.9	3.0	1.4
Studentized bootstrap, log scale	1.1	2.8	5.6	10.7	11.6	6.3	3.5	1.7
	1.1	2.5	5.0	10.1	10.8	5.7	2.9	1.3
Bootstrap percentile	1.8	3.6	6.5	11.6	14.6	8.9	5.9	3.3
	1.2	2.6	5.1	10.1	12.6	7.1	4.2	2.1
BC_a	1.9	4.0	6.9	12.3	14.0	8.3	5.3	3.0
	1.4	3.0	5.6	10.9	11.8	6.8	3.8	1.9
ABC	1.9	4.2	7.4	12.7	14.6	8.7	5.5	3.1
	1.3	3.0	5.7	11.0	12.1	6.8	3.7	1.9

stat.epfl.ch

Autumn 2022 – slide 159

Confidence interval lengths

Lengths of 95% confidence intervals for the first 1000 simulated samples in the numerical experiment with Gamma data.



stat.epfl.ch

Autumn 2022 – slide 160

Discussion

- ☐ Bootstrap confidence intervals usually under-cover (i.e., are too short).
- ☐ Normal, basic, and studentized intervals depend on scale.
- ☐ Percentile interval often too short but is transformation-invariant.
- ☐ Studentized intervals give best coverage overall, but
 - they depend on scale, can be sensitive to V ;
 - their lengths can be very variable;
 - they are best when V is approximately constant.
- ☐ Improved percentile intervals have same asymptotic error as Studentized intervals, but often are shorter, so give lower coverage probabilities.
- ☐ Caution: Edgeworth theory OK for smooth statistics, but beware rough statistics: must check output.
- ☐ Typically need $R > 1000$ for reliable estimation of quantiles.

stat.epfl.ch

Autumn 2022 – slide 161

Nonparametric delta method

- ☐ The **delta method** (Theorem 40) gives variance formulae for functions of averages.
- ☐ More generally we use the **nonparametric delta method**, which is based on the linear functional expansion

$$t(F) \doteq t(G) + \int L_t(x; G) dF(x),$$

where L_t , the first derivative of $t(\cdot)$ at G , is defined by

$$L_t(y; G) = \lim_{\varepsilon \rightarrow 0} \frac{t\{(1 - \varepsilon)G + \varepsilon H_y\} - t(G)}{\varepsilon} = \left. \frac{\partial t\{(1 - \varepsilon)G + \varepsilon H_y\}}{\partial \varepsilon} \right|_{\varepsilon=0},$$

with $H_y(u) \equiv H(u - y)$ the Heaviside function jumping from 0 to 1 at $u = y$.

- ☐ The **influence function value** $L_t(y; G)$ for the statistical functional t for an observation at y when the background distribution is G , satisfies $E_G\{L_t(Y; G)\} = 0$.
- ☐ If \hat{G} is based on a random sample y_1, \dots, y_n , then the j th **empirical influence value** is

$$l_j = L_t(y_j; \hat{G}),$$

and $E_{\hat{G}}\{L_t(Y; \hat{G})\} = n^{-1} \sum_j l_j = 0$.

- ☐ The influence function also plays an important role in robust statistics.

stat.epfl.ch

Autumn 2022 – slide 162

Nonparametric delta method II

- If we replace F by the EDF \hat{G} for a random sample Y_1, \dots, Y_n , then

$$t(\hat{G}) \doteq t(G) + \int L_t(x; G) d\hat{G}(x) = t(G) + \frac{1}{n} \sum_{j=1}^n L_t(Y_j; G),$$

has variance

$$\text{var}\{t(\hat{G})\} \doteq \frac{1}{n^2} \sum_{j=1}^n L_t^2(Y_j; G) = V_L,$$

say, which we estimate based on a sample y_1, \dots, y_n by $v_L = n^{-2} \sum l_j^2$.

Example 51 Apply the nonparametric delta method to the average \bar{Y} .

Example 52 Apply the nonparametric delta method to a statistic defined by an estimating equation, and hence find the variance of the ratio \bar{V}/\bar{U} for data pairs $Y = (U, V)$.

stat.epfl.ch

Autumn 2022 – slide 163

Example 51

- The population mean and its empirical version are

$$\theta = t(G) = \int x dG(x), \quad \hat{\theta} = t(\hat{G}) = \int x d\hat{G}(x) = n^{-1} \sum_{j=1}^n Y_j = \bar{Y}.$$

- If H_y puts unit mass at y , its 'density' is a Dirac delta function $\delta_y(x)$, and

$$\begin{aligned} \theta\{(1-\varepsilon)G + \varepsilon H_y\} &= \int x d\{(1-\varepsilon)G + \varepsilon H_y\}(x) \\ &= (1-\varepsilon) \int x dG(x) + \varepsilon \int x dH_y(x) = (1-\varepsilon)\theta(G) + \varepsilon y \end{aligned}$$

and therefore

$$L(y; G) = \lim_{\varepsilon \rightarrow 0} \frac{\theta\{(1-\varepsilon)G + \varepsilon H_y\} - \theta(G)}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} \frac{(1-\varepsilon)\theta(G) + \varepsilon y - \theta(G)}{\varepsilon} = y - \theta(G),$$

- Hence the empirical influence values and variance estimate are

$$l_j = L(y_j; \hat{G}) = y_j - \bar{y}, \quad v_L = \frac{1}{n^2} \sum (y_j - \bar{y})^2 = \frac{n-1}{n} n^{-1} s^2.$$

stat.epfl.ch

Autumn 2022 – note 1 of slide 163

Example 52

- The scalar parameter $\theta = t(G)$ is determined implicitly through the estimating equation

$$\int a(x; \theta) dG(x) = \int a\{x; t(G)\} dG(x) = 0.$$

We replace G by $G_\varepsilon = (1 - \varepsilon)G + \varepsilon H_y$ and see that

$$\begin{aligned} 0 &= \int a\{x; t(G_\varepsilon)\} dG_\varepsilon(x) \\ &= (1 - \varepsilon) \int a\{x; t(G_\varepsilon)\} dG(x) + \varepsilon \int a\{x; t(G_\varepsilon)\} dH_y(x) \\ &= (1 - \varepsilon) \int a\{x; t(G_\varepsilon)\} dG(x) + \varepsilon a\{y; t(G_\varepsilon)\}, \end{aligned}$$

and differentiation using the chain rule gives

$$0 = a\{y; t(G_\varepsilon)\} - \int a\{x; t(G_\varepsilon)\} dG(x) + \varepsilon a_\theta\{y; t(G_\varepsilon)\} \frac{\partial t(G_\varepsilon)}{\partial \varepsilon} + (1 - \varepsilon) \int a_\theta\{x; t(G_\varepsilon)\} \frac{\partial t(G_\varepsilon)}{\partial \varepsilon} dG(x),$$

which reduces to

$$0 = a\{y; t(G)\} + \int a_\theta\{x; t(G)\} dG(x) \frac{\partial t(G)}{\partial \varepsilon} \Big|_{\varepsilon=0}$$

on setting $\varepsilon = 0$. Hence

$$L_t(y; G) = \frac{\partial t(G_\varepsilon)}{\partial \varepsilon} \Big|_{\varepsilon=0} = \frac{a(y; \theta)}{-\int a_\theta(x; \theta) dG(x)}, \quad \text{where } a_\theta(x; \theta) = \frac{\partial a(x; \theta)}{\partial \theta}.$$

- In the case of the ratio and with $y = (u, v)$, we take $a(y; \theta) = v - \theta u$, so

$$\theta = \theta(G) = \int v dG(u, v) / \int u dG(u, v), \quad \hat{\theta} = \bar{v} / \bar{u},$$

and $a_\theta = -u$, so $l_j = (x_j - \hat{\theta} u_j) / \bar{u}$, giving

$$v_L = \frac{1}{n^2} \sum \left(\frac{x_j - \hat{\theta} u_j}{\bar{u}} \right)^2.$$

Comments

- For statistics involving only averages (ratio, correlation coefficient, ...), the nonparametric delta method retrieves the delta method.
- For example, the correlation coefficient may be written as a function of $\overline{xu} = n^{-1} \sum x_j u_j$, etc.:

$$\hat{\theta} = \frac{\overline{xu} - \bar{x}\bar{u}}{\left\{(\overline{x^2} - \bar{x}^2)(\overline{u^2} - \bar{u}^2)\right\}^{1/2}},$$

from which empirical influence values l_j can be derived, giving $v_L = 0.029$ for the handedness data, to be compared with $v = 0.043$ obtained by bootstrapping.

- v_L typically underestimates $\text{var}(\hat{\theta})$!
- The l_j can also be obtained by numerical differentiation if $t(\hat{G})$ is coded appropriately, or approximated using a jackknife method.

Thomas Bayes (1702–1761)



Bayes (1763/4) *Essay towards solving a problem in the doctrine of chances*. Philosophical Transactions of the Royal Society of London.

stat.epfl.ch

Autumn 2022 – slide 166

Bayesian vs frequentist inference

Observed data y^o assumed to be realisation of $Y \sim f(y; \theta) \equiv f(y | \theta)$, where $\theta \in \Theta$.

□ **Frequentist viewpoint:**

- some ‘true value’ of θ generated the data;
- this ‘true value’ of θ is treated as an unknown constant;
- probability statements compare y^o with outcomes in a suitable reference set \mathcal{S} .

□ **Bayesian viewpoint:**

- degrees of belief should (and can) be expressed using probability distributions;
- knowledge about θ prior to seeing y^o is expressed as a **prior density** $\pi(\theta)$;
- Bayes’ theorem

$$\pi(\theta | y^o) = \frac{\pi(\theta)f(y^o | \theta)}{\int \pi(\theta)f(y^o | \theta) d\theta}$$

should be used to convert $\pi(\theta)$ into a **posterior density** $\pi(\theta | y^o)$;

- probability statements are based on $\pi(\theta | y^o)$ and thus are conditioned on all observed quantities.

- The benefit is that statistics reduces to calculations of probabilities, at the cost of expressing prior information in distributional terms.

stat.epfl.ch

Autumn 2022 – slide 167

Example

Example 53 (a) Find the posterior density for the success probability θ based on a series of independent Bernoulli trials y_1, \dots, y_n , when the prior density is the **Beta density**

$$\pi(\theta) = \frac{\theta^{a-1}(1-\theta)^{b-1}}{B(a,b)}, \quad 0 < \theta < 1, \quad a, b > 0,$$

where $B(a,b) = \Gamma(a)\Gamma(b)/\Gamma(a+b)$ is the **beta function**, and

$$\Gamma(a) = \int_0^\infty u^{a-1} e^{-u} du$$

is the gamma function.

(b) Show how the mean and variance of θ are updated.

(c) Find the posterior density for predicting the result Z of the next trial.

Example 53

- Suppose that conditional on θ , the data y_1, \dots, y_n are a random sample from the Bernoulli distribution, for which $P(Y_j = 1) = \theta$ and $P(Y_j = 0) = 1 - \theta$, where $0 < \theta < 1$. The likelihood is

$$L(\theta) = f(y | \theta) = \prod_{j=1}^n \theta^{y_j} (1-\theta)^{1-y_j} = \theta^s (1-\theta)^{n-s}, \quad 0 < \theta < 1,$$

where $s = \sum y_j$.

- A natural prior here is the beta density with parameters a and b ,

$$\pi(\theta) = \frac{1}{B(a,b)} \theta^{a-1} (1-\theta)^{b-1}, \quad 0 < \theta < 1, \quad a, b > 0, \quad (8)$$

where $B(a,b)$ is the beta function $\Gamma(a)\Gamma(b)/\Gamma(a+b)$.

- The posterior density of θ conditional on the data is

$$\begin{aligned} \pi(\theta | y) &= \frac{\theta^{s+a-1} (1-\theta)^{n-s+b-1} / B(a,b)}{\int_0^1 \theta^{s+a-1} (1-\theta)^{n-s+b-1} d\theta / B(a,b)} \\ &\propto \theta^{s+a-1} (1-\theta)^{n-s+b-1}, \quad 0 < \theta < 1. \end{aligned} \quad (9)$$

As (8) has unit integral for all positive a and b , the constant normalizing (??) must be $B(a+s, b+n-s)$. Therefore

$$\pi(\theta | y) = \frac{1}{B(a+s, b+n-s)} \theta^{s+a-1} (1-\theta)^{n-s+b-1}, \quad 0 < \theta < 1.$$

- Thus the posterior density of θ has the same form as the prior: acquiring data has the effect of updating (a,b) to $(a+s, b+n-s)$. As the mean of the $B(a,b)$ density is $a/(a+b)$, the posterior mean is $(s+a)/(n+a+b)$, and this is roughly s/n in large samples. Hence the prior density inserts information equivalent to having seen a sample of $a+b$ observations, of which a were successes. If we were very sure that $\theta \doteq 1/2$, for example, we might take $a=b$ very large, giving a prior density tightly concentrated around $\theta = 1/2$, whereas taking smaller values of a and b would increase the prior uncertainty.

100 spins of a 5Fr coin

```

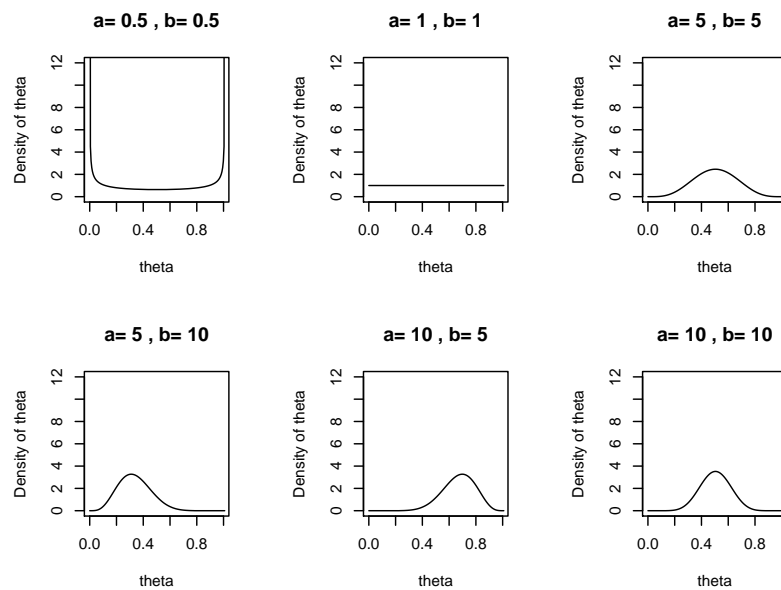
1 1 1 1 1 0 1 1 1 1 0 1 1 1 1 0 1 0 1 1
1 1 1 1 1 1 0 1 0 1 0 0 1 1 0 1 1 1 0 1
1 1 1 0 0 1 0 1 1 1 1 1 0 0 1 1 1 1 1 1
1 0 1 0 1 1 0 1 1 1 0 0 1 1 1 0 1 1 1 1
1 0 0 0 0 1 0 1 0 0 1 0 0 1 1 1 1 1 1 0

```

stat.epfl.ch

Autumn 2022 – slide 169

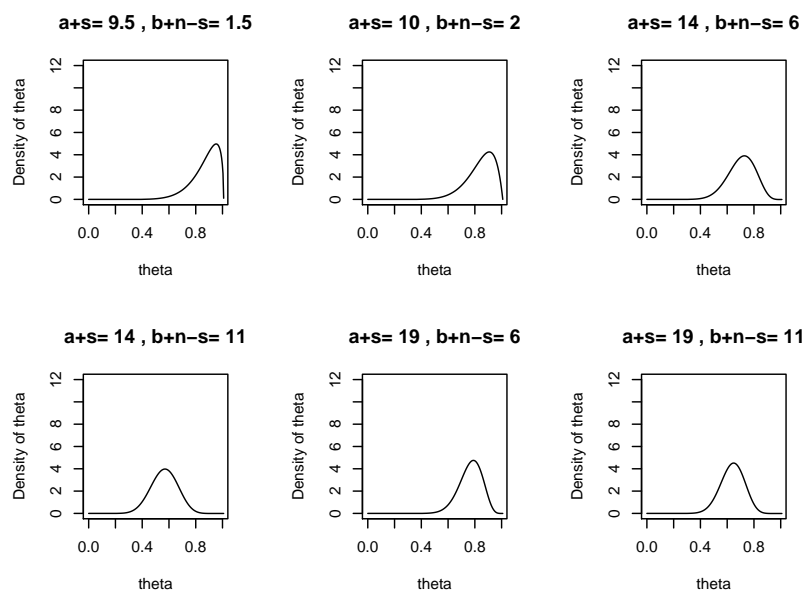
Beta prior densities



stat.epfl.ch

Autumn 2022 – slide 170

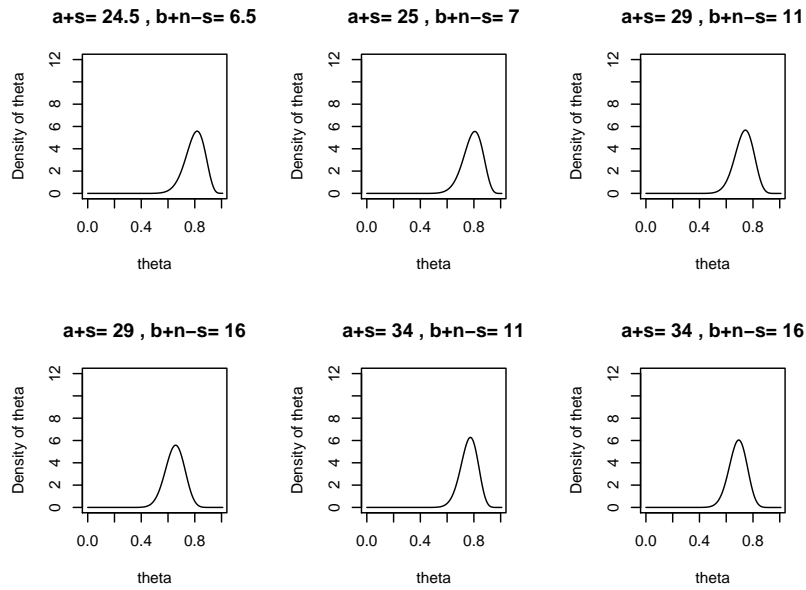
$n = 10, s = 9$



stat.epfl.ch

Autumn 2022 – slide 171

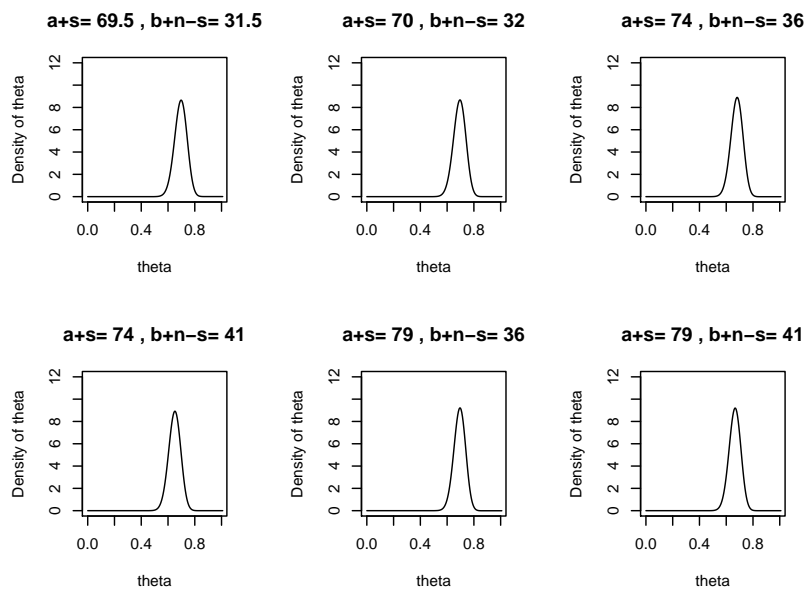
$n = 30, s = 24$



stat.epfl.ch

Autumn 2022 – slide 172

$n = 100, s = 69$



stat.epfl.ch

Autumn 2022 – slide 173

Link to likelihood

- In large samples the prior has less influence, because

$$\log \pi(\theta | y) = \log \pi(\theta) + \ell(\theta) - \log f(y),$$

where the terms on the right are successively $O(1)$, $O(n)$ and $O(n)$.

- Later we shall see that

$$f(y) \doteq \left(\frac{2\pi}{\hat{J}}\right)^{1/2} \pi(\hat{\theta}) e^{\ell(\hat{\theta})}$$

in terms of the MLE $\hat{\theta}$ and observed information \hat{J} , so

$$\pi(\theta | y) \doteq \frac{\pi(\theta)}{\pi(\hat{\theta})} \times \left(\frac{\hat{J}}{2\pi}\right)^{1/2} e^{\ell(\theta) - \ell(\hat{\theta})} \doteq \frac{\pi(\theta)}{\pi(\hat{\theta})} \times \left(\frac{\hat{J}}{2\pi}\right)^{1/2} e^{-\hat{J}(\hat{\theta} - \theta)^2/2},$$

giving the distributional approximation

$$\theta | y \sim \mathcal{N}(\hat{\theta}, \hat{J}^{-1}).$$

- Formal versions of this result, known as **Bernstein–von Mises theorems**, suggest that large-sample Bayesian and likelihood-based inferences will be similar.
- Hence we need to consider situations in which the prior may be appreciable relative to the information in the data, or in which standard likelihood approaches are unsuitable.

stat.epfl.ch

Autumn 2022 – slide 174

Conjugate priors

- Certain combinations of data model $f(y | \theta)$ and prior $\pi(\theta)$ give posterior densities of the same form as the prior.
- Example: $s \sim B(n, \theta)$ gives

$$\theta \sim \text{Beta}(a, b) \xrightarrow{s, n} \theta | y \sim \text{Beta}(a + s, b + n - s).$$

The beta density is the **conjugate prior** for binomial data.

- Conjugate priors greatly simplify computation and are widely used in modelling.
- Mixtures of conjugate priors are also conjugate (problem in Week 1).

Lemma 54 *An exponential family density*

$$f(y | \theta) = m(y) \exp[s(y)\varphi(\theta) - k\{\varphi(\theta)\}], \quad y \in \mathcal{Y}, \theta \in \Theta,$$

has conjugate prior

$$f(\theta; a, b) = h(a, b) \exp[a\varphi(\theta) - bk\{\varphi(\theta)\}], \quad \theta \in \Theta,$$

that depends on **hyperparameters** a, b .

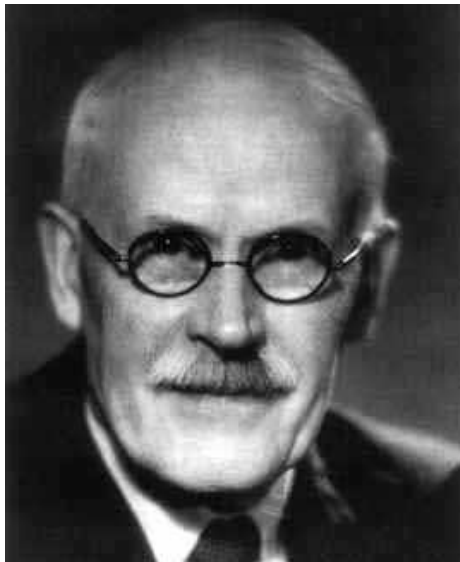
stat.epfl.ch

Autumn 2022 – slide 175

Two giants

Left: Harold Jeffreys (1891–1989), a geophysicist and astronomer who developed a (failed) theory of objective inference based on noninformative prior distributions.

Right: Ronald Alymer Fisher (1890–1962), a geneticist and statistician who developed a (failed) theory of objective inference based on the ‘fiducial’ distribution.



stat.epfl.ch

Autumn 2022 – slide 176

‘Ignorance’ about what?

Definition 55

- ☐ A **uniform prior** satisfies $\pi(\theta) \propto 1$ for $\theta \in \Theta$.
- ☐ An **improper prior** cannot be renormalised to have finite integral.
- ☐ The **Jeffreys prior** for a statistical model with Fisher information $\imath(\theta)$ is $\pi(\theta) \propto |\imath(\theta)|^{1/2}$.

Example 56 What does a uniform prior for $\theta \in (0, 1)$ imply for $\psi = \log\{\theta/(1 - \theta)\} \in \mathbb{R}$?

Lemma 57 The Jeffreys prior is invariant to smooth reparametrizations $\theta = \theta(\psi)$.

- ☐ Jeffreys priors were introduced to give ‘objective’ expressions of ignorance, and give uniform priors for location parameters, $1/\theta$ for scale parameters, etc.
- ☐ Jeffreys priors for the same θ based on different experiments might differ!
- ☐ Many other attempts to represent ‘ignorance’ have been made (e.g., by providing priors with minimal information), but none is seen as fully satisfactory.
- ☐ In practice ‘uninformative’ (i.e., flat but proper) priors are usually chosen and then sensitivity analyses performed.

stat.epfl.ch

Autumn 2022 – slide 177

Example 56

The probability of success in a Bernoulli trial lies in the interval $[0, 1]$, so if we are completely ignorant of its true value, the obvious prior to use is uniform on the unit interval: $\pi(\theta) = 1, 0 \leq \theta \leq 1$. But if we are completely ignorant of θ , we are also completely ignorant of $\psi = \log\{\theta/(1 - \theta)\}$, which takes values in the real line. The density implied for ψ by the uniform prior for θ is

$$\pi(\psi) = \pi\{\psi(\theta)\} \times \left| \frac{d\theta}{d\psi} \right| = \frac{e^\psi}{(1 + e^\psi)^2}, \quad -\infty < \psi < \infty :$$

the standard logistic density. Far from expressing ignorance about ψ , this density asserts that the prior probability of $|\psi| < 3$ is about 0.9.

stat.epfl.ch

Autumn 2022 – note 1 of slide 177

Lemma 57

- For a smooth reparametrization $\theta = \theta(\psi)$ in terms of ψ , the expected information for ψ is

$$\imath(\psi) = -E \left[\frac{d^2 \ell\{\theta(\psi)\}}{d\psi^2} \right] = -E \left\{ \frac{d^2 \ell(\theta)}{d\theta^2} \right\} \times \left| \frac{d\theta}{d\psi} \right|^2 = \imath(\theta) \times \left| \frac{d\theta}{d\psi} \right|^2.$$

Consequently $|\imath(\theta)|^{1/2} d\theta = |\imath(\psi)|^{1/2} d\psi$: the Jeffreys prior does behave consistently under reparametrization; furthermore such priors give widely-accepted solutions in some standard problems. When θ is vector, $|\imath(\theta)|$ is taken to be the determinant of $\imath(\theta)$.

- This prior was initially proposed with the aim of giving an ‘objective’ basis for inference, but after further paradoxes emerged its use was suggested for convenience, a matter of scientific convention rather than as a logically unassailable expression of ignorance about the parameter.

stat.epfl.ch

Autumn 2022 – note 2 of slide 177

High dimensions

Example 58 (Stein’s paradox) Let $Y_j \stackrel{\text{ind}}{\sim} \mathcal{N}(\mu_j, 1)$ for $j = 1, \dots, n$, and set $D = \sum Y_j^2$ and $\theta = \mu_1^2 + \dots + \mu_n^2$. Show that if the μ_j are independent a priori with flat priors, then

$$E(\theta | y) = D + n, \quad \text{but} \quad D \approx \theta + n + O_p(n^{1/2})$$

for any θ , which is absurd.

- Thus although flat priors may be sensible in low dimensions, they can lead to major problems in high dimensions.
- If we seek an uninformative prior for a scalar parameter ψ when nuisance parameters $\lambda_1, \dots, \lambda_p$ are orthogonal to ψ , we can set

$$\pi(\psi, \lambda) \propto \imath_{\psi\psi}^{1/2}(\psi, \lambda) \times g(\lambda),$$

where $\imath_{\psi\psi}(\psi, \lambda)$ is the (ψ, ψ) element of the Fisher information matrix and $g(\lambda)$ is an arbitrary function of the nuisance parameter.

stat.epfl.ch

Autumn 2022 – slide 178

Example 58

- If $y \mid \mu \sim \mathcal{N}(\mu, 1)$ and $\pi(\mu) \propto 1$, then symmetry of the normal density ϕ gives

$$\pi(\mu \mid y) = \frac{\phi(y - \mu)}{\int \phi(y - \mu) d\mu} = \frac{\phi(\mu - y)}{\int \phi(\mu - y) dy} = \phi(\mu - y),$$

so $\mu \mid y \sim \mathcal{N}(y, 1)$. If this is true independently for all the y_j , then

$$E(\theta \mid y) = \sum_{j=1}^n E(\mu_j^2 \mid y) = \sum_{j=1}^n \{E(\mu_j \mid y_j)^2 + \text{var}(\mu_j \mid y_j)\} = \sum_{j=1}^n (y_j^2 + 1) = D + n,$$

and its posterior variance is $\text{var}(\theta \mid y) = \sum_{j=1}^n \text{var}(\mu_j^2 \mid y_j) = 2n + 4D = O(n)$.

- On the other hand, for large n we have $D = \sum Y_j^2 \approx E(D) = \sum_{j=1}^n (\mu_j^2 + 1) = \theta + n$ and $\text{var}(D) = 2n + 4\theta = O(n)$.
- This implies that the posterior is placing probability in the wrong place asymptotically, i.e., around $D + n$ instead of around $D - n$. Hence the posterior probability that θ lies in any interval $D - n \pm a\sqrt{n}$ tends to zero.

stat.epfl.ch

Autumn 2022 – note 1 of slide 178

Matching priors

Definition 59 The **posterior α quantile** of a scalar parameter θ satisfies

$$P_{\theta \mid Y} \{\theta \leq \theta^\alpha(y) \mid y\} = \int_{-\infty}^{\theta^\alpha} \pi(\theta \mid y) d\theta = \alpha, \quad \alpha \in (0, 1)$$

- Consider random sample Y_1, \dots, Y_n with joint density $f(y \mid \theta)$, with prior $\pi(\theta)$ and $\theta \in \mathbb{R}^d$, and let $\hat{\theta}$ be the MLE and $\hat{\sigma}^2/n = \hat{j}^{-1}$ its asymptotic variance.
- Bayes and likelihood inferences will agree as $n \rightarrow \infty$, but is (approximate?) agreement achievable for small n ?
- If for every $\alpha \in (0, 1)$ and $\theta \in \Theta$ we could have

$$P_{Y \mid \theta} \{\theta^\alpha(Y) \geq \theta\} = \int I\{\theta^\alpha(y) \geq \theta\} f(y \mid \theta) dy = \alpha,$$

then Bayes and frequentist inference would agree perfectly, and we would have

- a Bayes/frequentist compromise;
- default priors for routine Bayesian use; and
- a basis for assessment of robustness of inference using other priors.

stat.epfl.ch

Autumn 2022 – slide 179

Edgeworth series

We use asymptotic approximations to compare the Bayesian and frequentist solutions.

Definition 60 Let X_1, \dots, X_n be a random sample of continuous variables with cumulant-generating function $K(u)$ and finite cumulants κ_r , let $\rho_r = \kappa_r / \kappa_2^{r/2}$ denote the r th standardized cumulant, and let $Z_n = (S_n - n\kappa_1) / (n\kappa_2)^{1/2}$ denote the standardized version of $S_n = X_1 + \dots + X_n$. Also let

$$\begin{aligned} H_1(z) &= z, \quad H_2(z) = z^2 - 1, \quad H_3(z) = z^3 - 3z, \quad H_4(z) = z^4 - 6z^2 + 3, \\ H_5(z) &= z^5 - 10z^3 + 15z, \quad H_6(z) = z^6 - 15z^4 + 45z^2 - 15 \end{aligned}$$

denote the Hermite polynomials. Then the **Edgeworth series** for the distribution of Z_n is

$$F_{Z_n}(z) = \Phi(z) - \phi(z) \left[\frac{\rho_3}{6n^{1/2}} H_2(z) + \frac{1}{n} \left\{ \frac{\rho_4}{24} H_3(z) + \frac{\rho_3^2}{72} H_5(z) \right\} + O(n^{-3/2}) \right],$$

and **Cornish–Fisher inversion** yields that the α quantile of $F_{Z_n}(z)$ equals

$$z_\alpha + \frac{\rho_3}{6n^{1/2}} H_2(z_\alpha) + \frac{1}{n} \left\{ \frac{\rho_4}{24} H_3(z_\alpha) + \frac{\rho_3^2}{36} (5z_\alpha - 2z_\alpha^3) \right\} + O(n^{-3/2}).$$

Matching: scalar θ

- We now compute Edgeworth series for the Bayesian quantity $n^{1/2}(\theta - \hat{\theta})/\hat{\sigma}$, conditional on y (so $\hat{\theta}(y), \hat{\sigma}(y)$ are constants), invert it to get the corresponding Cornish–Fisher series

$$\theta^\alpha(y) = \hat{\theta} - \frac{\hat{\sigma}}{n^{1/2}} z_\alpha + \frac{\hat{\sigma}}{n} \left\{ (z_\alpha^2 + 2) A_3(y) + A_1(y) \right\} + O(n^{-3/2}),$$

and then insert this expansion into

$$P_{Y|\theta} \{ \theta^\alpha(Y) \geq \theta \} = \int I \{ \theta^\alpha(y) \geq \theta \} f(y | \theta) dy.$$

- This gives

$$\alpha + \frac{\phi(z_\alpha)}{n^{1/2}} T_1(\pi, \theta) - \frac{z_\alpha \phi(z_\alpha)}{n} T_2(\pi, \theta) + O(n^{-3/2}),$$

where

$$T_1(\pi, \theta) = \frac{1}{\pi(\theta)} \frac{d}{d\theta} \left\{ \frac{\pi(\theta)}{\iota(\theta)^{1/2}} \right\}, \quad T_2 = 0 \iff \frac{d}{d\theta} \left\{ \frac{E_{Y|\theta}(\ell_\theta^3)}{\iota(\theta)^{3/2}} \right\} = 0.$$

- Choosing π to knock out T_1 will ensure matching to order n^{-1} , etc.

Discussion

- Clearly $T_1(\pi, \theta) \equiv 0$ if and only if

$$\pi(\theta) \propto \imath(\theta)^{1/2},$$

so the Jeffreys prior is matching to order n^{-1} for scalar θ .

- With the Jeffreys prior, T_2 vanishes if and only if

$$\frac{d}{d\theta} \left\{ \frac{E_{Y|\theta}(\ell_\theta^3)}{\imath(\theta)^{3/2}} \right\} = 0,$$

so even for scalar θ , higher-order matching is only possible in special cases.

- In the vector case, inferences for ψ match to order n^{-1} if ψ is orthogonal to the other parameters λ , and

$$\pi(\psi, \lambda) \propto \imath_{\psi\psi}^{1/2}(\psi, \lambda) \times g(\lambda).$$

- In general impossible to match for all parameters simultaneously—need separate (and incompatible) priors for each parameter.
- Higher order matching requires data-dependent priors.
- Kass and Wasserman (1996, JASA) give a general discussion of **reference priors**.

stat.epfl.ch

Autumn 2022 – slide 182

Inference

Once we have a prior, what about

- confidence sets?
- prediction?
- hypothesis tests?
- model comparison?
- model checking?

stat.epfl.ch

Autumn 2022 – slide 183

Confidence sets

- All measures of uncertainty are computed from the relevant posterior density.
- Posterior confidence bound for θ is quantile of $\pi(\theta | y)$:

$$P\{\theta \leq \theta^\alpha(y) | y\} = \int_{-\infty}^{\theta^\alpha(y)} \pi(\theta | y) d\theta = \alpha, \quad \alpha \in (0, 1),$$

giving $(1 - 2\alpha)$ posterior **credible set** $(\theta^\alpha(y), \theta^{1-\alpha}(y))$.

- In multiparameter case we use the marginal α quantile of ψ , $\psi^\alpha \equiv \psi^\alpha(y)$ as

$$P(\psi \leq \psi^\alpha | y) = \frac{\int_{-\infty}^{\psi^\alpha} \int f(y; \psi, \lambda) \pi(\psi, \lambda) d\lambda d\psi}{\iint f(y; \psi, \lambda) \pi(\psi, \lambda) d\lambda d\psi} \alpha, \quad \alpha \in (0, 1),$$

based on the marginal posterior density of ψ .

- A **highest posterior density (HPD) credible set** $\mathcal{C}_{1-\alpha}$ satisfies $P(\theta \in \mathcal{C}_{1-\alpha} | y) = 1 - \alpha$ and $\sup_{\theta \notin \mathcal{C}_{1-\alpha}} \pi(\theta | y) \leq \inf_{\theta \in \mathcal{C}_{1-\alpha}} \pi(\theta | y)$.
- Such intervals/sets are interpreted as probability statements about the parameter, with y fixed, contrary to frequentist confidence intervals.
- Likewise prediction intervals are based on the posterior predictive distribution $P(Z \leq z | y)$.

stat.epfl.ch

Autumn 2022 – slide 184

Example

Mortality rates r/m from cardiac surgery in 12 hospitals, showing the numbers of deaths r out of m operations.

<i>A</i>	0/47	<i>B</i>	18/148	<i>C</i>	8/119	<i>D</i>	46/810	<i>E</i>	8/211	<i>F</i>	13/196
<i>G</i>	9/148	<i>H</i>	31/215	<i>I</i>	14/207	<i>J</i>	8/97	<i>K</i>	29/256	<i>L</i>	24/360

Example 61 (Cardiac surgery data) A simple model for the data above treats the number of deaths r as binomial with mortality rate θ and denominator m . At hospital *A*, for example, $m = 47$ and $r = 0$, giving maximum likelihood estimate $\hat{\theta}_A = 0/47 = 0$, but it seems too optimistic to suppose that θ_A could be so small when the other rates are evidently larger. If we take a beta prior density with $a = b = 1$, the posterior density is beta with parameters $a + r = 1$ and $b + m - r = 48$. The 0.95 HPD credible interval is $(0, 6.05)\%$, while the equitailed credible interval uses the 0.025 and 0.975 quantiles of $\pi(\theta_A | y)$ and is $(0.05, 7.40)\%$.

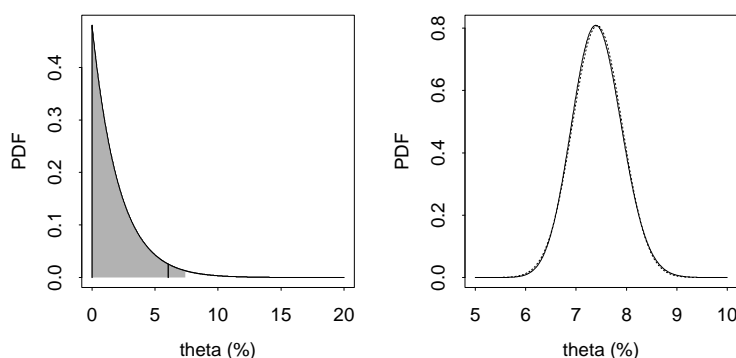
Show that the maximum posterior density estimator can be regarded as a penalized MLE.

stat.epfl.ch

Autumn 2022 – slide 185

Example

Cardiac surgery data. Left panel: posterior density for θ_A , showing boundaries of 0.95 highest posterior credible interval (vertical lines) and region between posterior 0.025 and 0.975 quantiles of $\pi(\theta_A | y)$ (shaded). Right panel: exact posterior beta density for overall mortality rate θ (solid) and normal approximation (dots).



stat.epfl.ch

Autumn 2022 – slide 186

Bayes factors

- ☐ Bayes factors compare competing models/hypotheses.
- ☐ Given prior probabilities $P(H_0)$ and $P(H_1)$ for two hypotheses, we compute

$$P(H_i | y) = \frac{P(y | H_i)P(H_i)}{P(y | H_0)P(H_0) + P(y | H_1)P(H_1)}, \quad i = 0, 1.$$

- ☐ Unlike in frequentist testing,
 - prior probabilities for the H_i must be specified, and
 - we compute the probability of each hypothesis given the data.
- ☐ To avoid specifying the prior probabilities we write

$$\frac{P(H_1 | y)}{P(H_0 | y)} = \frac{P(y | H_1)}{P(y | H_0)} \times \frac{P(H_1)}{P(H_0)} = B_{10} \times \frac{P(H_1)}{P(H_0)},$$

where B_{10} is the **Bayes factor**, and usually

$$P(y | H_i) = \int f(y | H_i, \theta_i) \pi(\theta_i | H_i) d\theta_i, \quad i = 0, 1.$$

stat.epfl.ch

Autumn 2022 – slide 187

Interpretation

- Often $2 \log B_{10}$ is used to summarise the evidence for H_1 , using a table like

B_{10}	$2 \log B_{10}$	Evidence for H_1
1–3	0–2	Hardly worth a mention
3–20	2–6	Positive
20–150	6–10	Strong
> 150	> 10	Very strong

- As $B_{10} = B_{01}^{-1}$, the evidence for H_0 is $2 \log B_{01} = -2 \log B_{10}$.
- Models $f(y | H, \theta)$ for n observations and $d \times 1$ parameter θ often compared using

$$\text{BIC} = -2\ell(\hat{\theta}) + d \log n,$$

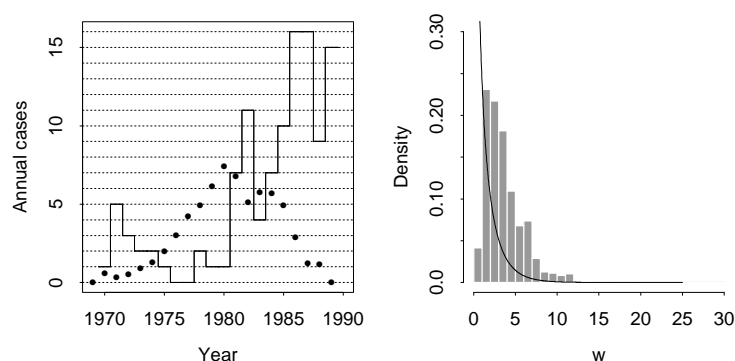
which can be derived by approximating the **model evidence** $P(y | H)$.

stat.epfl.ch

Autumn 2022 – slide 188

Example

Changepoint analysis for data on diarrhoea-associated haemolytic uraemic syndrome (HUS). Left: counts of cases of HUS treated in Birmingham, 1970–1989 (solid), and scaled likelihood ratio statistic $W_p(\tau)/10$ (blobs). Right: density of W , estimated from 10,000 simulations, and χ_1^2 density (solid).



stat.epfl.ch

Autumn 2022 – slide 189

Example

Example 62 (HUS data) The graph suggests a sharp rise in incidence around 1980. Suppose the annual counts y_1, \dots, y_n are realizations of independent Poisson variables with means λ_1 for $j = 1, \dots, \tau$ and λ_2 for $j = \tau + 1, \dots, n$. Here the changepoint τ can take values $1, \dots, n - 1$. Under H_0 , $\lambda_1 = \lambda_2 = \lambda$, that is, no change, and H_τ allows change after year τ . If we suppose that λ_1 and λ_2 have independent gamma prior densities with parameters γ and δ , then B_{10} can be computed for each τ .

There is very strong evidence for change in any year from 1976 to 1986, with most evidence for a change after 1980.

	1970	1971	1972	1973	1974	1975	1976	1977	1978	1979
y	1	5	3	2	2	1	0	0	2	1
$2 \log B_{\tau 0}, \gamma = \delta = 1$	4.9	-0.5	0.6	3.9	7.5	13	24	35	41	51
$2 \log B_{\tau 0}, \gamma = \delta = 0.01$	-1.3	-5.9	-4.5	-1.0	3.0	9.7	20	32	39	51
$2 \log B_{\tau 0}, \gamma = \delta = 0.0001$	-10	-15	-14	-10	-6.1	0.6	11	23	30	42

	1980	1981	1982	1983	1984	1985	1986	1987	1988	1989
y	1	7	11	4	7	10	16	16	9	15
$2 \log B_{\tau 0}, \gamma = \delta = 1$	63	55	38	42	40	31	11	-2.9	-5.3	0
$2 \log B_{\tau 0}, \gamma = \delta = 0.01$	64	57	40	47	46	38	18	1.8	1.2	0
$2 \log B_{\tau 0}, \gamma = \delta = 0.0001$	55	48	31	38	37	29	8.8	-7.1	-7.7	0

Nested models

- ☐ Often $\theta = (\psi, \lambda)$ and we want to compare $H_0 : \psi = \psi_0$ against $H_1 : \psi \neq \psi_0$.
- ☐ A prior density on θ will give

$$P(H_0) = \iint_{\{(\psi, \lambda) : \psi = \psi_0\}} \pi(\psi, \lambda) d\lambda d\psi = 0,$$

so the posterior odds in favour of H_1 are infinite for any dataset.

- ☐ To avoid we use prior densities weighted according to prior belief in H_0 and H_1 , giving overall prior

$$\pi(\psi, \lambda) = \delta(\psi - \psi_0) \pi(\psi_0, \lambda | H_0) P(H_0) + \pi(\psi, \lambda | H_1) P(H_1),$$

where

$$\int \pi(\psi_0, \lambda | H_0) d\lambda = \int \pi(\psi, \lambda | H_1) d\psi d\lambda = 1.$$

- ☐ Hence Bayes factors are more sensitive to the prior than are posterior densities.
- ☐ Improper priors cannot be used, as B_{10} depends on the ratio of the two arbitrary constants of proportionality in the priors.

Jeffreys–Lindley paradox

- Test $H_0 : \mu = 0$ against $H_1 : \mu \neq 0$ when $y_1, \dots, y_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$.
- Frequentist computes P-value $p_{\text{obs}} = \Phi(-n^{1/2}|\bar{y}|/\sigma)$.
- Bayesian writes $\pi_0 = P(H_0)$, supposes that under H_1 , $\mu \sim \mathcal{N}(0, \tau^2)$ and computes

$$B_{01} = \left(1 + n \frac{\tau^2}{\sigma^2}\right)^{1/2} \exp \left\{ -\frac{n\bar{y}^2}{2\sigma^2(1 + n^{-1}\sigma^2/\tau^2)} \right\}$$

- If $n\bar{y}^2/\sigma^2 = z_{\alpha/2}^2$, then $p_{\text{obs}} = \alpha$, but B_{01} gives increasingly strong evidence in favour of H_0 ; see the table, in which $\alpha = 0.01$:

n	1	10	100	1000	10^4	10^6	10^8
B_{01}	0.269	0.163	0.376	1.15	3.63	36.2	362

- The problem is that as $n \rightarrow \infty$, $\pi(\mu | H_1)$ is increasingly dispersed compared to $|\bar{y} - 0|$.
- To resolve this, note that we use tests when there is doubt about the hypotheses, i.e., sensible alternatives are $O(n^{-1/2})$ from the null, and if we take this account by setting $\tau^2 = \delta\sigma^2/n$, then the paradox dissipates, because (for example) with $\delta = 10$ and $\alpha = 0.05, 0.01, 0.001$, and 0.0001 , $B_{10} = 1.73, 6.2, 41.4$, and 293 , in broad agreement.

stat.epfl.ch

Autumn 2022 – slide 192

Model criticism

- Use marginal density $f(y)$ to check the model (and degree of agreement between $\pi(\theta)$ and $f(y | \theta)$). Simplest if

$$f(y) = f(y | s)f(a) \int f(t | a, \theta)\pi(\theta) d\theta,$$

where s is sufficient and a ancillary.

- Often leads to (Bayesian variants of) standard diagnostics (e.g., residuals, ...).
- Another measure of plausibility based on possible new dataset $Y_+ \sim f$ is

$$P\{f(Y_+) \leq f(y^o)\},$$

and yet another is based on **predictive diagnostics**, comparing a discrepancy measure $D_+ = d(Y_+, \theta)$ with its predictive distribution, i.e.,

$$P\{d(Y_+, \theta) \geq d(y, \theta) | y\},$$

where the averaging is over both Y_+ and the posterior distribution of θ .

- We choose $d(Y_+, \theta)$ to measure some key aspects of the data and model.

stat.epfl.ch

Autumn 2022 – slide 193

Prediction and model averaging

- Predict unobserved Z based on observed $Y = y$ from a single model by computing $f(z | y)$, but if there are several models, then

$$f(z | y) = \sum_{i=1}^k f(z | y, M_i) P(M_i | y),$$

which averages the posterior distributions of z under the different models, weighted according to their posterior probabilities

$$P(M_i | y) = \frac{f(y | M_i) P(M_i)}{\sum_{l=1}^k f(y | M_l) P(M_l)},$$

where

$$f(y | M_i) = \int f(y | \theta_i, M_i) \pi(\theta_i | M_i) d\theta_i,$$

$$f(z | M_i, y) = \frac{\int f(z | y, \theta_i, M_i) f(y | \theta_i, M_i) \pi(\theta_i | M_i) d\theta_i}{f(y | M_i)}.$$

- If we have all possible models, the main problem is computational ...

stat.epfl.ch

Autumn 2022 – slide 194

Example

Bayesian prediction using model averaging for the cement data. For each of the 16 possible subsets of covariates, the table shows the log Bayes factor in favour of that subset compared to the model with no covariates and gives the posterior probability of each model. The values of the posterior mean and scale parameters a and b are also shown for the six most plausible models; $(y_+ - a)/b$ has a posterior t density. For comparison, the residual sums of squares are also given.

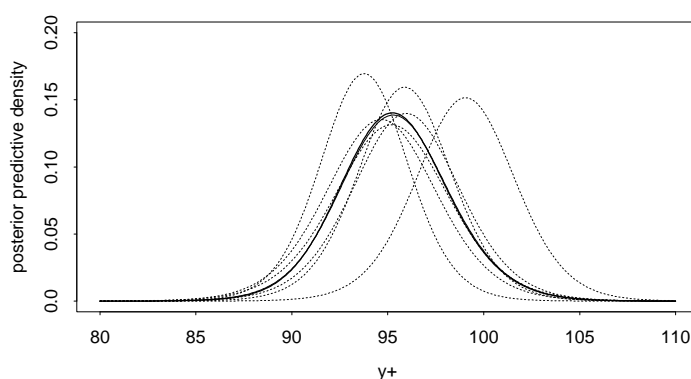
Model	RSS	$2 \log B_{10}$	$P(M y)$	a	b
----	2715.8	0.0	0.0000		
1---	1265.7	7.1	0.0000		
-2--	906.3	12.2	0.0000		
--3-	1939.4	0.6	0.0000		
---4	883.9	12.6	0.0000		
12--	57.9	45.7	0.2027	93.77	2.31
1-3-	1227.1	4.0	0.0000		
1--4	74.8	42.8	0.0480	99.05	2.58
-23-	415.4	19.3	0.0000		
-2-4	868.9	11.0	0.0000		
--34	175.7	31.3	0.0002		
123-	48.11	43.6	0.0716	95.96	2.80
12-4	47.97	47.2	0.4344	95.88	2.45
1-34	50.84	44.2	0.0986	94.66	2.89
-234	73.81	33.2	0.0004		
1234	47.86	45.0	0.1441	95.20	2.97

stat.epfl.ch

Autumn 2022 – slide 195

Example

Posterior predictive densities for cement data. Predictive densities for y_+ based on individual models are given as dotted curves, and the heavy curve is the averaged prediction from all 16 models.



stat.epfl.ch

Autumn 2022 – slide 196

Arguments for/against Bayes

☐ For:

- provides unified approach to inference—all unknowns, data, parameters, predictands are treated on the same footing;
- simple recipe — just apply Bayes' theorem and compute ...
- gives results similar to likelihood inferences (in large samples);
- argument based on axioms of 'rational behaviour' under uncertainty leads to 'coherent' (i.e., internally consistent) Bayes inference;

☐ Against:

- is it always (ever?) appropriate to treat data (whose model is checkable) on the same basis as the prior?
- Different priors may give different answers. Which is to be believed by a third party?
- How do we agree on a prior?
- External validity (in the frequency sense) with respect to reality is more important than internal consistency (one can be consistently wrong!)

☐ In any case, modelling can be flexible and general, provided computation is possible ...

stat.epfl.ch

Autumn 2022 – slide 197

Motivation

- We often want to approximate integrals such as those in the marginal posterior density

$$\pi(\psi | y) = \frac{\int f(y; \psi, \lambda) \pi(\psi, \lambda) d\lambda}{\iint f(y; \psi, \lambda) \pi(\psi, \lambda) d\lambda d\psi}$$

or the corresponding marginal posterior distribution function

$$P(\psi \leq \psi^0 | y) = \frac{\int_{-\infty}^{\psi^0} \int f(y; \psi, \lambda) \pi(\psi, \lambda) d\lambda d\psi}{\iint f(y; \psi, \lambda) \pi(\psi, \lambda) d\lambda d\psi}.$$

- Different approaches exist:
 - **deterministic** approximations include
 - ▷ quadrature rules — only work in low dimensions, not much used;
 - ▷ variational Bayes — provides numerical bounds on some integrals;
 - ▷ Laplace approximation — accurate analytical method with wide applications;
 - **Monte Carlo** approximations include
 - ▷ importance sampling — uses independent samples, can be unstable;
 - ▷ Markov chain Monte Carlo — widespread use in applications.

stat.epfl.ch

Autumn 2022 – slide 199

Laplace's method

Lemma 63 Let $h(u)$ be a smooth convex function defined for $u \in \mathbb{R}$, with a minimum at $u = \tilde{u}$, where $h'(\tilde{u}) = 0$ and $h''(\tilde{u}) > 0$, and let

$$I_n = \int_{-\infty}^{\infty} e^{-nh(u)} du.$$

Then

$$I_n = \left(\frac{2\pi}{nh_2} \right)^{1/2} e^{-nh(\tilde{u})} \times \left\{ 1 + n^{-1} \left(\frac{5h_3^2}{24h_2^3} - \frac{h_4}{8h_2^2} \right) + O(n^{-2}) \right\},$$

where $h_2 = h''(\tilde{u})$, etc. The leading term \tilde{I}_n is known as the **Laplace approximation** to I_n .

Comments:

- the error is relative, so the approximation is often very accurate far into the tails;
- \tilde{I}_n involves only h and its second derivative at \tilde{u} , so can be computed numerically;
- the series is asymptotic, so the partial sums may not converge, and including more than the leading term may give no improvements;
- most of the normal probability lies within ± 3 SD of the mean, so the limits of the integral don't matter (much) provided they lie outside the interval $\tilde{u} \pm 3(nh_2)^{-1/2}$;
- the exponent is written $-nh(u)$ only for formal justification of the approximation; in practice we set $n = 1$.

stat.epfl.ch

Autumn 2022 – slide 200

Note to Lemma 63

Close to \tilde{u} a Taylor series expansion gives $h(u) \doteq h(\tilde{u}) + \frac{1}{2}h_2(u - \tilde{u})^2$, so

$$\begin{aligned} I_n &\doteq e^{-nh(\tilde{u})} \int_{-\infty}^{\infty} e^{-nh_2(u-\tilde{u})^2/2} du \\ &= e^{-nh(\tilde{u})} \int_{-\infty}^{\infty} e^{-z^2/2} \frac{du}{dz} dz \\ &= \left(\frac{2\pi}{nh_2} \right)^{1/2} e^{-nh(\tilde{u})}, \end{aligned}$$

where the first and second equalities use the substitution $z = (nh_2)^{1/2}(u - \tilde{u})$ and the fact that the normal density has unit integral. A more detailed accounting gives the required result.

stat.epfl.ch

Autumn 2022 – note 1 of slide 200

Laplace's method: General case

Lemma 64 Let $h(u)$ be a smooth convex function defined for $u \in \mathbb{R}^d$, with a minimum at $u = \tilde{u}$, where $dh(\tilde{u})/du = 0$ and the hessian matrix

$$h_2 \equiv \frac{d^2 h(\tilde{u})}{du du^T}$$

is positive definite, and let

$$I_n = \int_{\mathbb{R}^d} e^{-nh(u)} du.$$

Then

$$I_n = \tilde{I}_n \{1 + O(n^{-1})\} = \left(\frac{2\pi}{n} \right)^{p/2} |h_2|^{-1/2} e^{-nh(\tilde{u})} \{1 + O(n^{-1})\}.$$

Example 65 Use Laplace approximation to derive the Bayesian information criterion.

stat.epfl.ch

Autumn 2022 – slide 201

Note to Example 65

□ Laplace approximation to $\log f(y)$ gives

$$\log \pi(\tilde{\theta}) + \log f(y | \tilde{\theta}) + \frac{p}{2} \log(2\pi/n) - \frac{1}{2} \log |\tilde{j}| + O(n^{-1}),$$

where $\tilde{\theta}$ maximises $\log \pi(\theta) + \log f(y | \theta)$ and $\tilde{j} = -n^{-1}$ times the hessian matrix of this function, evaluated at $\tilde{\theta}$.

□ Now $p \log(2\pi) - \log |\tilde{j}|$ is of order 1 as $n \rightarrow \infty$, and so is $\log \pi(\tilde{\theta})$, and $\tilde{\theta} = \hat{\theta} + O(n^{-1})$, so

$$-2 \log f(y) \doteq -2 \log f(y | \hat{\theta}) + p \log n + O(1) \approx \text{BIC}.$$

stat.epfl.ch

Autumn 2022 – note 1 of slide 201

Integral approximation

Lemma 66 Let

$$J_n(u_0) = \left(\frac{n}{2\pi}\right)^{1/2} \int_{-\infty}^{u_0} a(u) e^{-ng(u)} \{1 + O(n^{-1})\} du,$$

where $g(u)$ is a smooth convex function defined for $u \in \mathbb{R}$, and in addition to possessing the properties of h in Lemma 1, g satisfies $g(\tilde{u}) = 0$. Also let $a(u) > 0$. Then

$$J_n(u_0) = \Phi(n^{1/2}r_0^*) + O(n^{-1}),$$

where

$$r_0^* = r_0 + (r_0 n)^{-1} \log\left(\frac{v_0}{r_0}\right), \quad r_0 = \text{sign}(u_0 - \tilde{u})\{2g(u_0)\}^{1/2}, \quad v_0 = \frac{g'(u_0)}{a(u_0)}.$$

Example 67 Use the methods above to approximate the posterior conditional distribution

$$P(\theta \leq \theta_0 \mid y)$$

of a scalar parameter θ based on a random sample y_1, \dots, y_n from a regular model, and outline how posterior confidence intervals for θ are obtained.

stat.epfl.ch

Autumn 2022 – slide 202

Note to Lemma 66

- The first step is to change the variable of integration from u to $r(u) = \text{sign}(u - \tilde{u})\{2g(u)\}^{1/2}$; that is, $r^2/2 = g(u)$. Then $g'(u) = dg(u)/du$ and $r(u)$ have the same sign, and $r dr/du = g'(u)$, so

$$\begin{aligned} J_n(u_0) &= \left(\frac{n}{2\pi}\right)^{1/2} \int_{-\infty}^{r_0} a(u) \frac{r}{g'(u)} e^{-nr^2/2} \{1 + O(n^{-1})\} dr \\ &= \left(\frac{n}{2\pi}\right)^{1/2} \int_{-\infty}^{r_0} e^{-nr^2/2 + \log b(r)} \{1 + O(n^{-1})\} dr, \end{aligned}$$

where $b(r) = a(u)r/g'(u) > 0$ is regarded as a function of r .

- We now change variable again, from r to $r^* = r - (rn)^{-1} \log b(r)$, so

$$-nr^{*2} = -nr^2 + 2 \log b(r) - n^{-1}r^{-2} \{\log b(r)\}^2.$$

The Jacobian of the transformation and the third term in $-nr^{*2}$ contribute only to the error of $J_n(u_0)$, so

$$\begin{aligned} J_n(u_0) &= \left(\frac{n}{2\pi}\right)^{1/2} \int_{-\infty}^{r_0^*} e^{-nr^{*2}/2} \{1 + O(n^{-1})\} dr^* \\ &= \Phi(n^{1/2}r_0^*) + O(n^{-1}), \end{aligned} \tag{10}$$

where

$$r_0^* = r_0 + (r_0 n)^{-1} \log\left(\frac{v_0}{r_0}\right), \quad r_0 = \text{sign}(u_0 - \tilde{u})\{2g(u_0)\}^{1/2}, \quad v_0 = \frac{g'(u_0)}{a(u_0)}.$$

stat.epfl.ch

Autumn 2022 – note 1 of slide 202

Note to Example 67

- We write

$$P(\theta \leq \theta_0 | y) = \frac{\int_{-\infty}^{\theta_0} \pi(\theta) f(y | \theta) d\theta}{\int_{-\infty}^{\infty} \pi(\theta) f(y | \theta) d\theta}$$

and set $h(\theta) = -n^{-1}\{\ell(\theta) + \log \pi(\theta)\} = -\ell_m(\theta)/n$, say. This (scaled) modified log likelihood is maximised at $\tilde{\theta}$, which is the maximum a posteriori estimate of θ , and $h''(\theta) = -n^{-1}\ell_m''(\theta) = n^{-1}j(\theta) - n^{-1}(\log \pi)''(\theta)$.

- Laplace approximation of the denominator integral gives

$$\sqrt{\frac{2\pi}{nh_2}} \exp\{-nh(\tilde{\theta})\} \{1 + O(n^{-1})\},$$

where $h_2 = h''(\tilde{\theta})$, and inserting this into the expression for the posterior probability gives

$$P(\theta \leq \theta_0 | y) = \sqrt{\frac{nh_2}{2\pi}} \int_{-\infty}^{\theta_0} e^{-n\{h(\theta) - h(\tilde{\theta})\}} \{1 + O(n^{-1})\} d\theta,$$

to which we can apply Lemma 66 with $g(\theta) = h(\theta) - h(\tilde{\theta}) \geq 0$; this equals zero when $\theta = \tilde{\theta}$, and $a(\theta) = (nh_2)^{1/2}$. We take $u = \theta$, $u^0 = \theta_0$,

$$v_0 = g'(\theta_0)/(nh_2)^{1/2} = -n^{-1}\ell_m'(\theta_0)/\{-\ell_m''(\tilde{\theta})\}^{1/2}, \quad r_0 = \text{sign}(\theta_0 - \tilde{\theta}) \left[2\{\ell_m(\tilde{\theta}) - \ell_m(\theta_0)\}/n \right]^{1/2},$$

and therefore

$$n^{1/2}r_0^* = n^{1/2}r_0 - \frac{1}{n^{1/2}r_0} \log \left\{ \frac{-\ell_m'(\theta_0)/\{-\ell_m''(\tilde{\theta})\}^{1/2}}{n^{1/2}r_0} \right\}.$$

Hence we can simply set $n = 1$ and compute $r_0 = \text{sign}(\theta_0 - \tilde{\theta}) \left[2\{\ell_m(\tilde{\theta}) - \ell_m(\theta_0)\} \right]^{1/2}$.

- Hence we can write

$$P(\theta \leq \theta_0 | y) = \Phi\{r_B^*(\theta_0)\} \{1 + O(n^{-1})\},$$

where $r_B^*(\theta_0)$ is given by the expressions above with $n = 1$. We obtain confidence intervals by solving for θ_0 the equations

$$\alpha, 1 - \alpha = \Phi\{r_B^*(\theta_0)\}, \quad \text{or equivalently} \quad z_\alpha, z_{1-\alpha} = r_B^*(\theta_0).$$

- The likelihood root (almost) corresponds to setting $\pi(\theta) \propto 1$, so that $\tilde{\theta} = \hat{\theta}$ and $nh_2 = \hat{j}$, and then we get

$$r_0 = -\text{sign}(\hat{\theta} - \theta_0) \left[2\{\ell(\hat{\theta}) - \ell(\theta_0)\} \right]^{1/2}, \quad v_0 = -\hat{j}^{-1/2}\ell'(\theta_0).$$

This makes sense, because

$$P(\theta \leq \theta_0 | y) \doteq \Phi\{r_B^*(\theta_0)\}$$

is increasing in θ_0 , but the corresponding expression for a frequentist interval is decreasing in θ_0 . So we expect that $r_B^*(\theta_0) \doteq -r^*(\theta_0)$.

Integral approximation: General case

Lemma 68 Let $u = (u_1, u_2)$, where u_1 is scalar and u_2 a $p \times 1$ vector, and consider

$$J_n(u_1^0) = (2\pi)^{-(p+1)/2} c \int_{-\infty}^{u_1^0} \int \exp \{-ng(u_1, u_2)\} du_2 du_1, \quad (11)$$

where c is constant, the inner integral being over \mathbb{R}^p . Here g is supposed to have its previous smoothness properties, to be maximized at $(\tilde{u}_1, \tilde{u}_2)$, and satisfies $g(\tilde{u}_1, \tilde{u}_2) = 0$. Then

$$J_n(u_1^0) = \Phi(n^{1/2}r_0^*) + O(n^{-1}),$$

where $r_0^* = r_0 + (r_0 n)^{-1} \log \left(\frac{v_0}{r_0} \right)$, with

$$r_0 = \text{sign}(u_1^0 - \tilde{u}_1) \{2g(u_1^0, \tilde{u}_{20})\}^{1/2}, \quad v_0 = c^{-1} \frac{\partial g(u_1^0, \tilde{u}_{20})}{\partial u_1} |g_{22}(u_1^0, \tilde{u}_{20})|^{1/2},$$

where \tilde{u}_{20} is the maximizing value of u_2 when $u_1 = u_1^0$.

Multivariate case

- The computations of Example 67 can be extended to the multiparameter case using Lemmas 64 and 68, and give

$$P(\psi \leq \psi_0 | y) = \Phi\{r_B^*(\psi_0)\} \{1 + O(n^{-1})\},$$

where $r_B^*(\psi_0) = r_B(\psi_0) + r_B(\psi_0)^{-1} \log \{v_B(\psi_0)/r_B(\psi_0)\}$, with

$$r_B(\psi_0) = \text{sign}(\psi_0 - \tilde{\psi}) \left[2 \left\{ \ell_m(\tilde{\psi}, \tilde{\lambda}) - \ell_m(\psi_0, \tilde{\lambda}_{\psi_0}) \right\} \right]^{1/2},$$

$$v_B(\psi_0) = -\frac{\partial \ell_m(\psi_0, \tilde{\lambda}_{\psi_0})}{\partial \psi} \left\{ \frac{\left| -\frac{\partial^2 \ell_m(\psi_0, \tilde{\lambda}_{\psi_0})}{\partial \lambda \partial \lambda^T} \right|}{\left| -\frac{\partial^2 \ell_m(\tilde{\psi}, \tilde{\lambda})}{\partial \theta \partial \theta^T} \right|} \right\}^{1/2};$$

here $\tilde{\lambda}_{\psi_0}$ is the maximum *a posteriori* estimate of λ when ψ is fixed at ψ_0 .

- Often we find the derivatives numerically.
- There is a close link to maximum likelihood estimation, because $\tilde{\theta} = \hat{\theta} + O(n^{-1})$, so the order of error is not increased by using the MLEs instead of the MAPs — though the numerical approximations are not so good.

Frequentist aside

- In frequentist inference **saddlepoint approximation** is used to write conditional densities for exponential families as

$$f(t_1 | t_2; \psi) \doteq \left\{ \frac{|J_{\lambda\lambda}(\hat{\theta}_\psi)|}{2\pi|J(\hat{\theta})|} \right\}^{1/2} \exp \left\{ \ell(\hat{\theta}_\psi) - \ell(\hat{\theta}) \right\},$$

leading to

$$P(T_1 \leq t_1 | T_2 = t_2; \psi) \doteq \Phi\{r^*(\psi)\},$$

where $r^*(\psi) = r(\psi) + v(\psi)^{-1} \log\{r(\psi)/v(\psi)\}$, with

$$r(\psi) = \text{sign}(\hat{\psi} - \psi)[2\{\ell(\hat{\theta}) - \ell(\hat{\theta}_\psi)\}]^{1/2}, \quad v(\psi) = (\hat{\psi} - \psi) \left\{ \frac{|J(\hat{\theta})|}{|J_{\lambda\lambda}(\hat{\theta}_\psi)|} \right\}^{1/2}.$$

- Saddlepoint approximation involves writing the exponential family density as an integral of its Laplace transform (or equivalently its cumulant-generating function), and then approximating the resulting integral.
- The details are somewhat more painful, but the idea is similar to the Bayesian case.
- The approach sketched on slides 117–119 extends this to arbitrary regular models, by approximating them by exponential families.

stat.epfl.ch

Autumn 2022 – slide 205

Comments

- For successful approximation we must be able to write the integrand as

$$\exp \{ \log f(y; \theta) + \log \pi(\theta) \},$$

where the exponent is $O(n)$ and the integrand has one dominant mode.

- If so the methods can work well in fairly high dimensions, partly because the errors in numerator and denominator can cancel.
- However Monte Carlo methods are more flexible and in more general use . . .

stat.epfl.ch

Autumn 2022 – slide 206

Importance sampling

- Seek to estimate

$$\mu = \int m(\theta, y, z) \pi(\theta | y) d\theta,$$

where taking, for example,

- $m(\theta, y, z) = I(\theta \leq a)$ will give $\mu = P(\theta \leq a | y)$,
 - $m(\theta, y, z) = f(z | y, \theta)$ will give $\mu = f(z | y)$.
- If we can sample $\theta_1, \dots, \theta_R \stackrel{\text{iid}}{\sim} h(\theta)$, where the support of h includes that of $\pi(\theta | y)$, then we have an importance sampling estimator

$$\hat{\mu} = R^{-1} \sum_{r=1}^R m(\theta_r, y, z) \frac{\pi(\theta_r | y)}{h(\theta_r)} = R^{-1} \sum_{r=1}^R m(\theta_r, y, z) w(\theta_r),$$

where $w(\theta) = \pi(\theta | y) / h(\theta)$ is an importance sampling weight.

- Advantage of $\hat{\mu}$ over MCMC output is that its variance is readily obtained.
- Disadvantage is that choice of h is usually difficult. and especially if $\dim(\theta)$ is large, so huge samples are needed because most of the simulated θ_r receive zero weight and so are wasted.

Markov chain Monte Carlo

- Want to learn about distribution π of random variable $U \in \mathcal{U}$:
 - in Bayesian statistics, U is all unknowns and π is their posterior distribution conditioned on observed data y ;
 - in frequentist statistics U may be functions of the data y , and we seek to condition on other functions, e.g., to perform a conditional test.
- Construct a Markov chain $\{U^t\}$ with state space \mathcal{U} and transition kernel P , whose limiting distribution is π , i.e.,

$$P(U^t \in \mathcal{A} | u^0) \rightarrow \pi(\mathcal{A}) \quad t \rightarrow \infty, \quad u^0 \in \mathcal{U}, \mathcal{A} \subset \mathcal{U}.$$

- We then use P to simulate a realisation u^0, u^1, \dots, u^R of the chain, and hence get estimates such as

$$E_{\pi}\{g(U) | y\} = \int g(u) \pi(u | y) du \approx \frac{1}{R} \sum_{r=1}^R g(u^r), \quad \pi(\mathcal{A} | y) \approx \frac{1}{R} \sum_{r=1}^R I(u^r \in \mathcal{A}).$$

- Must choose P and u^0 so that
 - the distribution of U^t converges quickly to π (so minimise simulation effort);
 - u^0, u^1, \dots, u^R are as independent as possible (so have efficient estimation).

Markov chains

Definition 69

(a) A sequence U^0, U^1, U^2, \dots of elements of a set \mathcal{U} is a **Markov chain** if the conditional distribution of U^{t+1} given U^1, \dots, U^t depends only on U^t :

$$P(U^{t+1} \in \mathcal{A} \mid U^1, \dots, U^t) = P(U^{t+1} \in \mathcal{A} \mid U^t), \quad \mathcal{A} \subset \mathcal{U}.$$

We call \mathcal{U} the **state space** of the Markov chain.

(b) A Markov chain has **stationary transition probabilities** if the conditional distribution of U^{t+1} given U^t does not depend on t .

(c) The distribution of U^0 is called the **initial distribution**, and the conditional distribution

$$P(u, \mathcal{A}) = P(U^{t+1} \in \mathcal{A} \mid U^t = u)$$

is called the **transition probability distribution** (or **transition kernel**); this does not depend on t if the chain has stationary transition probabilities, and then we denote it by P .

(d) The **stationary** or **invariant** or **equilibrium** distribution of a Markov chain with transition kernel P satisfies

$$\pi(\mathcal{A}) = \int P(u, \mathcal{A}) \pi(du), \quad \mathcal{A} \subset \mathcal{U}.$$

stat.epfl.ch

Autumn 2022 – slide 209

Ergodicity and convergence

For the distribution of U^t to converge to a stationary distribution, the chain must satisfy three important properties:

- ☐ **irreducibility** — \mathcal{U} does not split into separate parts when we run the chain on it, so the kernel P allows us to reach any point of \mathcal{U} starting from anywhere else;
- ☐ **aperiodicity** — precludes the possibility of the ‘limiting’ distribution depending on the iteration number, i.e., eliminates possibilities like $a_n = (-1)^n$, which equals 1 if n is even and otherwise is odd;
- ☐ **positive recurrence** — every state is visited infinitely often, if the chain is run forever. This enables estimation of properties of that state.
- ☐ An irreducible, aperiodic, positive recurrent chain is called **ergodic**.
- ☐ The **ergodic theorem** states that an ergodic Markov chain has a unique stationary distribution π ,

$$P(U^t \in \mathcal{A} \mid U_0 = u) \rightarrow \pi(\mathcal{A}), \quad t \rightarrow \infty, \quad u \in \mathcal{U}, \mathcal{A} \subset \mathcal{U},$$

and if g is a real-valued function with $\int |g(u)| \pi(du) < \infty$, then

$$\frac{1}{R} \sum_{t=1}^R g(U^t) \xrightarrow{\text{a.s.}} \int g(u) \pi(du), \quad R \rightarrow \infty.$$

stat.epfl.ch

Autumn 2022 – slide 210

Detailed balance

- Modulo technical details (skipped here), the implication is that if we can find a transition kernel P with invariant distribution π , then we can generate samples (almost) from π .
- Why 'almost'? Because we run the chain for a finite number of steps, so in general our samples are not exactly from π .
- We now describe some standard recipes for building MCMC algorithms.
- For simplicity of exposition we take \mathcal{U} to be countable, so $P \equiv P(u, v)$ for $u, v \in \mathcal{U}$.
- A sufficient condition for invariance is **detailed balance**:

$$\pi(u)P(u, v) = \pi(v)P(v, u), \quad u, v \in \mathcal{U}.$$

- This guarantees invariance because

$$\begin{aligned} \int P(u, \mathcal{A})\pi(\mathrm{d}u) &= \sum_{v \in \mathcal{A}} \sum_{u \in \mathcal{U}} \pi(u)P(u, v) \\ &= \sum_{v \in \mathcal{A}} \sum_{u \in \mathcal{U}} \pi(v)P(v, u) \\ &= \sum_{v \in \mathcal{A}} \pi(v) \sum_{u \in \mathcal{U}} P(v, u) = \pi(\mathcal{A}) \times 1 = \pi(\mathcal{A}). \end{aligned}$$

Metropolis–Hastings algorithm

- A very general algorithm to estimate a target density π , with many variants.
- Hastings (1970) generalised an idea of Metropolis et al. (1953):
 - given a current value u of the chain, construct a candidate new value (a 'proposal') v by drawing from an arbitrary density $q(v | u)$;
 - accept the proposal as the next state of the chain with probability

$$a(u, v) = \min \left\{ 1, \frac{\pi(v)q(u | v)}{\pi(u)q(v | u)} \right\}$$

and otherwise leave u unchanged.

- The target density π is needed only up to the constant of proportionality, and only at u and the proposal v , so in particular the normalising constant is not needed.
- An important special case, the **Gibbs sampler**, updates each component u_i of u by successively writing $u = (u_i, u_{-i})$ and then replacing u_i with $v_i \sim \pi(u_i | u_{-i})$, where $\pi(u_i | u_{-i})$ is called the **full conditional density**.

Example 70 (Toy) Construct a Metropolis–Hastings algorithm with $\mathcal{N}(0, 1)$ target density and proposal distribution $q(v | u) = \sigma^{-1}\phi\{(v - u)/\sigma\}$.

Note: Detailed balance for the M-H algorithm

- First we note that

$$P(u, v) = q(v | u)a(u, v) + r(u)I(u = v),$$

where

$$r(u) = 1 - \int q(v | u)a(u, v) dv.$$

The first and second terms of $P(u, v)$ are the probability density for a move from u to v being proposed and accepted, and the probability that a move away from u is rejected.

- The Metropolis–Hastings update step satisfies detailed balance because

$$\begin{aligned}\pi(u)P(u, v) &= \pi(u)q(v | u) \min \left\{ 1, \frac{\pi(v)q(u | v)}{\pi(u)q(v | u)} \right\} + \pi(u)r(u)I(u = v) \\ &= \pi(v)q(u | v) \min \left\{ \frac{\pi(u)q(v | u)}{\pi(v)q(u | v)}, 1 \right\} + \pi(v)r(v)I(v = u) \\ &= \pi(v)P(v, u).\end{aligned}$$

Hence the corresponding Markov chain is reversible with equilibrium distribution π , provided it is irreducible and aperiodic.

Note to Example 70

- We need to work out the acceptance ratio

$$\frac{\pi(v)q(u | v)}{\pi(u)q(v | u)}$$

where

$$\pi(u) \propto e^{-u^2/2}, \quad q(u | v) = (2\pi\sigma^2)^{-1/2} e^{-(u-v)^2/2\sigma^2},$$

and this is

$$\frac{e^{-v^2/2} \times (2\pi\sigma^2)^{-1/2} e^{-(u-v)^2/2\sigma^2}}{e^{-u^2/2} \times (2\pi\sigma^2)^{-1/2} e^{-(v-u)^2/2\sigma^2}} = \exp\left\{\frac{1}{2}(u^2 - v^2)\right\},$$

so the move $u \mapsto v$ is accepted with probability $\min[1, \exp\{\frac{1}{2}(u^2 - v^2)\}]$.

- If $v^2 \leq u^2$ the acceptance ratio is greater than unity and the move is always accepted, whereas if $v^2 > u^2$ the move may not be accepted, and if $v^2 \gg u^2$ the move is very unlikely to be accepted.

- Note that

- we did not need the normalising constant for π to run the algorithm;
- the acceptance ratio does not depend on σ ;
- the acceptance probability does depend on σ . With $W \sim U(0, 1)$, it is

$$\begin{aligned}P(u \mapsto V | u) &= P(W \leq \min[1, \exp\{\frac{1}{2}(u^2 - v^2)\}] | u) \\ &= P(|V| \leq u | u) + \int_{\{v: |v| > |u|\}} e^{(u^2 - v^2)/2} \frac{1}{(2\pi\sigma^2)^{1/2}} e^{-(v-u)^2/2\sigma^2} du,\end{aligned}$$

which clearly depends on u and on σ .

Toy MH example: Code

```
toy.MH <- function(R=5000, sig=1, u0=-10, seed)
{
  set.seed(seed)
  u <- rep(u0,R)
  for (r in 2:R)
  {
    v <- rnorm(1, u[r-1], sig)
    log.ratio <- dnorm(v, log=T) + dnorm(v, mean=u[r-1], sd=sig, log=T) -
      dnorm(u[r-1], log=T) - dnorm(u[r-1], mean=v, sd=sig, log=T)
    a <- min( 1, exp(log.ratio) )
    u[r] <- u[r-1]
    if (runif(1)<=a) u[r] <- v
  }
  u
}

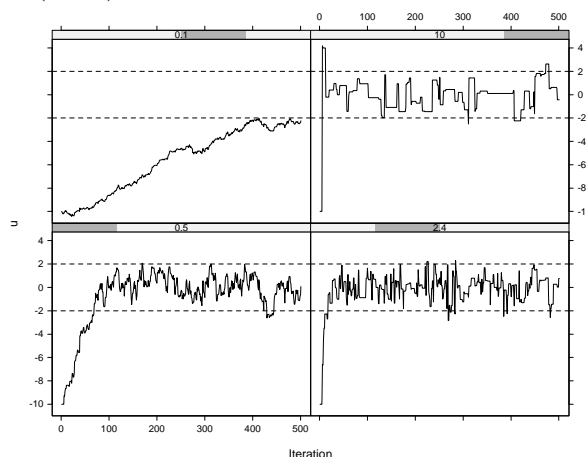
save.seed <- .Random.seed # use the same seed for each simulation
out1 <- toy.MH(sig=0.1, seed=save.seed)
out2 <- toy.MH(sig=0.5, seed=save.seed)
plot.ts(out1, ylim=c(-10,3), xlab="Iteration", ylab="u")
plot.ts(out2, ylim=c(-10,3), xlab="Iteration", ylab="u")
```

stat.epfl.ch

Autumn 2022 – slide 213

Toy MH example

Simulations from a Metropolis–Hastings algorithm with $\mathcal{N}(0, 1)$ target density, with $u^0 = -10$ and random walk proposal $v \sim \mathcal{N}(u, \sigma^2)$ with $\sigma = 0.1, 0.5, 2.4, 10$.



With $\sigma = 0.1, 0.5$, proposals often accepted but chain moves too slowly. With $\sigma = 10$ chain gets stuck for too long. Here $\sigma = 2.4$ seems best.

stat.epfl.ch

Autumn 2022 – slide 214

Proposal distributions

- In principle there is an (almost) completely free choice for the proposal distributions q_i , but just a few possibilities are typically used:
 - **Independence Metropolis–Hastings**, in which $q(v)$ is unrelated to u . Not much use in practice, but helpful for theoretical analysis.
 - **Random walk Metropolis**, in which $q(u, v) = q(v - u)$ and $q(\cdot)$ is a density symmetric about 0, giving

$$a(u, v) = \min \left\{ 1, \frac{\pi(v)}{\pi(u)} \right\}$$

because $q(u, v) = q(v, u)$. This amounts to setting $v = u + \varepsilon$, where $\varepsilon \sim q$.

- **Random walk Metropolis on the log scale**, applied when $u > 0$, in which random walk Metropolis is applied to $\log u$; then $q(v, u)/q(u, v) = v/u$ and so

$$a(u, v) = \min \left\{ 1, \frac{\pi(v)v}{\pi(u)u} \right\}.$$

Similar random walks can be applied to other transformations.

stat.epfl.ch

Autumn 2022 – slide 215

Toy Gibbs sampler

Example 71 Find the joint posterior density for the mean and standard deviation of a normal random sample of size n with prior distributions $\mu \sim \mathcal{N}(\xi, \kappa^{-1})$ and $\sigma^{-2} \sim \Gamma(\alpha, \beta)$.

stat.epfl.ch

Autumn 2022 – slide 216

Note to Example 71

The joint posterior is

$$\pi(\mu, \sigma^{-2} \mid y) \propto (\sigma^{-2})^{\alpha+n/2-1} \exp \left\{ -\frac{\beta}{\sigma^2} - \frac{\kappa(\mu - \xi)^2}{2} - \frac{\sum (y_j - \mu)^2}{2\sigma^2} \right\}$$

so the parameters are dependent *a posteriori* although they were independent *a priori*. The full conditional densities are

$$\begin{aligned} \mu \mid \sigma, y &\sim \mathcal{N} \left(\frac{\sum y_j + \sigma^2 \kappa \xi}{n + \kappa \sigma^2}, \frac{1}{n \sigma^{-2} + \kappa} \right), \\ \frac{1}{\sigma^2} \mid \mu, y &\sim \Gamma \left(\alpha + n/2, \beta + \sum (y_j - \mu)^2 / 2 \right), \end{aligned}$$

and the Gibbs sampler alternates updates of μ and of σ^{-2} using these two equations.

stat.epfl.ch

Autumn 2022 – note 1 of slide 216

Toy Gibbs example: Code

```
# Darwin's maize data in eighths of an inch
n <- 15
y <- c(49,-67,8,16,6,23,28,41,14,29,56,24,75,60,-48)

# Set (improper) prior parameters and number of iterations R
xi <- kappa <- alpha <- beta <- 0
R <- 10000

# Gibbs sampler with initial values mu=0, 1/sig^2=0.002
out <- matrix(NA,R,2)
out[1,] <- c(0, 0.002)
for (r in 2:R)
{
  new.mean <- (sum(y) + kappa*xi/out[r-1,2])/(n+kappa/out[r-1,2])
  new.var <- 1/(n*out[r-1,2] + kappa)
  out[r,1] <- rnorm(1, mean=new.mean, sd=sqrt(new.var))
  out[r,2] <- rgamma(1, rate=beta+sum((y-out[r,1])^2)/2, shape=alpha+n/2)
}

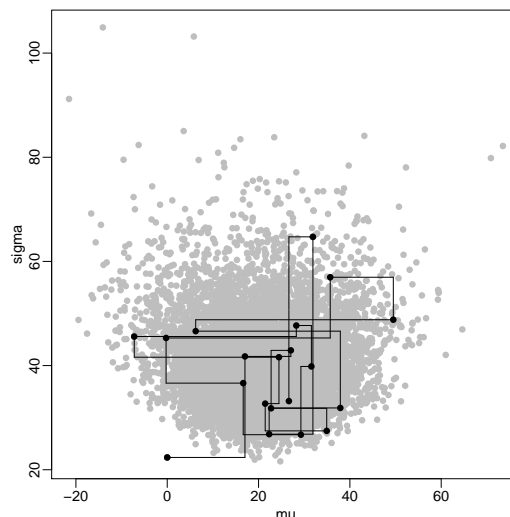
# posterior values of sigma
out[,2] <- sqrt(1/out[,2])
```

stat.epfl.ch

Autumn 2022 – slide 217

Toy Gibbs example

10,000 iterations of Gibbs sampler for (μ, σ) , with initial value $\mu = 0$; the $(\mu$ update, σ update) steps are shown for the first 20 iterations:

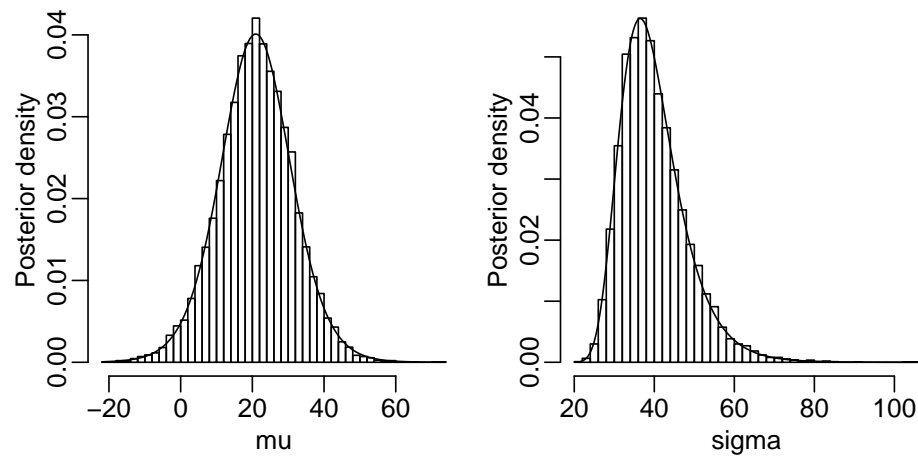


stat.epfl.ch

Autumn 2022 – slide 218

Toy Gibbs example

Marginal histograms and density estimates for μ and σ , based on 10,000 simulations:

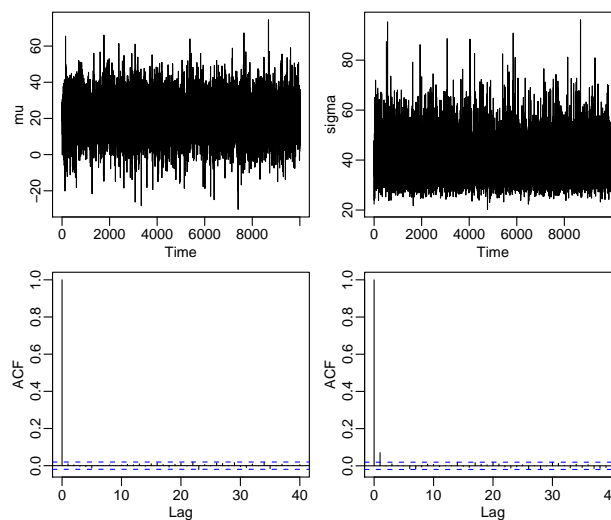


stat.epfl.ch

Autumn 2022 – slide 219

Toy Gibbs example

Time series of μ and σ , and correlograms. The series appear to be stationary with (very) low autocorrelation:



stat.epfl.ch

Autumn 2022 – slide 220

Toy Gibbs example: Estimation

- Any function of (μ, σ) can be estimated using the successive pairs $(\mu, \sigma)_1, \dots, (\mu, \sigma)_R$.
- For example, to compute $\psi = P(Y_+ \leq -50 \mid y)$ we can either add simulation of a new observation Y_+ to each iteration, giving $(\mu, \sigma, Y_+)_1, \dots, (\mu, \sigma, Y_+)_R$, or we can use conditioning to obtain the estimators

$$\hat{\psi}_1 = \frac{1}{R} \sum_{r=1}^R I(Y_{+,r} \leq -50), \quad \hat{\psi}_2 = \frac{1}{R} \sum_{r=1}^R \Phi\left(\frac{-50 - \mu_r}{\sigma_r}\right).$$

The maximum likelihood estimator of ψ is $\hat{\psi} = \Phi\{(-50 - \hat{\mu})/\hat{\sigma}\} = 0.030$, where $\hat{\mu}, \hat{\sigma}$ are the MLEs, but the Bayes estimator is $\hat{\psi}_2 = 0.045$, which is larger because it allows for the variability of the parameters (though it depends on the prior).

- Similar arguments apply to estimation of marginal densities, using either by a kernel density estimator or an unbiased estimator based on the full conditionals. For example,

$$\pi(\mu \mid y) \doteq R^{-1} \sum_{r=1}^R \frac{1}{h} K\left(\frac{\mu - \mu_r}{h}\right), \quad \pi(\mu \mid y) \doteq R^{-1} \sum_{r=1}^R \pi(\mu \mid \sigma_r, y),$$

where K is a kernel function with bandwidth h .

stat.epfl.ch

Autumn 2022 – slide 221

Discussion

- Update several variables at once by taking vector u_i — most useful if the components of u_i are conditionally independent given u_{-i} , which allows parallel updates.
- All the methods use the full conditionals $\pi(u_i \mid u_{-i})$: the Gibbs sampler draws from them, but the M-H algorithm only evaluates them at u and v .
- To ensure that the overall chain is ergodic we must make the chain reversible as a whole. In some cases this is obvious, but if not, and the kernels for updating different variables are P_1, \dots, P_m , then we might take

$$P = P_1 \cdots P_{m-1} P_m P_{m-1} \cdots P_1, \quad \text{or} \quad P = m^{-1} \sum_{i=1}^m P_i, \quad \text{or} \quad \frac{1}{m!} \sum_{\xi} \prod_{i=1}^m P_{\xi(i)},$$

where ξ is a random permutation of $\{1, \dots, m\}$.

- **Convergence diagnostics** are needed to check 'stationarity' of output — simple time series plots are helpful, but more sophisticated methods exist, often based on comparing multiple chains.
- There is a huge (and still growing) literature on all aspects of these methods.

stat.epfl.ch

Autumn 2022 – slide 222

Motivation

- Many types of data have layers of variation, which must be modelled:
 - disease incidence varies between regions of a country, and within regions it may vary due to effects of poverty, pollution, ...
 - success of surgical interventions may depend on patients (age/state of health) within surgeons (different experience/skill) within hospitals (different environments/skill of nursing staff)
- We think of populations from which patients, doctors, hospitals, ... are drawn, and this suggests modelling them using layers of randomness.
- This sort of construction is very common in modelling complex data, in both classical and Bayesian frameworks.
- Some theoretical justification is provided by the notion of exchangeability.
- First, an aside on graphical representations of complex models.

stat.epfl.ch

Autumn 2022 – slide 224

Graphical models

- Complex dependencies are often represented using graphs:
 - helps understanding;
 - transforming the type of graph can simplify certain computations.
- Graph language for generic variables Y_1, \dots, Y_n :
 - Y_j is represented by a **node** of the graph, so the node set is $\mathcal{J} = \{1, \dots, n\}$;
 - we define a **neighbourhood system** $\mathcal{N} = \{\mathcal{N}_j, j \in \mathcal{J}\}$ such that
 - ▷ the **neighbours** of j are the elements of $\mathcal{N}_j \subset \mathcal{J}$, where for each j the **neighbourhood** \mathcal{N}_j satisfies

$$(i) \quad j \notin \mathcal{N}_j, \quad (ii) \quad i \in \mathcal{N}_j \Leftrightarrow j \in \mathcal{N}_i,$$
 and let $\tilde{\mathcal{N}}_j = \mathcal{N}_j \cup \{j\}$;
 - the set of nodes and the neighbourhood structure $(\mathcal{J}, \mathcal{N})$ define the graph.

stat.epfl.ch

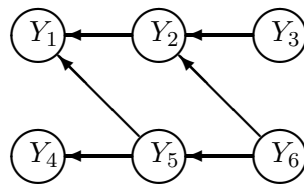
Autumn 2022 – slide 225

Directed acyclic graphs

Definition 72 A **directed acyclic graph (DAG)** is a graphical model that represents a hierarchical dependence structure:

- conditional dependence of Y_1 on Y_2 is represented by an arrow from the **parent** node Y_2 to the **child** node Y_1 ;
- Y_1 is a **descendent** of Y_3 if there is a chain of arrows from Y_3 to Y_1 ;
- it is **directed** because each arc is an arrow; and
- it is **acyclic** because it is impossible to start from a node, traverse a path by following arrows, and end up at the starting-point.

The decomposition $f(y) = f(y_1 | y_2, y_5)f(y_2 | y_3, y_6)f(y_3)f(y_4 | y_5)f(y_5 | y_6)f(y_6)$ gives:



Conditional independence graph

- Construct a **conditional independence graph** from a DAG, by adding edges between any parents that share a child and dropping the arrowheads.
- The conditional distribution of Y_j given Y_{-j} depends only on the variables $Y_{\mathcal{N}_j}$ directly linked to Y_j in the conditional independence graph:

$$f(y_j | y_{-j}) = f(y_j | y_{\mathcal{N}_j}).$$

- Why? For any DAG,

$$f(y) = \prod_{j \in \mathcal{J}} f(y_j | \text{parents of } y_j)$$

so

$$\begin{aligned}
 f(y_j | y_{-j}) &= \frac{f(y)}{\int f(y) dy_j} = \frac{\prod_{i \in \mathcal{J}} f(y_i | \text{parents of } y_i)}{\int \prod_{i \in \mathcal{J}} f(y_i | \text{parents of } y_i) dy_j} \\
 &\propto f(y_j | \text{parents of } y_j) \prod_{\{i: y_i \text{ is child of } y_j\}} f(y_i | \text{parents of } y_i) \\
 &\propto f(y_j | y_{\mathcal{N}_j}),
 \end{aligned}$$

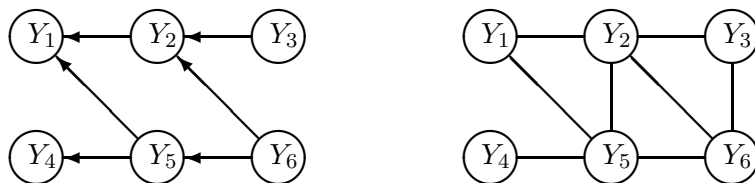
because terms without y_j cancel from the ratio.

Simplifying full conditional distributions

- The DAG and conditional independence graph help in constructing an MCMC sampler:
 - we use the model definition to write down the DAG;
 - we convert the DAG into a conditional independence graph;
 - we read the required conditional dependencies off from the conditional independence graph.
- The conditional independence graph (right) implies that

$$f(y_1 | y_{-1}) = f(y_1 | y_2, y_5), \quad f(y_2 | y_{-2}) = f(y_2 | y_1, y_3, y_5, y_6),$$

$$f(y_3 | y_{-3}) = f(y_3 | y_2, y_6), \quad f(y_4 | y_{-4}) = f(y_4 | y_5), \quad \dots$$



Exchangeability

Back to hierarchical models:

Definition 73 The random variables U_1, \dots, U_n are called **finitely exchangeable** if their density has the property

$$f(u_1, \dots, u_n) = f(u_{\xi(1)}, \dots, u_{\xi(n)})$$

for any permutation ξ of the set $\{1, \dots, n\}$. An infinite sequence U_1, U_2, \dots , is called **infinitely exchangeable** if every finite subset of it is finitely exchangeable.

- Variables are exchangeable if there is no reason to distinguish them.
- f is completely symmetric in its arguments and in probabilistic terms the U_1, \dots, U_n are indistinguishable (but not necessarily independent).

De Finetti's theorem

Theorem 74 (de Finetti) If U_1, U_2, \dots , is an infinitely exchangeable sequence of binary variables, taking values $u_j = 0, 1$, then for any n there is a distribution G such that

$$f(u_1, \dots, u_n) = \int_0^1 \prod_{j=1}^n \theta^{u_j} (1 - \theta)^{1-u_j} G(d\theta) \quad (12)$$

where

$$G(\theta) = \lim_{m \rightarrow \infty} P \{ m^{-1}(U_1 + \dots + U_m) \leq \theta \}, \quad \theta = \lim_{m \rightarrow \infty} m^{-1}(U_1 + \dots + U_m).$$

- Hence any set of exchangeable binary variables U_1, \dots, U_n that may be embedded within an infinite sequence may be modelled as if they were independent Bernoulli variables, conditional on their success probability θ , this having distribution G and being interpretable as the long-run proportion of successes.
- Similar theorems apply to other types of variables (continuous, ...).
- Thus a judgement that certain quantities are exchangeable implies that they may be represented as a random sample conditional on some θ — equivalent to using a prior distribution for θ .

stat.epfl.ch

Autumn 2022 – slide 230

Normal example

The following example illustrates properties of all hierarchical models.

Example 75 Suppose that v_1, \dots, v_n , σ^2 , μ_0 and τ^2 are known and

$$\begin{aligned} \mu &\sim \mathcal{N}(\mu_0, \tau^2), \\ \theta_1, \dots, \theta_n \mid \mu &\stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2), \\ y_j \mid \theta_j &\stackrel{\text{ind}}{\sim} \mathcal{N}(\theta_j, v_j), \quad j = 1, \dots, n. \end{aligned}$$

Here the **hyperparameters** μ_0 and τ^2 control the uncertainty at the top level of the hierarchy. Give the DAG and conditional independence graph for this model, and show that

$$\begin{aligned} E(\mu \mid y) &= \frac{\mu_0/\tau^2 + \sum y_j/(\sigma^2 + v_j)}{1/\tau^2 + \sum 1/(\sigma^2 + v_j)}, \quad \text{var}(\mu \mid y) = \frac{1}{1/\tau^2 + \sum 1/(\sigma^2 + v_j)}, \\ E(\theta_j \mid y) &= \frac{\sigma^2 y_j + v_j E(\mu \mid y)}{\sigma^2 + v_j}, \quad \text{var}(\theta_j \mid y) = \frac{1 + \text{var}(\mu \mid y)/\sigma^2}{1/v_j + 1/\sigma^2}. \end{aligned}$$

Discuss.

stat.epfl.ch

Autumn 2022 – slide 231

Note to Example 75

- The y_j have different variances, but their means θ_j are supposed indistinguishable and hence are modelled as exchangeable, being normal with unknown mean μ , and we can write

$$y_j = \mu_0 + (\mu - \mu_0) + (\theta_j - \mu) + (y_j - \theta_j),$$

where μ_0 is known, and as the y_j and θ_j are linear combinations of normal variables it is straightforward to check that

$$\begin{pmatrix} \mu \\ \theta \\ y \end{pmatrix} \sim \mathcal{N}_{2n+1} \left\{ \mu_0 \mathbf{1}_{2n+1}, \begin{pmatrix} \tau^2 & \tau^2 \mathbf{1}_n^\top & \tau^2 \mathbf{1}_n^\top \\ \tau^2 \mathbf{1}_n & \tau^2 \mathbf{1}_n \mathbf{1}_n^\top + \sigma^2 I_n & \tau^2 \mathbf{1}_n \mathbf{1}_n^\top + \sigma^2 I_n \\ \tau^2 \mathbf{1}_n & \tau^2 \mathbf{1}_n \mathbf{1}_n^\top + \sigma^2 I_n & V + \tau^2 \mathbf{1}_n \mathbf{1}_n^\top + \sigma^2 I_n \end{pmatrix} \right\}, \quad (13)$$

where $\mathbf{1}_n$ denotes the $n \times 1$ vector of ones and $V = \text{diag}(v_1, \dots, v_n)$.

- The most direct approach to computing the posterior distributions μ and θ given y is to write

$$\begin{pmatrix} \mu \\ \theta \\ y \end{pmatrix} \sim \mathcal{N}_{2n+1} \left\{ \mu_0 \mathbf{1}_{2n+1}, \begin{pmatrix} \Omega_{11} & \Omega_{12} \\ \Omega_{21} & \Omega_{22} \end{pmatrix} \right\},$$

where $\text{var}(y) = \Omega_{22}$. Then the posterior density of the parameters given y is also normal, with

$$\begin{pmatrix} \mu \\ \theta \end{pmatrix} | y \sim \mathcal{N}_{n+1} \{ \mu_0 \mathbf{1}_{n+1} + \Omega_{12} \Omega_{22}^{-1} (y - \mu_0 \mathbf{1}_n), \Omega_{11} - \Omega_{12} \Omega_{22}^{-1} \Omega_{21} \}. \quad (14)$$

We shall take a less messy and maybe more enlightening route, first computing the posterior distribution of μ , then that of θ given both μ and y , and then marginalising the latter over μ .

- Expression (13) shows that the joint density of μ and y is normal with covariance matrix

$$\begin{pmatrix} A & B^\top \\ B & C \end{pmatrix}, \quad A = \tau^2, \quad B = \tau^2 \mathbf{1}_n, \quad C = \tau^2 \mathbf{1}_n \mathbf{1}_n^\top + D, \quad D = \text{diag}(\sigma^2 + v_1, \dots, \sigma^2 + v_n).$$

The Woodbury formula gives

$$(D + \tau^2 \mathbf{1}_n \mathbf{1}_n^\top)^{-1} = D^{-1} - D^{-1} \mathbf{1}_n (\tau^{-2} + \mathbf{1}_n^\top D^{-1} \mathbf{1}_n)^{-1} \mathbf{1}_n D^{-1}$$

so with $a = \mathbf{1}_n^\top D^{-1} \mathbf{1}_n$ we have

$$\begin{aligned} A - B C^{-1} B^\top &= \tau^2 - \tau^2 \mathbf{1}_n^\top \{ D^{-1} - D^{-1} \mathbf{1}_n (\tau^{-2} + \mathbf{1}_n^\top D^{-1} \mathbf{1}_n)^{-1} \mathbf{1}_n D^{-1} \} \tau^2 \mathbf{1}_n \\ &= \tau^2 - \tau^4 \left\{ a - \frac{a^2}{\tau^{-2} + a} \right\} \\ &= (\tau^{-2} + a)^{-1}, \end{aligned}$$

which gives $\text{var}(\mu | y)$, and a simpler calculation using (14) with μ only gives the mean, resulting in

$$\mathbb{E}(\mu | y) = \frac{\mu_0/\tau^2 + \sum y_j/(\sigma^2 + v_j)}{1/\tau^2 + \sum 1/(\sigma^2 + v_j)}, \quad \text{var}(\mu | y) = \frac{1}{1/\tau^2 + \sum 1/(\sigma^2 + v_j)}.$$

The posterior mean of μ is a weighted average of its prior mean μ_0 and of the y_j , weighted according to their precisions. Typically τ^2 is taken to be very large, and then $\mathbb{E}(\mu | y)$ is essentially a weighted average of the data. Even when $v_j \rightarrow 0$ for all j there is still posterior uncertainty about μ , whose variance is σ^2/n because y_1, \dots, y_n is then a random sample from $N(\mu, \sigma^2)$.

Note 2 to Example 75

- To compute the posterior mean and variance of θ_j we note that the graph structure gives $f(\theta_j | \mu, y) = f(\theta_j | \mu, y_j)$. This simplifies the computation because we need only compute the joint distribution of (μ, θ_j, y_j) , and this is

$$\mathcal{N}_3 \left\{ 1_3 \mu_0, \begin{pmatrix} \tau^2 & \tau^2 & \tau^2 \\ \tau^2 & \tau^2 + \sigma^2 & \tau^2 + \sigma^2 \\ \tau^2 & \tau^2 + \sigma^2 & \tau^2 + \sigma^2 + v_j \end{pmatrix} \right\}$$

from which we obtain $\theta_j | \mu, y_j \sim \mathcal{N}\{(y_j/v_j + \mu/\sigma^2)/(1/v_j + 1/\sigma^2), (1/v_j + 1/\sigma^2)^{-1}\}$. As

$$E(\theta_j | y) = E\{E(\theta_j | \mu, y_j)\}, \quad \text{var}(\theta_j | y) = E\{\text{var}(\theta_j | \mu, y_j)\} + \text{var}\{E(\theta_j | \mu, y_j)\},$$

where the outer expectation and variance are over the distribution of μ given y , we finally obtain

$$E(\theta_j | y) = \frac{\sigma^2 y_j + v_j E(\mu | y)}{\sigma^2 + v_j}, \quad \text{var}(\theta_j | y) = \frac{1 + \text{var}(\mu | y)/\sigma^2}{1/v_j + 1/\sigma^2}.$$

- The posterior mean of θ_j is a weighted average of y_j and $E(\mu | y)$, showing shrinkage of y_j towards $E(\mu | y)$ by an amount that depends on v_j . As $v_j \rightarrow 0$, $E(\theta_j | y) \rightarrow y_j$, while as $v_j \rightarrow \infty$, $E(\theta_j | y) \rightarrow E(\mu | y)$. This is a characteristic feature of hierarchical models, in which there is a 'borrowing of strength' whereby all the data combine to estimate common parameters such as μ , while estimates of individual parameters such as the θ_j are shrunk towards common values by amounts that depend on the precisions v_j of the corresponding observations.

Example: Cardiac surgery data

<i>A</i>	0/47	<i>B</i>	18/148	<i>C</i>	8/119	<i>D</i>	46/810	<i>E</i>	8/211	<i>F</i>	13/196
<i>G</i>	9/148	<i>H</i>	31/215	<i>I</i>	14/207	<i>J</i>	8/97	<i>K</i>	29/256	<i>L</i>	24/360

Mortality rates r/m from cardiac surgery in 12 hospitals (numbers of deaths r out of m operations).

- Hierarchical model:

$$r_j | \theta_j \stackrel{\text{ind}}{\sim} B(m_j, \theta_j), \quad j = A, \dots, L, \quad \theta_A, \dots, \theta_L | \zeta \stackrel{\text{iid}}{\sim} f(\theta | \zeta), \quad \zeta \sim \pi(\zeta).$$

Conditional on θ_j , the number of deaths r_j at hospital j is binomial with probability θ_j and denominator m_j , the number of operations, which plays the same role as v_j^{-1} in the normal example above: when m_j is large then a death rate is relatively precisely known.

- Conditional on ζ , the θ_j are a random sample from a distribution $f(\theta | \zeta)$, and the prior distribution for ζ depends on fixed hyperparameters.
- We take $\beta_j = \log\{\theta_j/(1 - \theta_j)\} \sim N(\mu, \sigma^2)$, conditional on $\zeta = (\mu, \sigma^2)$, and $\mu \sim N(0, c^2)$ and $\sigma^2 \sim IG(a, b)$, with $a = b = 10^{-3}$, so σ^2 has prior mean one but variance 10^3 , and $c = 10^3$, giving μ prior variance 10^6 .

Example: Cardiac surgery data

- The joint density is

$$\left[\prod_j \binom{m_j}{r_j} \frac{e^{r_j \beta_j}}{(1 + e^{\beta_j})^{m_j}} \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\beta_j - \mu)^2 \right\} \right] \times \pi(\mu) \pi(\sigma^2),$$

so the full conditional densities for μ and σ^2 are normal and inverse gamma.

- We use a Metropolis–Hastings step for β , using a random walk proposal with

$$\beta'_j \sim \mathcal{N}\{\beta_j, d^2 \sigma^2 v_j / (\sigma^2 + v_j)\}, \quad v_j = \frac{m_j + 1}{(r_j + 1/2)(m_j - r_j + 1/2)},$$

where we choose d to optimise the algorithm.

- This normal approximation comes from Example 75, taking

$$\hat{\beta}_j \mid \beta \sim \mathcal{N}(\beta, v_j), \quad \beta_j \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2),$$

and then computing $\beta_j \mid \hat{\beta}_j$.

Example: Cardiac surgery data

```
cardiac.gibbs <- function(data, mu0=0, a=10^(-3), b=10^(-3), c=10^3, R=10^5, d=1)
{ # parameter is mu, sig2, beta
  card.update <- function(data, mu0, a, b, c, para)
  {
    sig2 <- para[2]
    beta <- para[-c(1,2)]
    n <- length(beta)
    mu <- rnorm( 1, (mu0/c^2 + sum(beta)/sig2)/(1/c^2+n/sig2),sqrt(1/(1/c^2+n/sig2)) )
    sig2 <- rigamma( a+n/2, b+0.5*sum((beta-mu)^2) )
    v <- (data$m+1)/((data$r+0.5)*(data$m-data$r+0.5))
    var.beta <- sig2*v/(v+sig2)
    beta.prop <- rnorm(n, beta, sd=d*sqrt(var.beta))
    acc.prob <- exp( data$r*beta.prop - data$m*log(1+exp(beta.prop)) -
                    0.5*(beta.prop-mu)^2/sig2 - data$r*beta +
                    data$m*log(1+exp(beta)) + 0.5*(beta-mu)^2/sig2 )
    acc.prob <- pmin(1,acc.prob) # use pmin and ifelse to do all
    beta <- ifelse(runif(n)<=acc.prob,beta.prop, beta) # acceptances/rejections at once
    c( mu, sig2, beta)
  }
  rigamma <- function(a, b) 1/rgamma(1, shape=a, rate=b)
  logit <- function(p) log(p/(1-p))
  out <- matrix(NA, 2+nrow(data), R)
  out[, 1] <- c(0, 1, rep(0,nrow(data)))
  for(r in 2:R)
    out[, r] <- card.update(data, mu0, a, b, c, out[,r-1])
  out
}

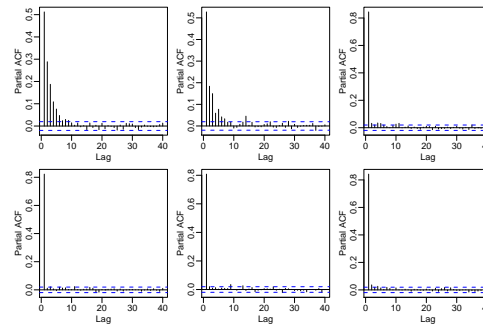
system.time( cardiac.sim <- cardiac.gibbs(cardiac, R=10^4, d=4) ) # around 3.5 seconds
acc.rate <- function(x) mean((diff(x)!=0))
apply(cardiac.sim,1,acc.rate) # compute acceptance rates for the proposals
```

Effect of d

Acceptance probabilities for different values of d :

d	0.1	0.5	1	2	3	5	10	20	30
μ	1	1	1	1	1	1	1	1	1
σ^2	1	1	1	1	1	1	1	1	1
β	0.95	0.82	0.7	0.5	0.37	0.25	0.12	0.06	0.05

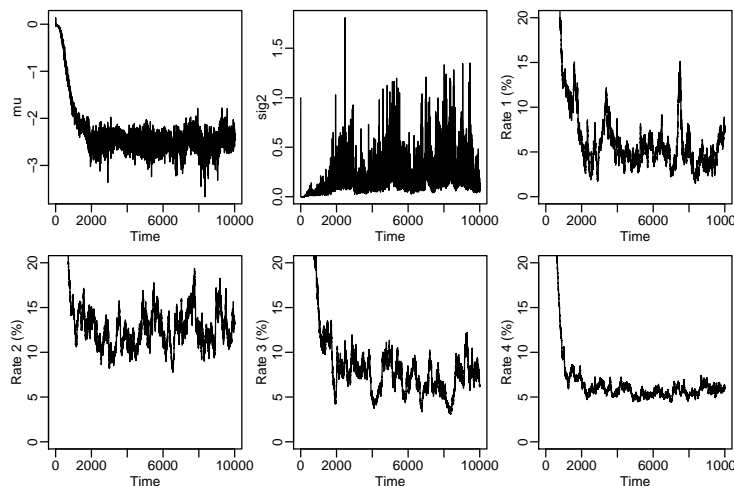
PACF for $d = 1$:



stat.epfl.ch

Autumn 2022 – slide 235

Effect of d : $d = 0.1$

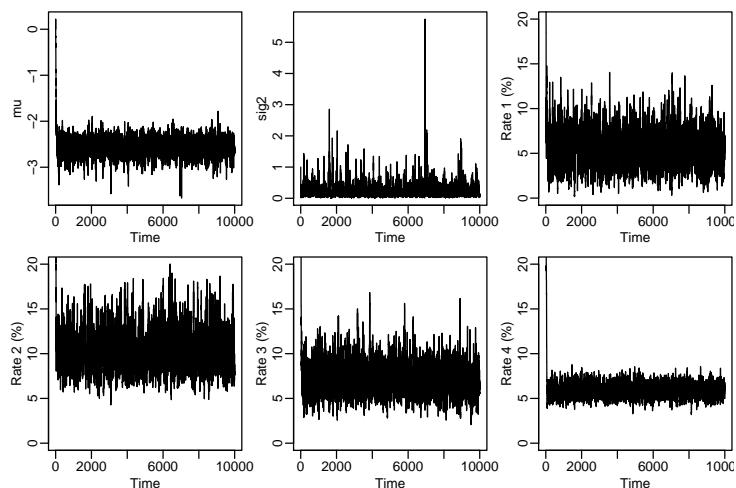


Taking $d = 0.1$ makes the acceptance probability too high, so the chain mixes too slowly.

stat.epfl.ch

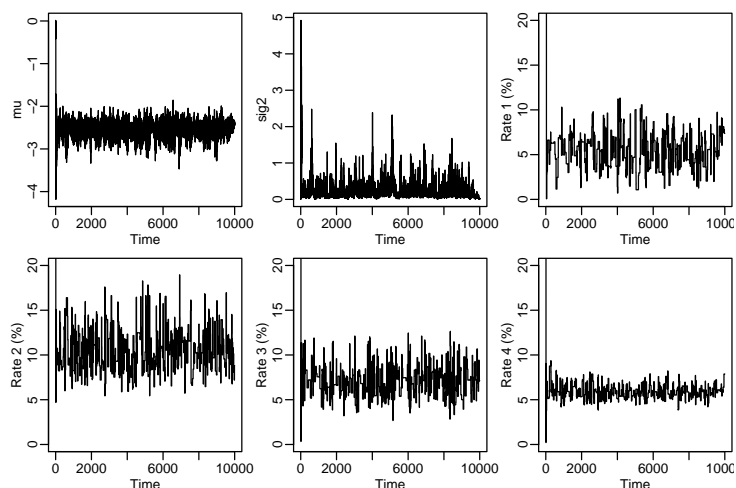
Autumn 2022 – slide 236

Effect of d : $d = 1$



Taking $d = 1$ is OK, but theory suggest that the acceptance rate should be around 0.2–0.4, so taking $d \approx 4$ seems somewhat better.

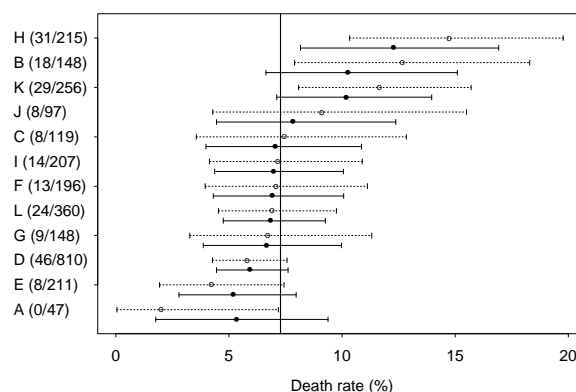
Effect of d : $d = 30$



Taking $d = 30$ makes the acceptance probability too low, so the chain sticks.

Example: Cardiac surgery data, effect of shrinkage

Posterior means and 0.95 equitailed credible intervals for separate analyses for each hospital are shown by hollow circles and dotted lines, while blobs and solid lines show the corresponding quantities for a hierarchical model. Note the shrinkage ('borrowing of strength') of the estimates for the hierarchical model towards the overall posterior mean rate, shown as the solid vertical line; the hierarchical intervals are slightly shorter than those for the simpler model.



stat.epfl.ch

Autumn 2022 – slide 239

Summary

- ☐ Graphical representation of dependence relations very useful.
- ☐ Hierarchical modelling allows us to fit complex models to data
- ☐ Key idea is to treat parameters as coming from a distribution, and to use the data to estimate the distribution
 - Appropriate when exchangeable elements are present
 - Not appropriate when we are interested in certain pre-specified parameters or where prior knowledge distinguishes them
 - Example of inappropriate use: economic modelling with countries of Europe treated as exchangeable
- ☐ Can be hard to count the number of parameters: the prior 'ties together' some parameters, so there are 'really' fewer—but how many?

stat.epfl.ch

Autumn 2022 – slide 240