Statistical Inference

Anthony Davison

©2024

http://stat.epfl.ch

1 Introduction	2
1.1 Background	3
1.2 Probability Revision	9
1.3 Statistics Revision	27
1.4 Bases for Uncertainty	47
2 Some Basic Concepts	59
2.1 Likelihood	60
2.2 Complications	65
2.3 Data Reduction	76
2.4 Inference	85
3 Likelihood Theory	91
3.1 Basic Results	92
3.2 Vector Parameter	108
3.3 Nuisance Parameters	113
4 Hypothesis Testing	125
4.1 Pure Significance Tests	126
4.2 Neyman-Pearson Approach	134
4.3 Multiple Testing	144
4.4 Post-Selection Inference	155
5 Bayesian Statistics	159

5.1 Introduction	160
5.2 Bayesian Inference	176
5.3 Bayesian Computation	193
5.4 Hierarchical Models	202
Appendix I: Monte Carlo Methods	215
Appendix II: Graphical Models	232

1 Introduction slide 2

slide 3

1.1 Background

Sta	arting 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 reach a decision to guide future actions.
	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

stat.epfl.ch Autumn 2024 – slide 4

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?
- ☐ For now we suppose these questions have satisfactory answers . . .

So	me statistical activities
	Conventionally divided into
	 design of investigations — how do we get reliable data to answer a question efficiently and securely?
	 descriptive statistics/exploratory data analysis — how can we get insight into a specific dataset?
	 inference — what can we learn about the properties of a 'population' underlying the data? decision analysis — what is the optimal decision in a given situation?
	to which we nowadays add
	 machine learning — algorithms, generally complex and computationally demanding, often used for prediction/decision-making.
stat.e	epfl.ch Autumn 2024 – slide 6
De	escriptive statistics
	In principle concerns only the data available, mainly involving
	 graphical summaries — histograms, boxplots, scatterplots,
	 numerical summaries — averages, variances, medians,
	Some summaries presuppose the existence of 'population' quantities (e.g., a density).
	We use probability models to analyse the properties of these summaries (e.g., formulation of a boxplot, 'is that difference significant?',).

stat.epfl.ch Autumn 2024 – slide 7

□ Even when we have 'all the data' (e.g., loyalty card transactions) we may want to ask 'what if?' questions, and these require further assumptions (e.g., temporal stability, future and current

customers are similar, ...).

Statistical inference		
	Use observed data to draw conclusions about a 'population' from which the data are assumed to be drawn, or about future data.	
	The 'population' and observed data are linked by concepts of probability.	
	Two distinct roles of probability in statistical analysis:	
	- as a description of variation in data ('aleatory probability', 'chance'), treating the observed data y as an outcome of a random process/probability model, perhaps	
	imposed by the investigator (via some sampling procedure);	
	- to formulate uncertainty ('epistemic probability') about the reality modelled in terms of the random experiment, based on y .	
	Most of the course concerns the formulation and expression of uncertainty.	
	We first revise some concepts from probability and basic statistics.	

Probability spaces

- \square Ordered triples (Ω, \mathcal{F}, P) consisting of
 - a set Ω of elementary outcomes ω corresponding to distinct potential outcomes of a random experiment;
 - an event space \mathcal{F} of subsets of Ω that satisfy (a) $\Omega \in \mathcal{F}$, (b) if $\mathcal{A} \in \mathcal{F}$, then $\mathcal{A}^c \in \mathcal{F}$, and (c) if $\mathcal{A}_1, \mathcal{A}_2, \ldots \in \mathcal{F}$, then $\bigcup A_j \in \mathcal{F}$;
 - a probability measure $P: \mathcal{F} \to [0,1]$ that satisfies (i) if $\mathcal{A} \in \mathcal{F}$, then $0 \leq P(\mathcal{A}) \leq 1$, (ii) $P(\Omega) = 1$, (iii) if $\mathcal{A}_1, \mathcal{A}_2, \ldots \in \mathcal{F}$ satisfy $\mathcal{A}_j \cap \mathcal{A}_k = \emptyset$ for $j \neq k$, then $P(\bigcup \mathcal{A}_j) = \sum P(\mathcal{A}_j)$.
- \square We call (Ω, \mathcal{F}) a measure space and any $A \in \mathcal{F}$ an event (measurable set).
- \square From these we deduce
 - the inclusion-exclusion formulae, and
 - computation of probabilities in simple problems using combinatorial formulae.
- \square If $P(\mathcal{B}) > 0$ we define conditional probabilities $P(\mathcal{A} \mid \mathcal{B}) = P(\mathcal{A} \cap \mathcal{B})/P(\mathcal{B})$, and derive
 - a new conditional probability distribution $P_{\mathcal{B}}(\mathcal{A}) = P(\mathcal{A} \mid \mathcal{B})$ for $\mathcal{A} \in \mathcal{F}$,
 - the law of total probability,
 - Bayes' theorem, and
 - the notion of independent events, for which $P(A \cap B) = P(A)P(B)$.

stat.epfl.ch Autumn 2024 – slide 10

Random variables

- Let (Ω, \mathcal{F}, P) be a probability space and $(\mathcal{X}, \mathcal{G})$ a measurable space. A random function X from Ω into \mathcal{X} has the property that $X^{-1}(\mathcal{C}) = \{\omega : X(\omega) \in \mathcal{C}\} \in \mathcal{F}$ for any $\mathcal{C} \in \mathcal{G}$, so $P(X \in \mathcal{C}) = P\{X^{-1}(\mathcal{C})\}$ is well-defined. Such a function is called measurable.
- If $\mathcal{X} = \mathbb{R}$ or \mathbb{R}^n we call X a random variable and there exists a cumulative distribution function (CDF) F such that $P\{X \in (-\infty, x_1] \times \cdots \times (-\infty, x_n]\} = F(x_1, \dots, x_n)$.
- \square A CDF increases from 0 when any of its arguments increases from $-\infty$ to $+\infty$.
- \Box F can be written as a sum of (sub-)distributions $F_{\rm ac}+F_{\rm dis}+F_{\rm sing}$, where
 - $F_{\rm ac}$ is absolutely continuous, i.e., there exists a non-negative probability density function (PDF) $f_{\rm ac}(x) = {\rm d}F_{\rm ac}(x)/{\rm d}x$,
 - $F_{\rm dis}$ is discrete, i.e., its probability mass function (PMF) $f_{\rm dis}(x)$ is positive only on a finite or countable set \mathcal{S} , and
 - $F_{\rm sing}$ is singular, and can be ignored (look up 'Cantor distribution' if interested).
- \square We call X continuous or discrete respectively if F_{dis} or F_{ac} is absent.
- ☐ If necessary we use **Lebesgue-Stieltjes integration**, whereby

$$P(X \in \mathcal{C}) = \int_{\mathcal{C}} dF(x) = \int_{\mathcal{C}} f_{ac}(x) dx + \sum_{x \in \mathcal{C} \cap \mathcal{S}} f_{dis}(x), \quad \mathcal{C} \subset \mathcal{X};$$

the notation \int_a^b is unwise because it doesn't distinguish $\mathcal{C} = [a,b]$ from $\mathcal{C} = (a,b)$.

New distributions and new random variables

 \square We define the **conditional distribution** of X given an event $\mathcal{B} \in \mathcal{F}$ by

$$P(X \in \mathcal{A} \mid \mathcal{B}) = P(\{X \in \mathcal{A}\} \cap \mathcal{B})/P(\mathcal{B}).$$

 \square If $Y = g(X) \in \mathcal{Y}$ and we write $g^{-1}(\mathcal{B}) = \{x : g(x) \in \mathcal{B}\}$ for $\mathcal{B} \subset \mathcal{Y}$, then

$$P(Y \in \mathcal{B}) = P\{g(X) \in \mathcal{B}\} = P\{X \in g^{-1}(\mathcal{B})\}.$$

 \square If X is continuous and Y = g(X) with g a smooth bijection, then (in obvious notation)

$$f_Y(y) = f_X\{g^{-1}(y)\} \left| \frac{\partial g^{-1}(y)}{\partial y} \right|,$$

where the last term is the Jacobian of the transformation.

 \square If $X = (X_1, X_2)$ is continuous, we obtain marginal and conditional densities

$$f_{X_2}(x_2) = \int f_{X_1, X_2}(x_1, x_2) \, \mathrm{d}x_1, \quad f_{X_1 \mid X_2}(x_1 \mid x_2) = \frac{f_{X_1, X_2}(x_1, x_2)}{f_{X_2}(x_2)},$$

with corresponding formulae in the discrete and mixed cases.

 \sqsupset X_1 and X_2 are independent $(X_1 \perp \!\!\! \perp X_2)$ iff $f_{X_1,X_2}(x_1,x_2) = f_{X_1}(x_1)f_{X_2}(x_2)$, $\forall x_1,x_2$.

stat.epfl.ch

Autumn 2024 - slide 12

Exchangeability

- □ Exchangeability is weaker than independence, often used to model variables that are indistinguishable in probabilistic terms, even if not independent.
- de Finetti proved that such variables must be constructed as $U_1, \ldots, U_n \mid \theta \stackrel{\text{iid}}{\sim} F_{\theta}$, where $\theta \sim G$ for distributions F_{θ} and G. The simplest theorem to this effect is the one below.

Definition 1 Random variables U_1, \ldots, U_n are finitely exchangeable if their density satisfies

$$f(u_1,\ldots,u_n)=f\left(u_{\xi(1)},\ldots,u_{\xi(n)}\right)$$

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.

Theorem 2 (de Finetti) If $U_1, U_2, ...$, is an infinitely exchangeable sequence of binary variables taking values in $\{0,1\}$, then for any n there is a distribution G such that

$$f(u_1, \dots, u_n) = \int_0^1 \prod_{j=1}^n \theta^{u_j} (1 - \theta)^{1 - u_j} G(d\theta)$$
 (1)

where

$$G(\theta) = \lim_{m \to \infty} P\left\{m^{-1}(U_1 + \dots + U_m) \le \theta\right\}, \quad \theta = \lim_{m \to \infty} m^{-1}(U_1 + \dots + U_m).$$

stat.epfl.ch

Terminology and notation

- □ PDFs and PMFs are not the same but we henceforth use the term density for both.
- \square $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} f$ means that the X_j are independent and all have density f, and we then call the X_j a random sample (of size n) from f.
- $\square \quad X_1,\ldots,X_n \overset{\mathrm{ind}}{\sim} f_1,\ldots,f_n$ means that the X_j are independent and $X_j \sim f_j.$
- $\square \quad X_1,\ldots,X_n \overset{\mathrm{ind}}{\sim} (\mu,\sigma^2) \text{ means that the } X_j \text{ are independent with mean } \mu \text{ and variance } \sigma^2 \text{ (with } 0<\sigma^2<\infty\text{)}.$ The X_j need not be normal or have the same distribution.
- \square $X_1,\ldots,X_n \stackrel{\mathrm{ind}}{\sim} (\mu_1,\ldots,\mu_n,\sigma_1^2,\ldots,\sigma_n^2)$ means that the X_j are independent with means μ_j and variances σ_j^2 (where $0<\sigma_j^2<\infty$).
- \square The p quantile of the distribution F of a scalar random variable X is

$$x_p = \inf\{x : F(x) \ge p\}, \quad 0$$

Usually $x_p = F^{-1}(p)$ for continuous X, but not for discrete (or mixed) X.

 $\hfill \square$ A standard normal variable $Z \sim \mathcal{N}(0,1)$ has PDF and CDF

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

and p quantile $z_p = \Phi^{-1}(p)$, so $X = \mu + \sigma Z \sim \mathcal{N}(\mu, \sigma^2)$ has p quantile $\mu + \sigma z_p$.

stat.epfl.ch Autumn 2024 – slide 14

Order statistics

 \square The order statistics of $X_1,\ldots,X_n\stackrel{\mathrm{iid}}{\sim} f$ are the ordered values

$$X_{(1)} \le X_{(2)} \le \dots \le X_{(n-1)} \le X_{(n)}.$$

 \square In particular, the minimum is $X_{(1)}$, the maximum is $X_{(n)}$, and the median is

$$X_{(m+1)}$$
 $(n=2m+1, \text{ odd}), \frac{1}{2}(X_{(m)}+X_{(m+1)})$ $(n=2m, \text{ even}).$

The median is the central value of X_1, \ldots, X_n .

 \square If f is continuous then the X_j must be distinct, and for $r=1,\ldots,n$ we have

$$P(X_{(r)} \le x) = \sum_{j=r}^{n} {n \choose j} F(x)^{j} \{1 - F(x)\}^{n-j},$$

$$f_{X_{(r)}}(x) = \frac{n!}{(r-1)! \, 1! \, (n-r)!} F(x)^{r-1} f(x) \{1 - F(x)\}^{n-r}.$$

 \square Joint densities can be obtained using the argument that gives $f_{X_{(r)}}(x)$, and in particular

$$f_{X_{(1)},\dots,X_{(n)}}(x_1,\dots,x_n) = n! f(x_1) \cdots f(x_n), \quad x_1 < \dots < x_n.$$

Example 3 Find the joint density of $X_{(2)}, \ldots, X_{(n-1)}$ given that $X_{(1)} = x_1$ and $X_{(n)} = x_n$.

Note: Densities of order statistics

- The event $X_{(r)} \leq x$ occurs iff at least r of the independent variables X_1, \ldots, X_n are less than or equal to x, and each of them does this with probability F(x). Hence the probability of the event is given by a binomial probability, and a little thought shows that this is the stated formula.
- The density can be obtained by differentiation of $P(X_{(r)} \leq x)$, whereupon one finds that almost all the terms cancel, giving the stated density. A more easily generalised argument is as follows: for the event $X_{(r)} \in [x, x+\mathrm{d}x)$, we need to split the sample into three groups of respective sizes r-1, 1 and n-r and 'probabilities' F(x), $f(x)\mathrm{d}x$, and 1-F(x). The corresponding multinomial 'probability' is

$$\frac{n!}{(r-1)! \times 1! \times (n-r)!} \{F(x)\}^{r-1} \times f(x) dx \times \{1 - F(x)\}^{n-r},$$

and dropping the dx gives the density function of $X_{(r)}$.

 \square For the joint density we divide the sample into n parts, each with one observation, and apply a version of the multinomial argument just given.

stat.epfl.ch

Autumn 2024 - note 1 of slide 15

Note to Example 3

The joint density of $X_{(1)}$ and $X_{(n)}$ is given by splitting the total n observations into three parts, with respective 'probabilities' $f(x_1)dx_1$, $F(x_n) - F(x_1)$ and $f(x_n)dx_n$ and sizes 1, n-2 and 1, giving

$$f_{X_{(1)},X_{(n)}}(x_1,x_n)\mathrm{d}x_1\mathrm{d}x_n = \frac{n!}{1!(n-2)!1!}f(x_1)\mathrm{d}x_1 \times \{F(x_n) - F(x_1)\}^{n-2} \times f(x_n)\mathrm{d}x_n, \quad x_1 < x_n.$$

We drop the dx_1dx_n to get the joint density.

 \square Hence the conditional density of $X_{(2)},\ldots,X_{(n-1)}$ given that $X_{(1)}=x_1$ and $X_{(n)}=x_n$ is

$$\frac{n!f(x_1)\cdots f(x_n)}{n!/(n-2)!\times f(x_1)\{F(x_n)-F(x_1)\}^{n-2}f(x_n)} = (n-2)!\prod_{j=2}^{n-1}\frac{f(x_j)}{F(x_n)-F(x_1)},$$

where $x_1 < x_2 < \cdots < x_{n-1} < x_n$. This is the joint density of the order statistics of a random sample of size n-2 from the truncated distribution $f(x)/\{F(x_n)-F(x_1)\}$, where $x_1 < x < x_n$

stat.epfl.ch

Autumn 2024 - note 2 of slide 15

Moments

 \square The expectation $\mathrm{E}\{g(X)\}$ of g(X) is defined if $\mathrm{E}\{|g(X)|\}<\infty$ as

$$E\{g(X)\} = \int_{\mathcal{X}} g(x) \, \mathrm{d}F(x).$$

 \square For scalar X we define moments $\mathrm{E}(X^r)$, mean $\mu=\mathrm{E}(X)$ and variance

$$var(X) = E[\{X - E(X)\}^2] = E(X^2) - E(X)^2 = E\{X(X - 1)\} + E(X) - E(X)^2.$$

- \square var(X) = 0 iff X is constant with probability one.
- \square For vector X we define the **mean vector** and **(co)variance matrix**

$$\mu = E(X), \quad cov(X_1, X_2) = E(X_1 X_2^T) - E(X_1)E(X_2)^T,$$

and write $var(X) = cov(X, X) = E\{(X - \mu)(X - \mu)^{T}\}.$

- \square The correlation, $\operatorname{corr}(X_1, X_2) = \operatorname{cov}(X_1, X_2) / \{\operatorname{var}(X_1)\operatorname{var}(X_2)\}^{1/2}$, is a measure of dependence between variables that does not depend on their units of measurement.
- \square Expectation $E(\cdot)$ is a linear operator, so it is easy to check that

$$E(a + BX) = a + BE(X), \quad cov(a + BX, c + DX) = Bvar(X)D^{T}.$$

stat.epfl.ch

Autumn 2024 - slide 16

Conditional moments

 $\ \square$ The conditional expectation of g(X,Y) given X=x is

$$E\{g(X,Y) \mid X = x\} = \int_{\mathcal{V}} g(x,y) \, dF(y \mid x),$$

which in the continuous and discrete cases equals

$$\int_{\mathcal{Y}} g(x, y) f_{Y|X}(y \mid x) \, \mathrm{d}y, \quad \sum_{y \in \mathcal{Y}} g(x, y) f_{Y|X}(y \mid x),$$

and other conditional moments are defined likewise.

- \square This is a function of x, so it defines a random variable $\tilde{g}(X) = \mathbb{E}\{g(X,Y) \mid X\}$.
- ☐ The law of total expectation (tower property) gives

$$\mathbb{E}\left\{g(X,Y)\right\} = \mathbb{E}_X\left[\mathbb{E}\left\{g(X,Y) \mid X=x\right\}\right],$$

$$\operatorname{var} \{g(X,Y)\} = \operatorname{E}_X \left[\operatorname{var} \{g(X,Y) \mid X = x\} \right] + \operatorname{var}_X \left[E\{g(X,Y) \mid X = x\} \right],$$

where E_X denotes expectation with respect to the marginal distribution of X, etc., with a similar expression (which you should give) for $cov\{g(X,Y),h(X,Y)\}$.

□ We ignore mathematical issues arising from conditioning on events of probability zero — look up 'Borel–Kolmogorov paradox' if interested.

stat.epfl.ch

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 $a^{\mathrm{T}}X \sim \mathcal{N}(a^{\mathrm{T}}\mu, a^{\mathrm{T}}\Omega a)$ for every constant vector $a_{n \times 1}$, and then

 \square $M_Y(t) = \exp(t^{\mathrm{T}}\mu + \frac{1}{2}t^{\mathrm{T}}\Omega t)$ and the mean vector and covariance matrix of X are

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

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

 \Box for any constants $a_{m\times 1}$ and $B_{m\times n}$

$$a + BX \sim \mathcal{N}_m \left(a + B\mu, B\Omega B^{\mathrm{T}} \right);$$

- \Box a + BX and c + DX are independent iff $B\Omega D^{\mathrm{T}} = 0$;
- \square X has a density on \mathbb{R}^n iff Ω is positive definite (i.e., has rank n), and then

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

 \square if $X^{\mathrm{T}}=(X_1^{\mathrm{T}},X_2^{\mathrm{T}})$, 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:

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

stat.epfl.ch Autumn 2024 – slide 18

MGFs and KGFs

 \square The moment-generating function (MGF) and cumulant-generating function (KGF) of a scalar random variable X are

$$M_X(t) = \mathbb{E}\left(e^{tX}\right), \quad K_X(t) = \log M_X(t), \quad t \in \mathcal{N} = \{t : M_X(t) < \infty\}.$$

 \square \mathcal{N} is non-empty, because $M_X(0)=1$, but the MGF and KGF are non-trivial only if \mathcal{N} contains an open neighbourhood of the origin, since then

$$M_X(t) = \mathbb{E}\left(\sum_{r=0}^{\infty} \frac{t^r X^r}{r!}\right) = \sum_{r=0}^{\infty} \frac{t^r}{r!} \mathbb{E}(X^r), \quad K_X(t) = \sum_{r=1}^{\infty} \frac{t^r}{r!} \kappa_r,$$

and one can obtain the moments $E(X^r)$ and cumulants κ_r by differentiation.

☐ In the vector case we define

$$M_X(t) = \mathbb{E}\left(e^{t^{\mathrm{T}}X}\right), \quad K_X(t) = \log M_X(t),$$

and differentiation with respect to the elements of $t=(t_1,\ldots,t_n)^{\rm T}$ gives the mean vector and covariance matrix of X.

- ☐ There is a 1–1 mapping between distributions and MGFs/KGFs (if the latter are non-trivial).
- \square KGFs for linear combinations are computed as $K_{a+BX}(t) = a^{\mathrm{T}}t + K_X(B^{\mathrm{T}}t)$.

Note: Moments and cumulants

- \square We consider scalar X, as the calculations for vector X are analogous.
- \square First note that $M_X(t)=1$ when t=0, since $\mathrm{E}(e^{tX})=\mathrm{E}(1)=1$; thus $0\in\mathcal{N}$ for any X.
- \square If $\mathcal N$ contains an open set (-a,a) for some a>0, and $\mu_r=\mathrm E(X^r)$ denotes the rth moment of X, then if |t|< a,

$$K_X(t) = \sum_{r=1}^{\infty} \frac{t^r \kappa_r}{r!} = \log M_X(t) = \log \left(\sum_{r=0}^{\infty} \frac{t^r \mu_r}{r!} \right) = \log(1+b) = b - b^2/2 + b^3/3 + \cdots,$$

where $b=t\mu_1+t^2\mu_2/2!+t^3\mu_3/3!+\cdots$. If we expand and compare coefficients of t,t^2,t^3,\ldots in the two expansions we get

$$\kappa_1 = \mu_1, \quad \kappa_2 = \mu_2 - \mu_1^2, \quad \kappa_3 = \mu_3 - 3\mu_2\mu_1 + 2\mu_1^3, \quad \kappa_4 = \mu_4 - 4\mu_3\mu_1 + 6\mu_2\mu_1^2 - 3\mu_1^4, \quad \dots,$$

so $\kappa_1 = \mathrm{E}(X), \; \kappa_2 = \mathrm{var}(X), \; \kappa_3 = \mathrm{E}\{(X - \mu_1)^3\}, \; \dots$

stat.epfl.ch

Autumn 2024 - note 1 of slide 19

Exponential tilting

 \square A baseline density f_0 with a non-trivial MGF can be used to construct a family of densities by exponential tilting, i.e.,

$$f(y;\varphi) = f_0(y) \exp \{ \varphi^{\mathrm{T}} s(y) - k(\varphi) \}, \quad y \in \mathcal{Y}, \varphi \in \mathcal{N},$$

where

$$\mathcal{N} = \{ \varphi : k(\varphi) < \infty \}$$

and individual members of the family are determined by the value of φ .

☐ Hölder's inequality gives

$$M\{\alpha\varphi_1 + (1-\alpha)\varphi_2\} \le M(\varphi_1)^{\alpha}M(\varphi_2)^{1-\alpha} < \infty, \quad 0 \le \alpha \le 1,$$

for any $\varphi_1, \varphi_2 \in \mathcal{N}$, so the set \mathcal{N} and the function k are both convex.

- \square This implies that $f(y;\varphi)$ is log-concave in φ , which is very useful for statistics.
- ☐ This construction leads to an elegant general theory putting many well-known distributions (Poisson, binomial, normal, ...) under the same roof.

Example 4 Investigate exponential tilting when $f_0(y)$ is uniform on $(0, 2\pi]$ with $s(y) = (\cos y, \sin y)^{\mathrm{T}}$.

stat.epfl.ch

Here $\mathcal{Y}=(0,2\pi]$ is finite, and s(y) has dimension 2 and is bounded, so with $(\varphi_1,\varphi_2)\in\mathbb{R}^2$,

$$\int f_0(y) \exp \{ \varphi^{\mathsf{T}} s(y) \} \, dy = \frac{1}{2\pi} \int_0^{2\pi} \exp(\varphi_1 \cos y + \varphi_2 \sin y) \, dy$$
$$= \frac{1}{2\pi} \int_0^{2\pi} \exp\{\theta_2 \cos(y - \theta_1) \} \, dy = I_0(\theta_2),$$

where $\theta_2 = (\varphi_1^2 + \varphi_2^2)^{1/2} \ge 0$, $\theta_1 = \tan^{-1}(\theta_2/\theta_1) \in (0, 2\pi]$, and $I_0(\theta_2)$ is a modified Bessel function of the first kind and order 0. Hence $\varphi_1 = \theta_2 \cos \theta_1$ and $\varphi_2 = \theta_2 \sin \theta_1$. Hence

$$k(\varphi) = \log I_0\{(\varphi_1^2 + \varphi_2^2)^{1/2}\}, \quad \varphi \in \mathcal{N} = \mathbb{R}^2.$$

This is the von Mises–Fisher distribution on the circle, which concentrates around θ_1 , with the degree of concentration determined by $\theta_2 \geq 0$; $\theta_2 = 0$ gives the uniform density.

stat.epfl.ch

Autumn 2024 - note 1 of slide 20

Exponential family models

If $\theta \in \Theta \subset \mathbb{R}^d$, where $\dim \Theta = d$, and there exists a $d \times 1$ function s = s(y) of data y and a parametrisation (i.e., a 1–1 function) $\varphi \equiv \varphi(\theta)$ such that

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

then this is an (d, d) exponential family of distributions, with

- canonical statistic S = s(Y),
- canonical parameter φ ,
- **cumulant generator** k, which is convex on $\mathcal{N} = \{\varphi : k(\varphi) < \infty\}$, and
- mean parameter $\mu \equiv \mu(\varphi) = E(S; \varphi) = \nabla k(\varphi)$, where $\nabla \cdot = \partial \cdot / \partial \varphi$.
- We suppose that there is no vector a such that a^TS is constant, and call the model a **minimal** representation if there is no vector a such that $a^T\varphi$ is constant.
- \square The cumulant-generating function for S is

$$K_S(t) = \log M_S(t) = k(\varphi + t) - k(\varphi), \quad t \in \mathcal{N}' \subset \mathbb{R}^d$$

where $0 \in \mathcal{N}'$. On writing $\nabla^2 \cdot = \partial^2 \cdot /\partial \varphi \partial \varphi^T$, one can check that

$$E(S) = \nabla k(\varphi), \quad var(S) = \nabla^2 k(\varphi).$$

stat.epfl.ch

Note: Cumulant-generating functions

 $\ \square$ The MGF for the canonical statistic S of an exponential family is

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

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

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

and therefore that if it is defined,

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

which yields $K_S(t) = k(\varphi + t) - k(\varphi)$.

 \square Now $M_S(0)=1$, $K_S(0)=0$, $\partial K_S(t)/\partial t=\nabla k(\varphi+t)$ and $\partial^2 K_S(t)/\partial t\partial t^{\mathrm{T}}=\nabla^2 k(\varphi+t)$, so

$$\mathrm{E}(S) = \left. \partial M_S(t) / \partial t \right|_{t=0} = \left. \partial e^{K_S(t)} / \partial t \right|_{t=0} = \left. \partial K_S(t) / \partial t \, e^{K_S(t)} \right|_{t=0} = \nabla k(\varphi).$$

A similar calculation for the variance gives

$$E(SS^{\mathrm{T}}) = \partial^{2} M_{S}(t) / \partial t \partial t^{\mathrm{T}} \big|_{t=0} = \nabla^{2} k(\varphi) + \nabla k(\varphi) \nabla k(\varphi)^{\mathrm{T}},$$

and thus

$$\operatorname{var}(S) = \operatorname{E}(SS^{\mathrm{\scriptscriptstyle T}}) - \operatorname{E}(S)\operatorname{E}(S)^{\mathrm{\scriptscriptstyle T}} = \nabla^2 k(\varphi) + \nabla k(\varphi)\nabla k(\varphi)^{\mathrm{\scriptscriptstyle T}} - \nabla k(\varphi)\nabla k(\varphi)^{\mathrm{\scriptscriptstyle T}} = \nabla^2 k(\varphi).$$

stat.epfl.ch

Autumn 2024 - note 1 of slide 21

Examples

Example 5 (Poisson sample) If $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} \operatorname{Poiss}(\theta)$, find the corresponding exponential family.

Example 6 (Satellite conjunction) A simple model for the position Y of a satellite in \mathbb{R}^2 relative to the origin is

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

where $d_1, d_2 > 0$ are known and $\psi > 0$, $0 < \lambda \le 2\pi$. Write the corresponding density

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

as an exponential family.

 \square NB: avoid confusion — exponential family \neq exponential distribution! The exponential distribution is just one example of an exponential family.

stat.epfl.ch

Independent Poisson Y_1, \ldots, Y_n have joint density

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

where $m(y) = (\prod y_i)^{-1}$. This is a (1,1) exponential family with

- \square canonical statistic $s = s(y) = \sum y_j$,
- \square canonical parameter $\log \theta = \varphi \in \mathcal{N} = \mathbb{R}$,
- \Box cumulant generator $k(\varphi) = n\theta = ne^{\varphi}$ and
- \square mean parameter $\mu = \nabla k(\varphi) = ne^{\varphi} = n\theta = \mathrm{E}(S)$.

Two standard parametrizations use the real parameter φ or the mean $\mu = ne^{\varphi} \in \mathbb{R}_+$.

stat.epfl.ch

Autumn 2024 - note 1 of slide 22

Note to Example 6

☐ The multivariate normal density is

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

and if Ω is known then the exponent can be written as

$$-\frac{1}{2}\log\{(2\pi)^n|\Omega|\} - \frac{1}{2}y^{\mathrm{T}}\Omega^{-1}y + y^{\mathrm{T}}\Omega^{-1}\mu - \frac{1}{2}\mu^{\mathrm{T}}\Omega^{-1}\mu = \log m(y) + s(y)^{\mathrm{T}}\varphi - k(\varphi),$$

where $s(y)=\Omega^{-1}y$, $\varphi=\mu$ and $k(\varphi)=\frac{1}{2}\varphi^{\mathrm{T}}\Omega^{-1}\varphi$. It is easy to check that $\nabla k(\varphi)=\Omega^{-1}\varphi=\mathrm{E}(S)$ and $\nabla^2 k(\varphi)=\Omega^{-1}=\mathrm{var}(S)$.

 $\hfill\Box$ In the satellite example $d=2,\ \Omega=D^{-1}$ is diagonal, and with $\theta^{\scriptscriptstyle {\rm T}}=(\psi,\lambda)$ we have

$$\varphi^{\mathrm{T}} = (\varphi_1, \varphi_2) = (\psi \cos \lambda, \psi \sin \lambda), \quad s(Y) = (d_1 Y_1, d_2 Y_2), \quad k(\varphi) = d_1 \varphi_1^2 / 2 + d_2 \varphi_2^2 / 2.$$

The θ parametrisation gives the polar coordinates of the mean φ , but these are clearly equivalent because there is a 1–1 mapping between them.

stat.epfl.ch

Autumn 2024 - note 2 of slide 22

Exponential family models II

- When $\dim s = d' > \dim \theta = d$ the model is called a (d', d) curved exponential family, and the $d' \times 1$ vector $\varphi(\theta)$ gives a d-dimensional sub-manifold of $\mathbb{R}^{d'}$.
- \square Exponential families are closed under sampling: the joint density of independent observations Y_1,\ldots,Y_n from an exponential family with the same $s(Y_j)^{\mathrm{T}}\varphi=S_i^{\mathrm{T}}\varphi$ is

$$\prod_{j=1}^{n} f(y_j; \theta) = \prod_{j=1}^{n} m(y_j) \exp\left\{s_j^{\mathrm{T}} \varphi - k_j(\varphi)\right\} = \prod_{j=1}^{n} m(y_j) \exp\left\{\left(\sum_{j=1}^{n} s_j\right)^{\mathrm{T}} \varphi - \sum_{j=1}^{n} k_j(\varphi)\right\},$$

so with $k_S(\varphi) = \sum_j k_j(\varphi)$, the density of $S = \sum_j S_j = \sum_j s(Y_j)$ is

$$f(s;\theta) = m^*(s)e^{s^{\mathrm{T}}\varphi - k_S(\varphi)}, \quad \text{with} \quad m^*(s) = \int_{\{y: \sum_j s(y_j) = s\}} \prod_{j=1}^n m(y_j) \,\mathrm{d}y.$$

This is an exponential family, with canonical statistic S, canonical parameter φ and cumulant generator $k_S(\varphi)$.

Example 7 (Satellite conjunction) Show that taking ψ known in Example 6 gives a (2,1) exponential family.

stat.epfl.ch Autumn 2024 – slide 23

Note to Example 7

We previously had

$$\varphi^{\mathrm{T}} = (\varphi_1, \varphi_2) = (\psi \cos \lambda, \psi \sin \lambda), \quad s(Y) = (d_1 Y_1, d_2 Y_2), \quad k(\varphi) = d_1 \varphi_1^2 / 2 + d_2 \varphi_2^2 / 2,$$

but with ψ known we can write

$$\varphi^{\mathrm{T}} = (\varphi_1, \varphi_2) = (\cos \lambda, \sin \lambda), \quad s(Y) = (\psi d_1 Y_1, \psi d_2 Y_2), \quad k(\varphi) = \psi^2 (d_1 \varphi_1^2 + d_2 \varphi_2^2)/2,$$

where λ is the only unknown parameter. This is a (2,1) exponential family because it cannot be written in terms of a scalar φ ; the mean traces a curve (a circle) as λ varies.

stat.epfl.ch

Autumn 2024 - note 1 of slide 23

Inequalities

 \square A real-valued convex function g defined on a vector space $\mathcal V$ has the property that for any $x,y\in\mathcal V$,

$$g\{tx + (1-t)y\} \le tg(x) + (1-t)g(y), \quad 0 \le t \le 1.$$

Equivalently, for all $y \in \mathcal{V}$, there exists a vector b(y) such that

$$g(x) \ge g(y) + b(y)^{\mathrm{T}}(x - y)$$

for all x. If g(x) is differentiable, then we can take b(y) = g'(y).

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

$$P\{h(X) \ge a\} \le E\{h(X)\}/a$$
, (basic inequality)

$$P(|X| \ge a) \le E(|X|)/a$$
, (Markov's inequality)

$$\mathrm{P}(|X| \geq a) \leq \mathrm{E}(X^2)/a^2,$$
 (Chebyshov's inequality)

$$E\{g(X)\} \ge g\{E(X)\}.$$
 (Jensen's inequality)

 \square On replacing X by X - E(X), Chebyshov's inequality gives

$$P\{|X - E(X)| \ge a\} \le var(X)/a^2.$$

stat.epfl.ch

Autumn 2024 - slide 24

Note: Inequalities

(a) Let Y = h(X). If $y \ge 0$, then for any a > 0, $y \ge yI(y \ge a) \ge aI(y \ge a)$. Therefore

$$\mathrm{E}\{h(X)\} = \mathrm{E}(Y) \geq \mathrm{E}\{YI(Y \geq a)\} \geq \mathrm{E}\{aI(Y \geq a)\} = a\mathrm{P}(Y \geq a) = a\mathrm{P}\{h(X) \geq a\},$$

and division by a > 0 gives the result.

- (b) Note that h(x) = |x| is a non-negative function on \mathbb{R} , and apply (a).
- (c) Note that $h(x) = x^2$ is a non-negative function on \mathbb{R} , and that $P(X^2 \ge a^2) = P(|X| \ge a)$.
- (d) A convex function has the property that, for all y, there exists a value b(y) such that $g(x) \geq g(y) + b(y)(x-y)$ for all x. If g(x) is differentiable, then we can take b(y) = g'(y). (Draw a graph if need be.) To prove this result, we take $y = \mathrm{E}(X)$, and then have

$$g(X) \ge g\{E(X)\} + b\{E(X)\}\{X - E(X)\},\$$

and taking expectations of this gives $E\{g(X)\} \ge g\{E(X)\}.$

stat.epfl.ch

Autumn 2024 - note 1 of slide 24

Modes of convergence

- \square Let X, X_1, X_2, \ldots have CDFs F, F_1, F_2, \ldots and let $\varepsilon > 0$ be arbitrary. Then
 - X_n converges to X almost surely, $X_n \xrightarrow{\text{a.s.}} X$, if $P(\lim_{n\to\infty} X_n = X) = 1$;
 - X_n converges to X in probability, $X_n \stackrel{P}{\longrightarrow} X$, if $\lim_{n\to\infty} P(|X_n X| > \varepsilon) = 0$;
 - X_n converges to X in distribution, $X_n \xrightarrow{D} X$, if $\lim_{n\to\infty} F_n(x) = F(x)$ at each point x where F(x) is continuous.
 - A sequence X_1, X_2, \ldots of estimators of a parameter θ is strongly consistent if $X_n \xrightarrow{a.s.} \theta$ and (weakly) consistent if $X_n \xrightarrow{P} \theta$.
- $\square \quad \stackrel{\mathrm{a.s.}}{\longrightarrow} \ \, \text{and} \ \, \stackrel{P}{\longrightarrow} \text{ , but not } \ \, \stackrel{D}{\longrightarrow} \text{ , require joint distributions of } (X_n,X) \text{ for every } n.$
- \square Let x_0,y_0 be constants, $X,Y,\{X_n\},\{Y_n\}$ rbe andom variables and $g(\cdot)$ and $h(\cdot,\cdot)$ continuous functions. Then

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

The last two lines are called the **continuous mapping theorem** (usually used with $\stackrel{P}{\longrightarrow}$) and **Slutsky's theorem**.

stat.epfl.ch Autumn 2024 – slide 25

Limit theorems

Theorem 8 (Weak law of large numbers, WLLN) If $X, X_1, X_2, ... \stackrel{\text{iid}}{\sim} F$ and E(X) is finite, then $\overline{X} = n^{-1}(X_1 + \cdots + X_n) \stackrel{P}{\longrightarrow} E(X)$.

Theorem 9 (Strong law of large numbers, SLLN) If $X, X_1, X_2, ... \stackrel{\text{iid}}{\sim} F$ and E(X) is finite, then $\overline{X} = n^{-1}(X_1 + \cdots + X_n) \stackrel{\text{a.s.}}{\longrightarrow} E(X)$.

Theorem 10 (Central limit theorem, CLT) If $X_1, X_2, \ldots \stackrel{\mathrm{iid}}{\sim} (\mu, \sigma^2)$ and $0 < \sigma^2 < \infty$, then

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

Theorem 11 ('Delta method') If $a_n(X_n - \mu) \stackrel{D}{\longrightarrow} Y$, $a_n, \mu \in \mathbb{R}$, $a_n \to \infty$ as $n \to \infty$, and g is continuously differentiable at μ with $g'(\mu) \neq 0$, then $a_n\{g(X_n) - g(\mu)\} \stackrel{D}{\longrightarrow} g'(\mu)Y$.

- The CLT provides the finite-sample approximation $Z_n \stackrel{.}{\sim} \mathcal{N}(\mu, \sigma^2/n)$, where $\stackrel{.}{\sim}$ means 'is approximately distributed as'.
- ☐ Many more general laws of large numbers and versions of the CLT exist.
- \square The delta method also applies with $X_n,Z\in\mathbb{R}^p,\ g(x):\mathbb{R}^p\to\mathbb{R}^q$ continuously differentiable and $g'(\mu)$ replaced by $J_g(\mu)=\partial g(\mu)/\partial \mu^{\mathrm{T}}.$

Statistical activities

- ☐ Planning of investigations
- ☐ Obtaining reliable data
- ☐ Exploratory data analysis/visualisation
- □ Model formulation
- ☐ Point estimation of a population parameter
- ☐ Interval estimation for a population parameter
- ☐ Hypothesis testing to assess whether observed data support a particular model
- ☐ **Prediction** of a future or unobserved random variable
- Decision analysis to choose an action based on data and the costs of potential actions

This course covers some aspects of those activities in red above.

Many inferential tasks can be formulated in decision-theoretic terms, but we shall mostly avoid this.

stat.epfl.ch

Autumn 2024 - slide 28

Statistical models

- Use observed data to draw conclusions about a 'population', i.e., a model from which the data are assumed to be drawn, or about future data.
- \square A statistical model is a family of probability distributions for data y in a sample space \mathcal{Y} .
- \square A parametric model (family of models) $f \equiv f(y; \theta)$ or equivalently $F \equiv f(y; \theta)$ is determined by parameters $\theta \in \Theta \subset \mathbb{R}^d$, for fixed finite d.
- \square If no such θ exists, F is nonparametric, and then the parameter is often determined by F through a statistical functional $\theta = t(F)$, e.g.,

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

- ☐ Parameters have different roles (which can change during an investigation):
 - interest parameters represent targets of inference (e.g., the mean of a population, the slope of a line, a baseline blood pressure) with direct substantive interpretations;
 - nuisance parameters are needed to complete a model specification, but are not themselves of main concern.
- \square A parametric model should have a 1–1 map from θ to $f(\cdot;\theta)$, so parameters identify models.

stat.epfl.ch

М	odel formulation
	Two broad types of statistical model:
	 substantive — based on fundamental subject-matter theory (e.g., quantum theory, Mendelian genetics, Navier–Stokes equations);
	 empirical — a convenient, adequately realistic, representation of data variation;
	 and of course a broad spectrum between them.
	We aim that
	 primary questions/issues are encapsulated in the interest parameter;
	 secondary aspects can be accounted for, often via nuisance parameters;
	 variation in the data is realistically modelled, leading to reasonable statements of uncertainty;
	 any special feature of the data or 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 2024 – slide 30

Point estimation

- An estimator of a parameter $\theta \in \Theta$ based on data Y is a random variable $\tilde{\theta} = \tilde{\theta}(Y)$ taking values in Θ . A specific value is an estimate $\tilde{\theta}(y)$.
- \square An M(aximisation)-estimator is computed using a function $\rho(y;\theta')$ as

$$\tilde{\theta} = \operatorname{argmax}_{\theta'} \frac{1}{n} \sum_{j=1}^{n} \rho(Y_j; \theta').$$

Often $\tilde{\theta}$ also solves

$$\frac{1}{n}\sum_{j=1}^{n}\nabla\rho(Y_j;\theta')=0$$

and is then called a **Z**(ero)-estimator.

- \Box Equivalently we could minimise the loss function $-\rho$ with respect to θ .
- \square If the true underlying model is g, then $\widetilde{\theta}$ is replaced by θ_g , where

$$\theta_g = \operatorname{argmax}_{\theta'} \int \rho(y; \theta') g(y) \, dy, \quad \int \nabla \rho(y; \theta_g) g(y) \, dy = 0.$$

Clearly if $g(y) = f(y; \theta)$, then we want $\theta_g = \theta$, uniquely.

stat.epfl.ch Autumn 2024 – slide 32

Examples

- \square Some examples (for a *d*-dimensional parameter θ):
 - maximum likelihood estimation has $\rho(y; \theta') = \log f(y; \theta')$;
 - method of moments estimation has $h(y) = (y, y^2, \dots, y^d)^T$, $\mu(\theta') = E\{h(Y)\}$, and

$$-\rho(y; \theta') = \{h(y) - \mu(\theta')\}^{\mathrm{T}} \{h(y) - \mu(\theta')\};$$

– generalized method of moments estimation (widely used in econometrics) also has a symmetric positive definite $d \times d$ matrix $w(\theta')$ and

$$-\rho(y; \theta') = \{h(y) - \mu(\theta')\}^{\mathrm{T}} w(\theta') \{h(y) - \mu(\theta')\};$$

- least squares estimation is method of moments estimation with $h(y_j) = y_j$ and $\mu_j(\theta') = \mathrm{E}(Y_j) = x_j^\mathrm{T} \theta';$
- score-matching estimation (unfortunate misnomer) with $Y \sim g$ has

$$-\rho(y; \theta') = \left\{ \nabla_y \log f(y; \theta) - \nabla_y \log g(y) \right\}^2.$$

☐ There are many (many!) other approaches to estimation.

Examples

Example 12 Discuss maximum likelihood estimation of the parameters of the normal distribution.

Example 13 Discuss moment estimation of the parameters of the Weibull distribution.

Example 14 Show that under mild (but not entirely trivial) conditions on the density g, the population version of the score-matching estimator is

$$\mathrm{argmin}_{\theta} \mathbf{E} \left[\left\{ \nabla_y \log f(Y;\theta) \right\}^2 + 2 \nabla_y^2 \log f(Y;\theta) \right],$$

and give the sample version.

The density function of a normal random variable with mean μ and variance σ^2 is $(2\pi\sigma^2)^{-1/2}\exp\{-(y-\mu)^2/(2\sigma^2)\}$, so here $\theta_{2\times 1}=(\mu,\sigma^2)^{\rm T}\in\mathbb{R}\times\mathbb{R}_+$, and the likelihood for a random sample y_1,\ldots,y_n equals

$$L(\theta) = f(y; \theta) = \prod_{j=1}^{n} f(y_j; \theta) = \prod_{j=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_j - \mu)^2}{2\sigma^2}\right\}.$$

Therefore the log likelihood is

$$\ell(\mu, \sigma) = -\frac{n}{2}\log(2\pi) - \frac{n}{2}\log\sigma^2 - \frac{1}{2\sigma^2}\sum_{j=1}^n (y_j - \mu)^2, \quad \mu \in \mathbb{R}, \sigma^2 > 0.$$

Its first derivatives are

$$\frac{\partial \ell}{\partial \mu} = \sigma^{-2} \sum_{j=1}^{n} (y_j - \mu), \quad \frac{\partial \ell}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{j=1}^{n} (y_j - \mu)^2,$$

and its (negative) second derivatives are

$$\frac{\partial^2 \ell}{\partial \mu^2} = -\frac{n}{\sigma^2}, \quad \frac{\partial^2 \ell}{\partial \mu \partial \sigma^2} = -\frac{n}{\sigma^4} (\overline{y} - \mu), \quad \frac{\partial^2 \ell}{\partial (\sigma^2)^2} = \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \sum_{j=1}^n (y_j - \mu)^2.$$

☐ To obtain the MLEs, we solve simultaneously the equations

$$\begin{pmatrix} \frac{\partial \ell(\mu, \sigma^2)}{\partial \mu} \\ \frac{\partial \ell(\mu, \sigma^2)}{\partial \sigma^2} \end{pmatrix} = \begin{pmatrix} \sigma^{-2} \sum_{j=1}^n (y_j - \mu) \\ -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{j=1}^n (y_j - \mu)^2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

Now

$$\frac{\partial \ell(\widehat{\mu}, \widehat{\sigma}^2)}{\partial \mu} = 0 \Rightarrow \frac{1}{\widehat{\sigma}^2} \sum_{j=1}^n (y_j - \widehat{\mu}) = 0 \Rightarrow n\widehat{\mu} = \sum_{j=1}^n y_j \Rightarrow \widehat{\mu} = n^{-1} \sum_{j=1}^n y_j = \overline{y}$$

and

$$\frac{\partial \ell(\widehat{\mu}, \widehat{\sigma}^2)}{\partial \sigma^2} = 0 \Rightarrow \frac{n}{2\widehat{\sigma}^2} = \frac{1}{2\widehat{\sigma}^4} \sum_{j=1}^n (y_j - \widehat{\mu})^2 \Rightarrow \widehat{\sigma}^2 = n^{-1} \sum_{j=1}^n (y_j - \widehat{\mu})^2 = n^{-1} \sum_{j=1}^n (y_j - \overline{y})^2.$$

The first of these has the sole solution $\widehat{\mu}=\overline{y}$ for all values of σ^2 , and therefore $\ell(\widehat{\mu},\sigma^2)$ is unimodal with maximum at $\widehat{\sigma}^2=n^{-1}\sum(y_j-\overline{y})^2$. At the point $(\widehat{\mu},\widehat{\sigma}^2)$, the hessian matrix is diagonal with elements $\mathrm{diag}\{n/\widehat{\sigma}^2,n/(2\widehat{\sigma}^4)\}$, and so is positive definite. Hence $\widehat{\mu}=\overline{y}$ and $\widehat{\sigma}^2=n^{-1}\sum(y_j-\overline{y})^2$ are the sole solutions to the likelihood equation, and therefore are the maximum likelihood estimates.

stat.epfl.ch

Autumn 2024 - note 1 of slide 34

 \square A Weibull variable X has CDF $F(x)=1-e^{-(\lambda x)^{\alpha}}$, for x>0 and $\lambda,\alpha>0$, and is exponential when $\alpha=1$. Note that $W=(\lambda X)^{\alpha}\sim \exp(1)$, so

$$\mathrm{E}(X^r) = \mathrm{E}\{(W^{1/\alpha}/\lambda)^r\} = \lambda^{-r}\mathrm{E}(W^{r/\alpha}) = \lambda^{-r}\int_0^\infty w^{r/\alpha}e^{-w}\,\mathrm{d}w = \lambda^{-r}\Gamma(1+r/\alpha),$$

where $\Gamma(\cdot)$ is the gamma function. Hence with $\theta = (\lambda, \alpha)$ the moment estimators solve

$$\overline{Y} = \mu_1(\theta) = \lambda^{-1} \Gamma(1 + 1/\alpha), \quad \overline{Y^2} = \mu_2(\theta) = \lambda^{-2} \Gamma(1 + 2/\alpha), \quad \lambda, \alpha > 0,$$

i.e.,

$$\overline{Y^2}/(\overline{Y})^2 = \Gamma(1+2/\tilde{\alpha})/\Gamma(1+1/\tilde{\alpha})^2, \quad \tilde{\lambda} = \Gamma(1+1/\tilde{\alpha})/\overline{Y}.$$

stat.epfl.ch

Autumn 2024 - note 2 of slide 34

Note to Example 14

- \square Score-matching can be useful when $\log f(y;\theta) = h(y;\theta) k(\theta)$ with $k(\theta)$ intractable. It is a misnomer because the standard use of the term 'score' is for the derivative of the log likelihood with respect to θ (not y).
- \square On writing $\log f(y;\theta) = \ell(\theta)$ for brevity we can write

$$\left\{\nabla_y \log f(y;\theta) - \nabla_y \log g(y)\right\}^2 = \left\{\nabla_y \ell(\theta)\right\}^2 - 2\nabla_y \ell(\theta)\nabla_y \log g(y) + \left\{\nabla_y \log g(y)\right\}^2,$$

we see that the population version of the estimator is

$$\theta_g = \operatorname{argmin}_{\theta} \int \{\nabla_y \ell(\theta)\}^2 g(y) \, dy - 2 \int \{\nabla_y \ell(\theta) \nabla_y \log g(y)\} g(y) \, dy,$$

because θ does not appear in the third term of the square. As g is unknown, the second integral here appears intractable, but as $g(y)\nabla_y\log g(y)=\nabla_y g(y)$, we have

$$\int \nabla_y \ell(\theta) \nabla_y \log g(y) g(y) dy = \int \nabla_y \ell(\theta) \nabla_y g(y) dy$$

and integration by parts gives

$$\int \nabla_y \ell(\theta) \nabla_y g(y) \, dy = \left[\nabla_y \ell(\theta) g(y) \right] - \int \nabla_y^2 \ell(\theta) g(y) \, dy$$
$$= -E \left\{ \nabla_y^2 \log f(Y; \theta) \right\},$$

when the first integration term is identically zero. Hence

$$\theta_g = \operatorname{argmin}_{\theta} \mathbb{E} \left[\left\{ \nabla_y \log f(Y; \theta) \right\}^2 + 2 \nabla_y^2 \log f(Y; \theta) \right],$$

whose sample version,

$$\tilde{\theta} = \operatorname{argmin}_{\theta} \sum_{j=1}^{n} \left[\left\{ \nabla_{y} \log f(Y_{j}; \theta) \right\}^{2} + 2 \nabla_{y}^{2} \log f(Y_{j}; \theta) \right],$$

can be computed from the sample.

☐ Weighted versions can be used to kill the first term of the integral, when it is non-zero (exercise).

stat.epfl.ch

Autumn 2024 - note 3 of slide 34

Comparison of point estimators

- ☐ There are two generic bases for comparing point estimators:
 - asymptotic what happens when $n \to \infty$?
 - finite-sample what happens for sample sizes met in practice?
- Consistency is a key asymptotic criterion: does $\tilde{\theta}$ approach θ_g when $n \to \infty$?

Definition 15 An estimator $\tilde{\theta}$ of θ_q is (weakly) consistent if $\tilde{\theta} \stackrel{P}{\longrightarrow} \theta_q$ as $n \to \infty$.

☐ Consistency is necessary but not sufficient for an estimator to be good, because

$$\tilde{\theta} \xrightarrow{P} \theta_q \Rightarrow \tilde{\theta}^* = \tilde{\theta} + 10^6 / \sqrt{\log \log n} \xrightarrow{P} \theta_q, \quad n \to \infty,$$

but $\tilde{\theta}^*$ is (probably) useless: consistency can be considered a 'safety net'.

 \square Obviously we would like $ilde{ heta}$ to be 'suitably close' to $heta_g$, by minimising

$$\mathrm{MSE}(\tilde{\theta}; \theta_g) = \mathrm{E}\left\{ (\tilde{\theta} - \theta_g)^2 \right\}, \quad \mathrm{MAD}(\tilde{\theta}; \theta_g) = \mathrm{E}\left(|\tilde{\theta} - \theta_g| \right),$$

or other measures of distance (loss functions), asymptotically or in finite samples.

stat.epfl.ch Autumn 2024 – slide 35

Bias-variance and other tradeoffs

 \square Using the bias $b(\tilde{\theta}; \theta_g) = E(\tilde{\theta}) - \theta_g$, the mean square error can be expressed as

$$MSE(\tilde{\theta}; \theta_q) = b(\tilde{\theta}; \theta_q)^2 + var(\tilde{\theta}),$$

so we must balance ('trade off') the bias and the variance when choosing $\tilde{\theta}$.

- \square In simple problems we could insist that the estimator is **unbiased**, i.e., $b(\tilde{\theta}; \theta_g) \equiv 0$, but this is usually artificial because
 - many good estimators are biased, and some unbiased estimators are useless;
 - it may be impossible to find an unbiased estimator; and
 - other properties may be more desirable (e.g., robustness).

An exception is meta-analysis, which involves combining different estimators with possibly very varied sample sizes, in which case we want them to estimate the same thing!

Example 16 The method of moments estimator of a scalar θ based on a random sample $Y_1,\ldots,Y_n\stackrel{\mathrm{iid}}{\sim} (\mu,\sigma^2)$ with sample average \overline{Y} solves the equation $\mu(\theta)=\overline{Y}$. Show that if $\mu(\cdot)$ has two smooth derivatives and is 1–1, then the estimator is consistent and asymptotically normal, with bias and variance both of order n^{-1} .

- As the function $\mu(\cdot)$ is smooth and 1–1, it has a differentiable inverse, and thus by the continuous mapping theorem, $\tilde{\theta} = \mu^{-1}(\overline{Y}) \stackrel{P}{\longrightarrow} \mu^{-1}\{\mu(\theta)\} = \theta$, i.e., $\tilde{\theta}$ is consistent. For simplicity of notation write $g(x) = \mu^{-1}(x)$ below.
- \square Now $\overline{Y} = \mu + \sigma n^{-1/2} Z_n$, where $Z_n = (\overline{Y} \mu)/(\sigma^2/n)^{1/2} \stackrel{D}{\longrightarrow} Z \sim \mathcal{N}(0,1)$, and we have

$$g(\overline{Y}) = g(\mu) + g'(\mu)\sigma n^{-1/2}Z_n + \frac{\sigma^2}{2}n^{-1}g''(\mu + \sigma n^{-1/2}Z_n')Z_n^2,$$

where $Z_n' \in (0, Z_n)$, i.e.,

$$\tilde{\theta} = \theta + n^{-1/2} \sigma g'(\mu) Z_n + n^{-1} A_n,$$

say, where A_n is a random variable of order 1. Taking expectations gives

$$b(\tilde{\theta}; \theta) = E(\tilde{\theta}) - \theta = n^{-1}E(A_n) = O(n^{-1}),$$

under mild further conditions on g''.

□ Now

$$n^{1/2}(\tilde{\theta} - \theta)/\{\sigma g'(\mu)\} = Z_n + n^{-1/2}A'_n \xrightarrow{D} Z,$$

using this (or the delta method), so in large samples we have

$$\tilde{\theta} \stackrel{\cdot}{\sim} \mathcal{N}\{\theta, \sigma^2 g'(\mu)^2/n\}.$$

stat.epfl.ch

Autumn 2024 - note 1 of slide 36

Efficiency and the Cramèr-Rao lower bound

Definition 17 If $\tilde{\theta}_1$ and $\tilde{\theta}_2$ are estimators of scalar θ , then the relative efficiency of $\tilde{\theta}_1$ compared to $\tilde{\theta}_2$ can be defined as

$$\frac{\mathrm{MSE}(\tilde{\theta}_2; \theta)}{\mathrm{MSE}(\tilde{\theta}_1; \theta)}.$$

In large samples the squared bias is often negligible compared to the variance, and we define the asymptotic relative efficiency as $var(\tilde{\theta}_2)/var(\tilde{\theta}_1)$. Similar expressions apply if the parameter has dimension d.

Under mild conditions on the underlying model, a scalar estimator $\tilde{\theta}$ based on $Y \sim f(y; \theta)$ satisfies the Cramèr–Rao lower bound,

$$\operatorname{var}(\tilde{\theta}) \ge \frac{\{1 + \nabla b(\tilde{\theta}; \theta)\}^2}{i(\theta)},$$

where $i(\theta)$ is defined on the next slide. This bound applies for any sample size n. Moreover

- as $n \to \infty$ the lower bound $\to 1/i(\theta)$, the asymptotic variance of the maximum likelihood estimator, which hence is most efficient in large samples; and
- a similar result applies for vector θ .

stat.epfl.ch

Bartlett identities

- □ For data $Y \sim f(y; \theta)$ we define the log likelihood function $\ell(\theta) = \log f(Y; \theta)$ and $d \times 1$ score vector $U(\theta) = \nabla \ell(\theta)$.
- \Box If we can differentiate with respect to θ under the integral sign, we get the Bartlett identities:

$$0 = \int \nabla \log f(y; \theta) \times f(y; \theta) \, \mathrm{d}y,$$

$$0 = \int \nabla^2 \log f(y; \theta) \times f(y; \theta) \, \mathrm{d}y + \int \nabla \log f(y; \theta) \, \nabla^{\mathrm{\scriptscriptstyle T}} \log f(y; \theta) \times f(y; \theta) \, \mathrm{d}y,$$

giving the moments of $U(\theta)$, viz

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

where $var\{U(\theta)\} = i(\theta)$ is the $d \times d$ Fisher (or expected) information matrix.

- \square We write $i_1(\theta)$ for the Fisher information for a single observation of a random sample Y_1, \ldots, Y_n , and then that in the sample is $i(\theta) = ni_1(\theta)$.
- \Box Later we shall see that in large samples, the maximum likelihood estimator $\widehat{\theta}$ satisfies

$$\widehat{\theta} \stackrel{\cdot}{\sim} \mathcal{N}_d \left\{ \theta, \imath(\theta)^{-1} \right\}.$$

stat.epfl.ch Autumn 2024 – slide 38

Note: Bartlett identities

 \Box For any θ we have $1 = \int f(y; \theta) \, \mathrm{d}y$, so provided we can exchange the order of integration and differentiation we have

$$0 = \nabla \int f(y; \theta) \, dy = \int \nabla f(y; \theta) \, dy = \int \nabla f(y; \theta) \frac{f(y; \theta)}{f(y; \theta)} \, dy = \int \nabla \log f(y; \theta) \, f(y; \theta) \, dy.$$

- \Box 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 $U(\theta)$.
- \square For independent data Y_1,\ldots,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

$$\operatorname{var}\{U(\theta)\} = \sum_{j=1}^{n} \operatorname{var}\{U_{j}(\theta)\} = \sum_{j=1}^{n} \operatorname{E}\{U_{j}(\theta)U_{j}^{\mathrm{T}}(\theta)\} = \sum_{j=1}^{n} -\operatorname{E}\{\nabla^{\mathrm{T}}U_{j}(\theta)\} = -\operatorname{E}\{\nabla^{\mathrm{T}}U(\theta)\}$$

and this equals $\mathrm{E}\left\{-\nabla^2\ell(\theta)\right\}=\imath(\theta)$, and this in turn equals $n\imath_1(\theta)$.

stat.epfl.ch

Autumn 2024 - note 1 of slide 38

Note: CRLB

□ We have

$$E(\tilde{\theta}) = \int \tilde{\theta}(y) f(y; \theta) dy = \theta + b(\tilde{\theta}; \theta),$$

and differentiation with respect to θ gives (setting $b'(\theta) = db(\tilde{\theta}; \theta)/d\theta$)

$$1 + b'(\theta) = \int \tilde{\theta}(y) df(y; \theta) / d\theta dy = \int \tilde{\theta}(y) \nabla \ell(\theta) f(y; \theta) dy = E\{\tilde{\theta}U(\theta)\} = \cos\{\tilde{\theta}, U(\theta)\},$$

because $U(\theta)$ has mean zero. Hence the definition of correlation gives

$$\operatorname{cov}\{\tilde{\theta}, U(\theta)\}^2 = \{1 + b'(\theta)\}^2 \le \operatorname{var}(\tilde{\theta})\operatorname{var}\{U(\theta)\} = \operatorname{var}(\tilde{\theta})\iota(\theta),$$

which gives the result.

If the bias is of order n^{-1} , so too is its derivative, so in large samples we obtain

$$\operatorname{var}(\tilde{\theta}) \ge i(\theta)^{-1} = \operatorname{var}(\hat{\theta}).$$

stat.epfl.ch

Autumn 2024 - note 2 of slide 38

Pivots

- \square Point estimation does not express uncertainty we need to assess how well the observed data y^o support different possible values of a parameter.
- □ We aim to find subsets of the parameter space that contain the 'true' parameter with a specified probability when the parameter of interest is scalar, these subsets are usually intervals.
- ☐ Pivots are useful in finding such subsets.

Definition 18 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., one that does not depend on θ).

Example 19 If $M = \max(Y_1, \dots, Y_n)$, where $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} U(0, \theta)$, show that $Q_1 = M/\theta$ is a pivot and find a pivot based on \overline{Y} .

stat.epfl.ch

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

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

SO

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

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

□ If $Y \sim U(0,\theta)$, then $E(Y) = \theta/2$ and $var(Y) = \theta^2/12$. Hence \overline{Y} has mean $\theta/2$ and variance $\theta^2/(12n)$, and for large n, $\overline{Y} \sim \mathcal{N}\{\theta/2, \theta^2/(12n)\}$ using the central limit theorem. Therefore

$$Q_2 = \frac{\overline{Y} - \theta/2}{\sqrt{\theta^2/(12n)}} = (3n)^{1/2} (2\overline{Y}/\theta - 1) \stackrel{.}{\sim} \mathcal{N}(0, 1).$$

Thus Q_2 depends on both data and θ , and has an (approximately) known distribution: hence Q_2 is an (approximate) pivot.

 \square As $Y/\theta \sim U(0,1)$, we see that we could use simulation to compute the exact distribution of Q_2 , and thus obtain an exact pivot (apart from simulation error). This is called a bootstrap calculation, about which more later.

stat.epfl.ch

Autumn 2024 - note 1 of slide 39

Confidence intervals

Definition 20 Let $Y = (Y_1, \ldots, Y_n)$ 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.

- \square L = l(Y) and U = u(Y) are computed from the data. They do not depend on θ .
- \square In a continuous setting (so < gives the same probabilities as \le), 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

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

- 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.
- \square We often take standard values of α , such that $1 \alpha = 0.9, 0.95, 0.99, \dots$
- \square A weaker requirement is $P(L \le \theta \le U) \ge 1 \alpha$, giving confidence level at least 1α .

stat.epfl.ch

Construction of a CI

- \square 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

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

into the form

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

where the bounds $L = l(Y; \alpha_L, \alpha_U)$, $U = u(Y; \alpha_L, \alpha_U)$ do not depend on θ ;

- then we replace Y by its observed value y° to get a realisation of the CI.
- \square Going from quantiles of Q to L,U is known as inverting the pivot it is convenient if Q is monotone in θ for each Y.
- Often we have an approximate pivot $(\widehat{\theta} \theta)/V^{1/2} \sim \mathcal{N}(0,1)$, where V estimates $\mathrm{var}(\widehat{\theta})$ and $V^{1/2}$ is called a standard error. The resulting (approximate) 95% interval is $\widehat{\theta} \pm 1.96V^{1/2}$.

Example 21 In Example 19, find CIs based on Q_1 and on Q_2 .

stat.epfl.ch Autumn 2024 – slide 41

Note to Example 21

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

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

and a little algebra gives that

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

SO

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

 \qed For $Q_2=(3n)^{1/2}(2\overline{Y}/\theta-1)\stackrel{.}{\sim}\mathcal{N}(0,1)$, the quantiles are $z_{1-\alpha_L}$ and z_{α_U} , so

$$P\{z_{\alpha_U} \le (3n)^{1/2} (2\overline{Y}/\theta - 1) \le z_{1-\alpha_L}\} = 1 - \alpha_L - \alpha_U,$$

and hence we obtain

$$L = \frac{2\overline{Y}}{1 + z_{1-\alpha_L}/(3n)^{1/2}}, \quad U = \frac{2\overline{Y}}{1 + z_{\alpha_U}/(3n)^{1/2}};$$

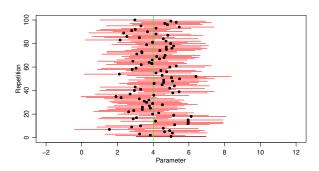
note that for large n these are $L\approx 2\overline{Y}\{1-z_{1-\alpha_L}/(3n)^{1/2}\}$ and $U\approx 2\overline{Y}\{1-z_{\alpha_U}/(3n)^{1/2}\}$.

stat.epfl.ch

Autumn 2024 - note 1 of slide 41

Interpretation of a CI

- \square (L,U) is a random interval that contains θ with probability $1-\alpha$.
- \square We imagine an infinity of possible datasets from the experiment that resulted in (L, U).
- \Box Our CI based on y° is regarded as randomly chosen from the resulting infinity of CIs.
- \square Although we do not know if $\theta \in (l(y^{\mathrm{o}}; \alpha_L, \alpha_U), u(y^{\mathrm{o}}; \alpha_L, \alpha_U))$, the event $\theta \in (L, U)$ has probability 1α 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 2024 – slide 42

More about CIs

- Almost invariably CIs are two-sided and equi-tailed, i.e., $\alpha_L = \alpha_U = \alpha$, but one-sided CIs of form $(-\infty, U)$ or (L, ∞) are sometimes required:
 - compute a two-sided interval with $\alpha_L = \alpha_U = \alpha$, then replace the unwanted limit by $\pm \infty$ (or another value if required in the context).
- ☐ For a two-sided CI we define the lower- and upper-tail errors

$$P(\theta < L), P(U < \theta)$$

and if these equal the required value for each possible α_L, α_U , then the **empirical coverage** of the CI exactly equals the desired value:

- this occurs when the distribution of the corresponding pivot is known, but in practice this
 distribution is usually approximate, and then we use simulation to assess if and when CIs are
 adequate;
- it's better to consider the two errors separately, as their sum may be OK even when they are individually incorrect;
- these errors are properties of the CI procedure, not of individual intervals!

Prediction

- \square Prediction refers to 'estimation' of unobserved (future, latent, ...) random variables Y_+ .
- In parametric cases we often base prediction (or tolerance) intervals on existing data Y by finding a pivot that depends on both Y_+ and Y, and predicting Y_+ using this pivot, e.g., using its mean or median.

Example 22 If $Y_1, \ldots, Y_n, Y_+ \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, give prediction limits and a predictor for Y_+ based on the other variables.

Example 23 (Conformal prediction) Suppose we seek a prediction interval for the outcome of an ML algorithm. In the simplest case, with Y_1, \ldots, Y_n, Y_+ real-valued and exchangeable, $\beta \in (0,1)$, $m = \lceil (n+1)\beta \rceil$ and q_β equal to the mth order statistic of Y_1, \ldots, Y_n , show that

$$P(Y_+ \le q_\beta) \ge \beta$$
,

and deduce that $P(q_{\alpha} < Y_{+} \le q_{1-\alpha}) \ge 1 - 2\alpha$.

stat.epfl.ch Autumn 2024 – slide 44

Note to Example 22

$$Q = \frac{Y_{+} - \overline{Y}}{\{(1 + 1/n)S^{2}\}^{1/2}} \sim t_{n-1},$$

leading to two-sided equi-tailed $(1-2\alpha)$ prediction interval

$$\overline{Y} \pm (1 + 1/n)^{1/2} St_{n-1} (1 - \alpha).$$

Note that even as $n \to \infty$ this interval does not vanish, rather it approaches $\mu \pm \sigma z_{1-\alpha}$.

- \square The Y_j are replaced by y_j^{o} to give the realisation of the interval.
- $\hfill\Box$ One obvious scalar predictor \widehat{Y}_+ is given by taking the median for Q, i.e., solving

$$q_{0.5} = \frac{\widehat{Y}_{+} - \overline{Y}}{\{(1+1/n)S^2\}^{1/2}},$$

where in this case $q_{0.5}=0$, giving $\widehat{Y}_+=\overline{Y}$ and realised value $\overline{y}^{\mathrm{o}}.$

stat.epfl.ch

Autumn 2024 - note 1 of slide 44

Let q_{β}^+ denote the mth order statistic of $\mathcal{Y}_+ = \{Y_1, \dots, Y_n, Y_+\}$, and note that under exchangeability Y_+ equals any of the order statistics of \mathcal{Y}_+ with probability 1/(n+1). Therefore

$$P(Y_{+} \le q_{\beta}^{+}) = m/(n+1) = \lceil (n+1)\beta \rceil/(n+1) \ge (n+1)\beta/(n+1) = \beta.$$

 \square Now suppose that m=2 and $Y_+ \leq q^+_\beta$, so using an obvious notation $\mathcal Y$ can be represented as

$$\bullet \le + \le \bullet \le \cdots$$
 or $+ \le \bullet \le \bullet \le \cdots$.

In both cases $q_{\beta} \geq q_{\beta}^+$, so $Y_+ \leq q_{\beta}^+$ implies that $Y_+ \leq q_{\beta}$, and conversely. This holds for any m, so

$$P(Y_{+} \le q_{\beta}) = P(Y_{+} \le q_{\beta}^{+}) \ge \beta.$$

Finally

$$P(q_{\alpha} < Y_{+} \le q_{1-\alpha}) = P(Y_{+} \le q_{1-\alpha}) - P(Y^{+} \le q_{\alpha}) \ge 1 - \alpha - \alpha = 1 - 2\alpha,$$

as required.

- \square For this argument to be practical we must have $1 \le m \le n$, so if β is too small or too large, then we must replace the corresponding limit by $\pm \infty$, which does not usually give a useful interval.
- In applications the data are of form (X,Y) and we train a prediction algorithm \widehat{f} using a training subset of $\mathcal{Y} = \{(X_1,Y_1),\ldots,(X_n,Y_n)\}$, giving residuals $Y_j \widehat{f}(X_j)$ for a test subset of \mathcal{Y} disjoint from the training set, and then apply the argument above to these residuals and $Y_+ \widehat{f}(X_+)$.

stat.epfl.ch

Autumn 2024 - note 2 of slide 44

Hypothesis testing

- A **statistical hypothesis** is an assertion about the population underlying some data, or equivalently a restriction on possible models for the data, such as:
 - the population has mean μ_0 ;
 - the population is $\mathcal{N}(\mu_0, \sigma_0^2)$, with both parameters specified;
 - the population is $\mathcal{N}(\mu, \sigma^2)$, with the parameters unspecified;
 - the data are sampled from the discrete uniform distribution on $\{1, \ldots, 9\}$;
 - the population density is symmetric about some μ ;
 - the population mean $\mu(x)$ increases when a covariate x increases.
- ☐ These are assertions about populations, not about data, but they have implications for data.
- ☐ Sometimes the distribution is fully specified, but not always.
- ☐ Some, but not all, hypotheses concern parameters.
- ☐ A hypothesis test uses a stochastic 'argument by contradiction' to make an inference about a statistical hypothesis: we assume that the hypothesis is true, and attempt to use our data to disprove it.

stat.epfl.ch

Elements of a test		
	A null hypothesis H_0 to be tested.	
	A test statistic T , large values of which suggest that H_0 is false, and with observed value $t_{\rm obs}$.	
	A P-value	
	$p_{\mathrm{obs}} = P_0(T \ge t_{\mathrm{obs}}),$	
	where the null distribution $P_0(\cdot)$ denotes a probability computed under H_0 .	
	The smaller $p_{ m obs}$ is, the more we doubt that H_0 is true.	
	Tests on parameters are often based on pivots: if $\theta=\theta_0$, then $T= q(Y;\theta_0) $ has a known distribution G_0 , say, and observing a value $t_{\rm obs}= q(y^{\rm o};\theta_0) $ that is unusual relative to G_0 'contradicts' H_0 .	
	In other cases we choose a test statistic that seems plausible, such as Pearson's statistic,	
	$T = \sum_{k=1}^{K} (O_k - E_k)^2 / E_k,$	
	used to check whether observed counts O_k in K categories agree with their expectations $E_k = \mathrm{E}(O_k)$ computed under H_0 .	

 \square In any case we need to know (or be able to approximate) the distribution of T under H_0 . Autumn 2024 - slide 46

1.4 Bases for Uncertainty

slide 47

Uncertainty

stat.epfl.ch

- ☐ Essentially three bases for statements of uncertainty:
 - a frequentist (sampling theory) inference compares y with a set $S \subset \mathcal{Y}$ of other data that might have been observed in a hypothetical sampling experiment;
 - a Bayesian (inverse probability) inference expresses uncertainty via a prior probability density and uses Bayes' theorem to update this in light of the data;
 - 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.
- \square A frequentist should choose the reference set (aka recognisable subset) $\mathcal S$ of the sample space \mathcal{Y} thoughtfully.

Example 24 (Measuring machines) A physical quantity θ can be measured with two machines, both giving normal observations $Y \sim \mathcal{N}(\theta, \sigma_m^2)$. A measurement from machine 1 has variance $\sigma_1^2 = 1$, and one from machine 2 has variance $\sigma_2^2=100$. A machine is chosen by tossing a fair coin, giving M=1,2 with equal probabilities. Thus $\mathcal{Y}=\{(y,m):y\in\mathbb{R},m\in\{1,2\}\}.$

If we observe (y, m) = (0, 1), then clearly we can ignore the fact that we might have observed m = 2, i.e., we should take $S_1=\{(y,1):y\in\mathbb{R}\}$ rather than $S_2=\{(y,2):y\in\mathbb{R}\}$ or $S=\mathcal{Y}.$

Autumn 2024 - slide 48 stat.epfl.ch

Comments on sampling theory inference

- \square 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 performed with respect to \mathcal{S} .
- \square We choose $\mathcal S$ to ensure that the probability calculation is **relevant** to the data actually observed. For example, if y^{o} 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

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

for every $\theta \in \Theta$ and every $\alpha \in (0,1)$. Hence if such intervals are used infinitely often, then

- although any particular interval either does or does not contain θ ,
- it was drawn from a population of intervals with error probability exactly α .
- \square Bayesians object that inferences should only be based on the dataset y° <u>actually observed</u>, so the reference set S is irrelevant.

Example 25 What would the confidence intervals look like in Example 24? How would the image on slide 42 change? What hypothetical repetitions form the reference sets?

stat.epfl.ch Autumn 2024 – slide 49

Bayesian inference

- \Box Our observed data y^{o} are assumed to be a realisation from a density $f(y \mid \theta)$.
- If we can summarise information about θ , separately from $y^{\rm o}$, in a prior density $f(\theta)$, then we base all our uncertainty statements on the posterior density given by Bayes' theorem,

$$f(\theta \mid y^{o}) = \frac{f(y^{o} \mid \theta)f(\theta)}{\int f(y^{o} \mid \theta)f(\theta) d\theta}.$$

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

here θ is regarded as random and y° as fixed.

 \square A point estimate $\tilde{\theta}(y^{\circ})$ of θ is obtained by minimising a posterior expected loss, i.e.,

$$\tilde{\theta}(y^{\mathrm{o}}) = \mathrm{argmin}_{\tilde{\theta}} \mathrm{E} \left\{ L(\theta, \tilde{\theta}) \mid y^{\mathrm{o}} \right\} = \mathrm{argmin}_{\tilde{\theta}} \int L(\theta, \tilde{\theta}) f(\theta \mid y^{\mathrm{o}}) \, \mathrm{d}\theta,$$

where the loss function $L(\theta, \tilde{\theta}) \geq 0$ measures the loss when θ is estimated by $\tilde{\theta}$.

Example 26 Perform Bayesian inference based on $Y_1, \ldots, Y_n \mid \theta \stackrel{\text{iid}}{\sim} U(0, \theta)$ with a Pareto(a, b) prior for θ .

☐ In situations like this, where the support of the density depends on a parameter, it is useful to include an indicator function when writing down the density, viz

$$f(y \mid \theta) = \theta^{-1} I(0 < y < \theta), \quad y \in \mathbb{R}, \theta > 0.$$

As a function of y for fixed θ , its support is the set $(0,\theta)$, but as a function of θ for fixed y, its support is (y,∞) . Sketch these to appreciate the difference.

 \Box The prior density is $f(\theta) = ab^a/\theta^{a+1}I(\theta > b)$ for a, b > 0, and the joint density of the data is

$$f(y \mid \theta) = f(y_1, \dots, y_n \mid \theta) = \prod_{j=1}^n f(y_j \mid \theta) = \prod_{j=1}^n I(0 < y_j < \theta)\theta^{-1} = \theta^{-n}I(0 < m < \theta).$$

where $m = \max(y_1, \dots, y_n)$, so the posterior density is proportional to

$$f(\theta \mid y) \propto f(y_1, \dots, y_n \mid \theta) f(\theta) = \theta^{-n} I(0 < m < \theta) \frac{ab^a}{\theta^{a+1}} I(\theta > b) \propto \theta^{-(A+1)} I(\theta > B),$$

where A=a+n and $B=\max(m,b)$. There are two possibilities here: the prior gives a lower bound b for θ , and if m < b then there is no reason to update this lower bound, but if m > b then clearly $\theta > m > b$, so the bound must be increased at least to m.

- The posterior density has support on (B,∞) and is proportional to $\theta^{-(A+1)}$, so it is $\operatorname{Pareto}(A=a+n,B=\max(y_1,\ldots,y_n,b))$. The p quantile of this distribution satisfies $p=1-(B/\theta_p)^A$, i.e., $\theta_p=B(1-p)^{-1/A}$, which depends on the data and prior; of course 0< p<1.
- ☐ To get a point estimate we might take loss function

$$L(\tilde{\theta}, \theta) = |\tilde{\theta} - \theta| = (\tilde{\theta} - \theta)I(\tilde{\theta} > \theta) + (\theta - \tilde{\theta})I(\theta > \tilde{\theta}).$$

and a standard computation shows that this is minimised at $\tilde{\theta} = \theta_{1/2} = B2^{1/A}$.

stat.epfl.ch

Autumn 2024 - note 1 of slide 50

Comments on Bayesian inference

Often Bayesian models are formulated using a judgement that some variables/observations are exchangeable, as de Finetti theorems then imply that we can write

$$Y_1, \ldots, Y_n \mid \theta \stackrel{\text{iid}}{\sim} f(y; \theta), \qquad \theta \sim f(\theta).$$

- ☐ In general, 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.
- Objectively specifying prior 'ignorance' is problematic and can lead to paradoxes, especially in high dimensions.
- ☐ (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 2024 – slide 51

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.

Example: Shoe data

- \square 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.
- \square The m=10 differences d_1,\ldots,d_m have average $\overline{d}=0.41$.

Boy	Material		Difference
	Α	В	d
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 2024 – slide 53

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
Α	l_j	r_j
В	$\theta + l_j$	$\theta + r_j$

- We observe either $(\theta + l_j, r_j)$ or $(l_j, r_j + \theta)$ so the difference D_j of B and A for boy j is $\theta + l_j r_j$ or $\theta + r_j l_j$. These are equally likely, so we can write $D_j = \theta + I_j c_j$, where
 - θ 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.
- \square If we observe $(\theta + l_j, r_j)$ for boy j, then we cannot observe $(l_j, \theta + r_j)$, which is said to be counterfactual.

Example: Shoe data III

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

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

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

$$E(I_i) = 0, \quad var(I_i) = 1.$$

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

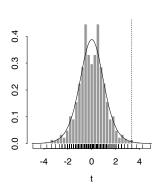
$$S^{2} = \frac{1}{m(m-1)} \sum_{j=1}^{m} (D_{j} - \overline{D})^{2}.$$

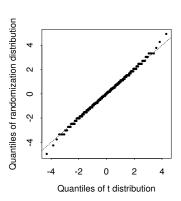
- \square and S^2 can be computed from the observed data, so the standardized quantity $Z=(\overline{D}-\theta)/S$ is an approximate pivot.
- \square If there was no difference between B and A (i.e., $\theta=0$), then $T=\overline{D}/S$ would be symmetrically distributed, as positive and negative values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

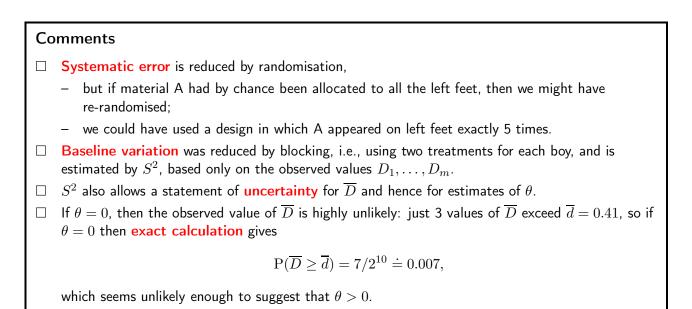
stat.epfl.ch Autumn 2024 – slide 55

Example: Shoe data IV

Randomization distribution of $T=\overline{D}/S$ for the shoes data, i.e., setting $\theta=0$, 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 2024 – slide 57

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

Big picture summary

designed experiments.

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 repeated sampling (frequentist) approach, which invokes hypothetical repetitions of the data-generating mechanism, or
- a Bayesian approach, which requires that 'prior information' on unknown quantities be expressed as a probability distribution, 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.