

# Statistics for Data Science: Week 6

Myrto Limnios and Rajita Chandak

Institute of Mathematics – EPFL

`rajita.chandak@epfl.ch`, `myrto.limnios@epfl.ch`



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

Consider the simplest situation:

$$\Theta_0 = \{\theta_0\} \quad \& \quad \Theta_1 = \{\theta_1\}$$

## The Neyman-Pearson Lemma - Continuous Case

Let  $\mathbf{Y}$  have joint density/frequency  $f \in \{f_0, f_1\}$  and suppose we wish to test

$$H_0 : f = f_0 \quad \text{vs} \quad H_1 : f = f_1.$$

If  $\Lambda(\mathbf{Y}) = f_1(\mathbf{Y})/f_0(\mathbf{Y})$  is a continuous random variable, then there exists a  $k > 0$  such that

$$\mathbb{P}_0[\Lambda(\mathbf{Y}) \geq k] = \alpha$$

and the test whose test function is given by

$$\delta(\mathbf{Y}) = \mathbf{1}\{\Lambda(\mathbf{Y}) \geq k\},$$

is a *most powerful (MP)* test of  $H_0$  versus  $H_1$  at significance level  $\alpha$ .

## Proof.

Use obvious notation  $\mathbb{E}_0, \mathbb{E}_1, \mathbb{P}_0, \mathbb{P}_1$  corresponding to  $H_0$  or  $H_1$ . Let  $G_0(t) = \mathbb{P}_0[\Lambda \leq t]$ . By assumption,  $G_0$  is a differentiable distribution function, and so is onto  $[0, 1]$ . Consequently, the set  $\mathcal{K}_{1-\alpha} = \{t : G_0(t) = 1 - \alpha\}$  is non-empty for any  $\alpha \in (0, 1)$ . Setting  $k = \inf\{t \in \mathcal{K}_{1-\alpha}\}$  we will have  $\mathbb{P}_0[\Lambda \geq k] = \alpha$  and  $k$  is simply the  $1 - \alpha$  quantile of the distribution  $G_0$ . Consequently,

$$\mathbb{P}_0[\delta = 1] = \alpha \quad (\text{since } \mathbb{P}_0[\delta = 1] = \mathbb{P}_0[\Lambda \geq k])$$

and therefore  $\delta \in \mathcal{D}(\{\theta_0\}, \alpha)$  (i.e.  $\delta$  indeed respects the level  $\alpha$ ).

To show that  $\delta$  is also most powerful, it suffices to prove that if  $\psi$  is any function with  $\psi(\mathbf{y}) \in \{0, 1\}$ , then

$$\mathbb{E}_0[\psi(\mathbf{Y})] \leq \underbrace{\mathbb{E}_0[\delta(\mathbf{Y})]}_{=\alpha \text{ (by first part of proof)}} \implies \underbrace{\mathbb{E}_1[\psi(\mathbf{Y})]}_{\beta_1(\psi)} \leq \underbrace{\mathbb{E}_1[\delta(\mathbf{Y})]}_{\beta_1(\delta)}.$$

(recall that  $\beta_1(\delta) = 1 - \mathbb{P}_1[\delta = 0] = \mathbb{P}_1[\delta = 1] = \mathbb{E}_1[\delta]$ ).

WLOG assume that  $f_0$  and  $f_1$  are density functions. Note that

$$f_1(\mathbf{y}) - k \cdot f_0(\mathbf{y}) \geq 0 \text{ if } \delta(\mathbf{y}) = 1 \quad \& \quad f_1(\mathbf{y}) - k \cdot f_0(\mathbf{y}) < 0 \text{ if } \delta(\mathbf{y}) = 0.$$

Therefore, since  $\psi$  can only take the values 0 or 1,

$$\begin{aligned} \psi(\mathbf{y})(f_1(\mathbf{y}) - k \cdot f_0(\mathbf{y})) &\leq \delta(\mathbf{y})(f_1(\mathbf{y}) - k \cdot f_0(\mathbf{y})) \\ \int_{\mathbb{R}^n} \psi(\mathbf{y})(f_1(\mathbf{y}) - k \cdot f_0(\mathbf{y})) d\mathbf{y} &\leq \int_{\mathbb{R}^n} \delta(\mathbf{y})(f_1(\mathbf{y}) - k \cdot f_0(\mathbf{y})) d\mathbf{y} \end{aligned}$$

Rearranging the terms yields

$$\begin{aligned} \int_{\mathbb{R}^n} (\psi(\mathbf{y}) - \delta(\mathbf{y})) f_1(\mathbf{y}) d\mathbf{y} &\leq k \int_{\mathbb{R}^n} (\psi(\mathbf{y}) - \delta(\mathbf{y})) f_0(\mathbf{y}) d\mathbf{y} \\ \implies \mathbb{E}_1[\psi(\mathbf{Y})] - \mathbb{E}_1[\delta(\mathbf{Y})] &\leq k (\mathbb{E}_0[\psi(\mathbf{y})] - \mathbb{E}_0[\delta(\mathbf{Y})]) \end{aligned}$$

But  $k > 0$  by assumption, so when  $\mathbb{E}_0[\psi(\mathbf{Y})] \leq \mathbb{E}_0[\delta(\mathbf{Y})]$  the RHS is negative, i.e.  $\delta$  is an MP test of  $H_0$  vs  $H_1$  at level  $\alpha$ . □

- Basically we reject if the likelihood of  $\theta_0$  is  $k$  times higher than the likelihood of  $\theta_1$ . This is called a likelihood ratio test, and  $\Lambda$  is the likelihood ratio statistic: *how much more plausible is the alternative than the null?*
- When  $\Lambda$  is a continuous RV, the choice of  $k$  is essentially unique. That is, if  $k'$  is such that  $\delta' = \mathbf{1}\{\Lambda \geq k'\} \in \mathcal{D}(\{\theta_0\}, \alpha)$ , then  $\delta = \delta'$  almost surely.
- The resulting most powerful test is not necessarily unique.
- Unless  $\Lambda$  is continuous, the most powerful test is not necessarily guaranteed to exist.
- The problem if  $\Lambda$  is a RV with a discontinuous dist is that there may exist no  $k$  for which the equation  $\mathbb{P}_0[\Lambda \geq k] = \alpha$  has a solution.
- In any case, typically the distribution of the test statistic converges to a continuous limit with large  $n$ , so these problems become inessential.

## Example (Poisson Distribution)

Let  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Poisson}(\mu)$  and for  $\mu_1 > \mu_0$  consider the hypotheses:

$$H_0 : \mu = \mu_0 \quad \text{vs} \quad H_1 : \mu = \mu_1.$$

These correspond to the Higgs example by setting  $\mu_0 = b$  and  $\mu_1 = b + s$ . Applying the Neyman-Pearson lemma gives a test statistic

$$\delta(Y_1, \dots, Y_n) = \mathbf{1} \left\{ \sum_{i=1}^n Y_i > q_{1-\alpha} \right\},$$

provided  $\alpha$  is such that  $G_0(q_{1-\alpha}) = \mathbb{P}_{\mu_0}[\tau(Y_1, \dots, Y_n) \leq q_{1-\alpha}] \stackrel{!}{=} 1 - \alpha$ .

Since the  $Y_i$  are independent, one can easily show that

$$\tau(Y_1, \dots, Y_n) \stackrel{H_0}{\sim} \text{Poisson}(n\mu_0).$$

This being a discrete distribution, the only  $\alpha$  for which we get an MP test are

$$e^{-n\mu_0}, e^{-n\mu_0} (1 + n\mu_0), e^{-n\mu_0} \left( 1 + n\mu_0 + \frac{(n\mu_0)^2}{2} \right), \dots \text{ and so on}$$

Nevertheless notice that as  $n \rightarrow \infty$ , these values become dense near the origin.

When  $\{\Theta_0, \Theta_1\}$  are not singletons, choosing a **most powerful test** is a **much stronger requirement**:

- 1 It should respect the level for all  $\theta \in \Theta_0$ , i.e.

$$\delta \in \mathcal{D}(\Theta_0, \alpha) = \{\delta : \mathcal{Y}^n \rightarrow \{0, 1\} : \mathbb{E}_\theta[\delta] \leq \alpha, \forall \theta \in \Theta_0\}$$

- 2 It should be most powerful for all  $\theta \in \Theta_1$  (i.e. for all possible simple alternatives),

$$\mathbb{E}_\theta[\delta] \geq \mathbb{E}_\theta[\delta'] \quad \forall \theta \in \Theta_1 \quad \& \quad \delta' \in \mathcal{D}(\Theta_0, \alpha)$$

Unfortunately UMP tests rarely exist. **Why?**

→ Consider  $H_0 : \theta = \theta_0$  vs  $H_1 : \theta \neq \theta_0$

- A UMP test must be MP test for any  $\theta \neq \theta_0$ .
- But the form of the MP test typically differs for  $\theta_1 > \theta_0$  and  $\theta_1 < \theta_0$ !  
→ e.g. recall exponential mean example

## Example (No UMP test exists)

Let  $Y_1, \dots, Y_n \sim \text{Bernoulli}(\theta)$  and suppose we want to test:

$$H_0 : \theta = \theta_0 \quad \text{vs} \quad H_1 : \theta \neq \theta_0$$

at some level  $\alpha$ . To this aim, consider first

$$H'_0 : \theta = \theta_0 \quad \text{vs} \quad H'_1 : \theta = \theta_1$$

Neyman-Pearson lemma gives test statistics

$$T = \frac{f(\mathbf{Y}; \theta_1)}{f(\mathbf{Y}; \theta_0)} = \left( \frac{1 - \theta_1}{1 - \theta_0} \right)^n \left( \frac{\theta_1(1 - \theta_0)}{\theta_0(1 - \theta_1)} \right)^{\sum_{i=1}^n Y_i}$$

- If  $\theta_1 > \theta_0$  then  $T$  increasing in  $\sum_{i=1}^n Y_i$   
     $\rightarrow$  MP test would reject for large values of  $\sum_{i=1}^n Y_i$
- If  $\theta_1 < \theta_0$  then  $T$  decreasing in  $\sum_{i=1}^n Y_i$   
     $\rightarrow$  MP test would reject for small values of  $\sum_{i=1}^n Y_i$



So what can we do for more general  $\{\Theta_0, \Theta_1\}$ ?

- **One sided hypotheses:** when  $\Theta_0$  is an interval of the form  $(-\infty, \theta_0]$  or  $[\theta_0, +\infty)$  and  $\Theta_1 = \Theta_0^c$ , there are often uniformly most powerful tests depending on the underlying model.
  - For example, in one-parameter exponential families, one simply uses the Neyman-Pearson lemma, taking the null to be  $\theta = \theta_0$  and the alternative  $\theta = \theta_1$  for any  $\theta_1 \in \Theta_1$  (the form of the test depends only on the direction of the null and the boundary of the null).
  - This generalises to families admitting a so-called “monotone likelihood ratio”
  - In the absence of the “monotone likelihood ratio” property, one can seek **locally most powerful tests**, near the hypothesis boundary. It can be shown that the score function (derivative of the loglikelihood) at the boundary  $\theta_0$  can serve as a test statistic to this aim.
- **General hypothesis pairs:** we need to abandon optimality, and search for sensible tests. But the **likelihood ratio** idea can serve us well in this pursuit.

Consider now the multiparameter case  $\theta \in \mathbb{R}^p$  with general  $\Theta_0, \Theta_1$

- As noted optimality breaks down.
- But we can still seek general-purpose approaches.

**The idea:** Combine Neyman-Pearson paradigm with Maximum Likelihood

### Definition (Likelihood Ratio)

The *likelihood ratio statistic* corresponding to the pair of hypotheses  $H_0 : \theta \in \Theta_0$  vs  $H_1 : \theta \in \Theta_1$  is defined to be

$$\Lambda(\mathbf{Y}) = \frac{\sup_{\theta \in \Theta} f(\mathbf{Y}; \theta)}{\sup_{\theta \in \Theta_0} f(\mathbf{Y}; \theta)} = \frac{\sup_{\theta \in \Theta} L(\theta)}{\sup_{\theta \in \Theta_0} L(\theta)}$$

- Intuition: choose the “most favourable”  $\theta \in \Theta_0$  (in favour of  $H_0$ ) and compare it against the “most favourable”  $\theta \in \Theta_1$  (in favour of  $H_1$ ) in a simple vs simple setting (applying NP-lemma)
- Typically  $\Theta_0$  is a lower dimensional subspace of  $\Theta_1$ , so taking sup over  $\Theta$  (rather than  $\Theta_1$ ) incurs no loss. In this case  $\Theta_0 \cap \Theta_1 \neq \emptyset$ , but  $\text{Leb}(\Theta_0 \cap \Theta_1) = 0$ , which suffices.

## Example

Let  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$  where both  $\mu$  and  $\sigma^2$  are unknown. Consider:

$$H_0 : \mu = \mu_0 \quad \text{vs} \quad H_1 : \mu \neq \mu_0$$

$$\Lambda(\mathbf{Y}) = \frac{\sup_{(\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}^+} f(\mathbf{Y}; \mu, \sigma^2)}{\sup_{(\mu, \sigma^2) \in \{\mu_0\} \times \mathbb{R}^+} f(\mathbf{Y}; \mu, \sigma^2)} = \left( \frac{\sum_{i=1}^n (Y_i - \mu_0)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \right)^{\frac{n}{2}} = \left( \frac{\hat{\sigma}_0^2}{\hat{\sigma}^2} \right)^{\frac{n}{2}}$$

So reject when  $\Lambda \geq k$ , where  $k$  is s.t.  $\mathbb{P}_0[\Lambda \geq k] = \alpha$ . **Distribution of  $\Lambda$ ?** By monotonicity look only at

$$\begin{aligned} \frac{\sum_{i=1}^n (Y_i - \mu_0)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} &= 1 + \frac{n(\bar{Y} - \mu_0)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 + \frac{1}{n-1} \left( \frac{n(\bar{Y} - \mu_0)^2}{S^2} \right) \\ &= 1 + \frac{T^2}{n-1} \end{aligned}$$

With  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$  and  $T = \sqrt{n}(\bar{Y} - \mu_0)/S \stackrel{H_0}{\sim} t_{n-1}$ .

So  $T^2 \stackrel{H_0}{\sim} F_{1, n-1}$  and  $k$  may be chosen appropriately.

## Example

Let  $Y_1, \dots, Y_m \stackrel{iid}{\sim} \text{Exp}(\lambda)$  and  $Z_1, \dots, Z_n \stackrel{iid}{\sim} \text{Exp}(\theta)$ . Assume  $\mathbf{Y}$  indep  $\mathbf{Z}$ .

Consider:  $H_0 : \theta = \lambda$  vs  $H_1 : \theta \neq \lambda$

i.e.  $(\theta, \lambda) \in R_+^2$  against  $(\theta, \lambda) \in 45 \text{ degree line}$

Unrestricted MLEs:  $\hat{\lambda} = 1/\bar{Y}$  &  $\hat{\theta} = 1/\bar{Z}$   
 $\sup_{(\lambda, \theta) \in \mathbb{R}_+^2} f(\mathbf{Y}, \mathbf{Z}; \lambda, \theta)$

$$\begin{aligned} \text{Restricted MLEs: } \hat{\lambda}_0 = \hat{\theta}_0 &= \left[ \frac{m\bar{Y} + n\bar{Z}}{m+n} \right]^{-1} \\ \sup_{(\lambda, \theta) \in \{(y, z) \in \mathbb{R}_+^2 : y=z\}} f(\mathbf{Y}, \mathbf{Z}; \lambda, \theta) \\ \implies \Lambda &= \left( \frac{m}{m+n} + \frac{n}{n+m} \frac{\bar{Z}}{\bar{Y}} \right)^m \left( \frac{n}{n+m} + \frac{m}{m+n} \frac{\bar{Y}}{\bar{Z}} \right)^n \end{aligned}$$

Depends on  $T = \bar{Y}/\bar{Z}$  and can make  $\Lambda$  large/small by varying  $T$ .

$\hookrightarrow$  But  $T \stackrel{H_0}{\sim} F_{2m, 2n}$  so given  $\alpha$  we may find the critical value  $k$ .

More often than not,  $\text{dist}(\Lambda)$  intractable

$\hookrightarrow$  (and no simple dependence on  $T$  with tractable distribution either)

Consider asymptotic approximations?

Setup

- $\Theta$  open subset of  $\mathbb{R}^p$
- either  $\Theta_0 = \{\theta_0\}$  or  $\Theta_0$  open subset of  $\mathbb{R}^s$ , where  $s < p$
- Concentrate on  $\mathbf{Y} = (Y_1, \dots, Y_n)$  has iid components.
- Initially restrict attention to  $H_0 : \theta = \theta_0$  vs  $H_1 : \theta \neq \theta_0$ . LR becomes:

$$\Lambda_n(\mathbf{Y}) = \prod_{i=1}^n \frac{f(Y_i; \hat{\theta}_n)}{f(Y_i; \theta_0)}$$

where  $\hat{\theta}_n$  is the MLE of  $\theta$ .

- Impose regularity conditions from MLE asymptotics

Theorem (Wilks' Theorem, case  $p = 1$ )

Let  $Y_1, \dots, Y_n$  be iid random variables with density (frequency) depending on  $\theta \in \mathbb{R}$  and satisfying conditions (A1)-(A6), with  $\mathcal{I}_1(\theta) = \mathcal{J}_1(\theta)$ . If the MLE sequence  $\hat{\theta}_n$  is consistent for  $\theta$ , then the likelihood ratio statistic  $\Lambda_n$  for  $H_0 : \theta = \theta_0$  satisfies

$$2 \log \Lambda_n \xrightarrow{d} V \sim \chi_1^2$$

when  $H_0$  is true.

- Obviously, knowing approximate distribution of  $2 \log \Lambda_n$  is as good as knowing approximate distribution of  $\Lambda_n$  for the purposes of testing (by monotonicity and rejection method).
- Theorem extends immediately and trivially to the case of general  $p$  and for a hypothesis pair  $H_0 : \theta = \theta_0$  vs  $H_1 : \theta \neq \theta_0$ .  
(i.e. when null hypothesis is simple)

## Proof (\*).

Under the conditions of the theorem and when  $H_0$  is true,

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} \mathcal{N}(0, \mathcal{I}_1^{-1}(\theta_0))$$

Now take logarithms and expand in a Taylor series around  $\hat{\theta}_n$ ,

$$\begin{aligned} \log \Lambda_n &= \sum_{i=1}^n [\ell(Y_i; \hat{\theta}_n) - \ell(Y_i; \theta_0)] = \sum_{i=1}^n [\ell(Y_i; \hat{\theta}_n) - \ell(Y_i; \hat{\theta}_n)] + \\ &\quad + (\theta_0 - \hat{\theta}_n) \sum_{i=1}^n \ell'(Y_i; \hat{\theta}_n) - \frac{1}{2}(\hat{\theta}_n - \theta_0)^2 \sum_{i=1}^n \ell''(Y_i; \theta_n^*) \\ &= -\frac{1}{2}n(\hat{\theta}_n - \theta_0)^2 \frac{1}{n} \sum_{i=1}^n \ell''(Y_i; \theta_n^*) \end{aligned}$$

where  $\theta_n^*$  lies between  $\hat{\theta}_n$  and  $\theta_0$ .

If  $H_0$  is true, and since  $\hat{\theta}_n$  is a consistent sequence,  $\theta_n^*$  is sandwiched so

$$\theta_n^* \xrightarrow{P} \theta_0.$$

Hence under assumptions (A1)-(A6), and when  $H_0$  is true, a first order Taylor expansion about  $\theta_0$ , the continuous mapping theorem and the LLN give

$$\frac{1}{n} \sum_{i=1}^n \ell''(Y_i; \theta_n^*) \xrightarrow{P} -\mathbb{E}_{\theta_0}[\ell''(Y_i; \theta_0)] = \mathcal{I}_1(\theta_0)$$

On the other hand, by the continuous mapping theorem,

$$n(\hat{\theta}_n - \theta_0)^2 \xrightarrow{d} \frac{V}{\mathcal{I}_1(\theta_0)}$$

Applying Slutsky's theorem now yields the result. □ □



### Theorem (Wilk's theorem, general $p$ , general $s \leq p$ )

Let  $Y_1, \dots, Y_n$  be iid random variables with density (frequency) depending on  $\theta \in \mathbb{R}^p$  and satisfying conditions (B1)-(B6), with  $\mathcal{I}_1(\theta) = \mathcal{J}_1(\theta)$ . If the MLE sequence  $\hat{\theta}_n$  is consistent for  $\theta$ , then the likelihood ratio statistic  $\Lambda_n$  for  $H_0 : \{\theta_j = \theta_{j,0}\}_{j=1}^s$  satisfies  $2 \log \Lambda_n \xrightarrow{d} V \sim \chi_s^2$  when  $H_0$  is true.

#### Comments:

- Note that it may potentially be that  $s < p$ , and this is accommodated by the theorem
- Hypotheses of the form  $H_0 : \{g_j(\theta) = a_j\}_{j=1}^s$ , for  $g_j$  differentiable real functions, can also be handled by Wilks' theorem:
  - Define  $(\phi_1, \dots, \phi_p) = g(\theta) = (g_1(\theta), \dots, g_p(\theta))$
  - $g_{s+1}, \dots, g_p$  defined so that  $\theta \mapsto g(\theta)$  is 1-1
  - Apply theorem with parameter  $\phi$

Many other tests possible. For example:

- Wald's test

- ↪ For a simple null, may compare the unrestricted MLE with the MLE under the null. Large deviations indicate evidence against null hypothesis. Distributions are approximated for large  $n$  via the asymptotic normality of MLEs.

- Score Test

- ↪ For a simple null, if the null hypothesis is false, then the loglikelihood gradient at the null should not be close to zero, at least when  $n$  reasonably large: so measure its deviations from zero. Use asymptotics for distributions (under conditions we end up with a  $\chi^2$ )

- ...

So far focussed on Neyman-Pearson Framework:

- ❶ Fix a significance level  $\alpha$  for the test
- ❷ Consider rules  $\delta$  respecting this significance level
  - We choose one of those rules,  $\delta^*$ , based on power considerations
- ❸ We reject at level  $\alpha$  if  $\delta^*(\mathbf{y}) = 1$ .

Useful for attempting to determine optimal test statistics

What if we already have a given form of test statistic in mind? (e.g. LRT)

→ A different perspective on testing (used more in practice) says:

Rather than consider a family of test functions respecting level  $\alpha$ ...

... consider family of test functions indexed by  $\alpha$

- ❶ Fix a family  $\{\delta_\alpha\}_{\alpha \in (0,1)}$  of decision rules, with  $\delta_\alpha$  having level  $\alpha$ 
  - for a given  $\mathbf{y}$  some of these rules reject the null, while others do not
- ❷ Which is the smallest  $\alpha$  for which  $H_0$  is rejected given  $\mathbf{y}$ ?

## Definition ( $p$ -Value)

Let  $\{\delta_\alpha\}_{\alpha \in (0,1)}$  be a family of test functions satisfying

$$\alpha_1 < \alpha_2 \implies \{\mathbf{y} \in \mathcal{Y}^n : \delta_{\alpha_1}(\mathbf{y}) = 1\} \subseteq \{\mathbf{y} \in \mathcal{Y}^n : \delta_{\alpha_2}(\mathbf{y}) = 1\}.$$

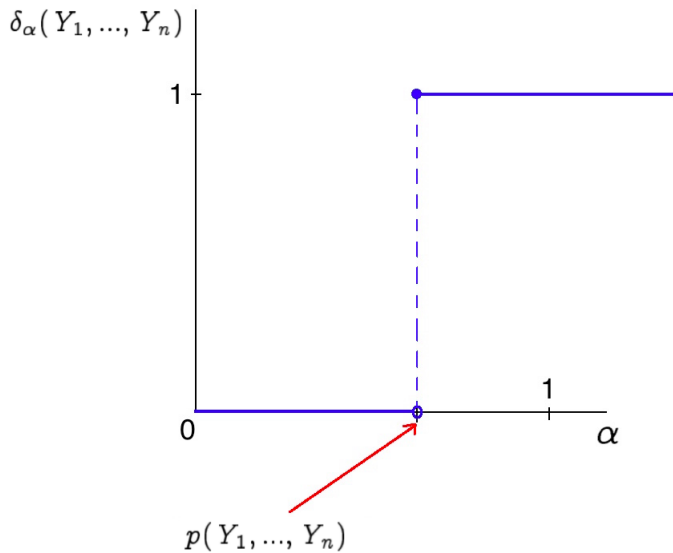
The  $p$ -value (or observed significance level) of the family  $\{\delta_\alpha\}$  is

$$p(\mathbf{y}) = \inf\{\alpha : \delta_\alpha(\mathbf{y}) = 1\}$$

$\hookrightarrow$  The  $p$ -value is the smallest value of  $\alpha$  for which the null would be rejected at level  $\alpha$ , given  $\mathbf{Y} = \mathbf{y}$ .

Most usual setup:

- Have a single test statistic  $T$
- Construct family  $\delta_\alpha(\mathbf{y}) = \mathbf{1}\{T(\mathbf{y}) > k_\alpha\}$
- If  $\mathbb{P}_{H_0}[T \leq t] = G(t)$  then  $p(\mathbf{y}) = \mathbb{P}_{H_0}[T(\mathbf{Y}) \geq T(\mathbf{y})] = 1 - G(T(\mathbf{y}))$



Notice: contrary to NP-framework did not make explicit decision!

- We simply reported a  $p$ -value
- The  $p$ -value is used as a measure of evidence against  $H_0$ 
  - ↪ Small  $p$ -value provides evidence against  $H_0$
  - ↪ Large  $p$ -value provides no evidence against  $H_0$
- How small does “small” mean?
  - ↪ Depends on the specific problem...

Intuition:

- Recall that extreme values of test statistics are those that are “inconsistent” with null (NP-framework)
- $p$ -value is probability of observing a value of the test statistic as extreme as or more extreme than the one we observed, under the null
- If this probability is small, then we have witnessed something quite unusual under the null hypothesis
- Gives evidence against the null hypothesis

## Example (Normal Mean)

Let  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$  where both  $\mu$  and  $\sigma^2$  are unknown. Consider:

$$H_0 : \mu = 0 \quad \text{vs} \quad H_1 : \mu \neq 0$$

Likelihood ratio test: reject when  $T^2$  large,  $T = \sqrt{n}\bar{Y}/S \stackrel{H_0}{\sim} t_{n-1}$ .

Since  $T^2 \stackrel{H_0}{\sim} F_{1,n-1}$ ,  $p$ -value is

$$p(\mathbf{y}) = \mathbb{P}_{H_0}[T^2(\mathbf{Y}) \geq T^2(\mathbf{y})] = 1 - G_{F_{1,n-1}}(T^2(\mathbf{y}))$$

Consider two samples (datasets),

$$\mathbf{y} = (0.66, 0.28, -0.99, 0.007, -0.29, -1.88, -1.24, 0.94, 0.53, -1.2)$$

$$\mathbf{y}' = (1.4, 0.48, 2.86, 1.02, -1.38, 1.42, 2.11, 2.77, 1.02, 1.87)$$

Obtain  $p(\mathbf{y}) = 0.32$  while  $p(\mathbf{y}') = 0.006$ .

- Reporting a  $p$ -value does not necessarily mean making a decision
- A small  $p$ -value can simply reflect our “confidence” in rejecting a null

Recall example: **Statisticians working for Trump** gather iid sample  $\mathbf{Y}$  from Florida with  $Y_i = \mathbf{1}\{\text{vote Biden}\}$ . Trumps team want to test

$$\begin{cases} H_0 : \text{Trump wins Florida} \\ H_1 : \text{Biden wins Florida} \end{cases}$$

- Will statisticians decide for Trump?
- Perhaps better to report  $p$ -value to him and let him decide...

**What if statisticians working for newspaper, not Trump?**

- Something easier to interpret than test/ $p$ -value?



## A Glance Back at Point Estimation

- Let  $Y_1, \dots, Y_n$  be iid random variables with density (frequency)  $f(\cdot; \theta)$ .
- Problem with point estimation:  $\mathbb{P}_\theta[\hat{\theta} = \theta]$  typically small (if not zero)
  - ↪ always attach an estimator of variability, e.g. standard error
  - ↪ interpretation?
- Hypothesis tests may provide way to interpret estimator's variability within the setup of a particular problem
  - ↪ e.g. if observe  $\hat{P}[\text{Biden wins}] = 0.52$  can actually see what  $p$ -value we get when testing  $H_0 : P[\text{Biden wins}] \geq 1/2$ .
- Something more directly interpretable?

Back to our example: [What do pollsters do in newspapers?](#)

- ↪ They announce their point estimate (e.g. 0.52)
- ↪ They give upper and lower *confidence limits*

What are these and how are they interpreted?

Simple underlying idea:

- Instead of estimating  $\theta$  by a single value
- Present a whole range of values for  $\theta$  that are consistent with the data  
→ In the sense that they could have produced the data

### Definition (Confidence Interval)

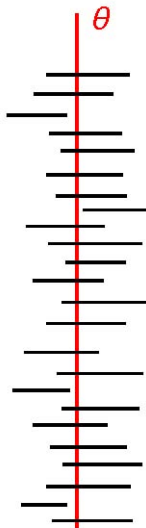
Let  $\mathbf{Y} = (Y_1, \dots, Y_n)$  be random variables with joint distribution depending on  $\theta \in \mathbb{R}$  and let  $L(\mathbf{Y})$  and  $U(\mathbf{Y})$  be two statistics with  $L(\mathbf{Y}) < U(\mathbf{Y})$  a.s. Then, the random interval  $[L(\mathbf{Y}), U(\mathbf{Y})]$  is called a  $100(1 - \alpha)\%$  confidence interval for  $\theta$  if

$$\mathbb{P}_\theta[L(\mathbf{Y}) \leq \theta \leq U(\mathbf{Y})] \geq 1 - \alpha$$

for all  $\theta \in \Theta$ , with equality for at least one value of  $\theta$ .

- $1 - \alpha$  is called the coverage probability or confidence level
- Beware of interpretation!

- Probability statement is **NOT** made about  $\theta$ , which is constant.
- Statement is about interval: probability that the interval contains the true value is at least  $1 - \alpha$ .
- Given any realization  $\mathbf{Y} = \mathbf{y}$ , the interval  $[L(\mathbf{y}), U(\mathbf{y})]$  will either contain or not contain  $\theta$ .
- Interpretation: if we construct intervals with this method, then we expect that  $100(1 - \alpha)\%$  of the time our intervals will engulf the true value.



## Example (The example that says all)

Let  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \mathcal{N}(\mu, 1)$ . Then  $\sqrt{n}(\bar{Y} - \mu) \sim \mathcal{N}(0, 1)$ , so that

$$\mathbb{P}_\mu[-1.96 \leq \sqrt{n}(\bar{Y} - \mu) \leq 1.96] = 0.95$$

and since

$$-1.96 \leq \sqrt{n}(\bar{Y} - \mu) \leq 1.96 \iff \bar{Y} - 1.96/\sqrt{n} \leq \mu \leq \bar{Y} + 1.96/\sqrt{n}$$

we obviously have

$$\mathbb{P}_\mu \left[ \bar{Y} - \frac{1.96}{\sqrt{n}} \leq \mu \leq \bar{Y} + \frac{1.96}{\sqrt{n}} \right] = 0.95$$

So that the random interval  $[L(\mathbf{Y}), U(\mathbf{Y})] = \left[ \bar{Y} - \frac{1.96}{\sqrt{n}}, \bar{Y} + \frac{1.96}{\sqrt{n}} \right]$  is a 95% confidence interval for  $\mu$ .

Central Limit Theorem: same argument can yield approximate 95% CI when  $Y_1, \dots, Y_n$  are iid,  $\mathbb{E}Y_i = \mu$  and  $\text{var}(Y_i) = 1$ , regardless of their distribution.

## Example (continued)

Notice that the interval is centred at  $\bar{Y}$ , the MLE of  $\mu$ . It's often thus written:

$$\bar{Y} \pm z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

### Observations:

- The length of the interval is  $2z_{1-\alpha/2}\sigma/\sqrt{n}$ , which depends on  $\sigma^2$ ,  $n$  and  $\alpha$ .
- The parameter  $\sigma^2$  is beyond our control.
- We can nevertheless control  $n$  and  $1 - \alpha$ . Increasing  $n$ , the length of the interval decays as  $1/\sqrt{n}$ .
- Reducing  $\alpha$  (i.e. increasing  $1 - \alpha$ ) increases the length of the interval (the dependence is quite non-linear, and 5% is chosen as a “sweet spot”).

What can we learn from previous example?

### Definition (Pivot)

A random function  $g(\mathbf{Y}, \theta)$  is said to be a *pivotal quantity* (or simply a *pivot*) if it is a function both of  $\mathbf{Y}$  and  $\theta$  whose distribution does not depend on  $\theta$ .

$\hookrightarrow \sqrt{n}(\bar{Y} - \mu) \sim \mathcal{N}(0, 1)$  is a pivot in previous example

Why is a pivot useful?

- $\forall \alpha \in (0, 1)$  we can find constants  $a < b$  independent of  $\theta$ , such that

$$\mathbb{P}_{\theta}[a \leq g(\mathbf{Y}, \theta) \leq b] = 1 - \alpha \quad \forall \theta \in \Theta$$

- If  $g(\mathbf{Y}, \theta)$  can be manipulated then the above yields a CI

## Example

Let  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \mathcal{U}[0, \theta]$ . Recall that MLE  $\hat{\theta}$  is  $\hat{\theta} = Y_{(n)}$ , with distribution

$$\mathbb{P}_{\theta} [Y_{(n)} \leq x] = F_{Y_{(n)}}(x) = \left(\frac{x}{\theta}\right)^n \implies \mathbb{P}_{\theta} \left[\frac{Y_{(n)}}{\theta} \leq y\right] = y^n$$

→ Hence  $Y_{(n)}/\theta$  is a pivot for  $\theta$ . Can now choose  $a < b$  such that

$$\mathbb{P}_{\theta} \left[ a \leq \frac{Y_{(n)}}{\theta} \leq b \right] = 1 - \alpha$$

→ But there are  $\infty$ -many such choices!

↪ Idea: choose pair  $(a, b)$  that minimizes interval's length!

Solution can be seen to be  $a = \alpha^{1/n}$  and  $b = 1$ , yielding

$$\left[ Y_{(n)}, \frac{Y_{(n)}}{\alpha^{1/n}} \right]$$

Pivotal method extends to construction of CI for  $\theta_k$ , when

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_k, \dots, \theta_p) \in \mathbb{R}^p$$

and the remaining coordinates are also unknown.  $\rightarrow$  Pivotal quantity should now be function  $g(\mathbf{Y}; \theta_k)$  which

- ① Depends on  $\mathbf{Y}$ ,  $\theta_k$ , but no other parameters
- ② Has a distribution independent of any of the parameters

$\rightarrow$  e.g.: CI for normal mean, when variance unknown

$\rightarrow$  Main difficulties with pivotal method:

- Hard to find exact pivots in general problems
- Exact distributions may be intractable

Resort to asymptotic approximations...

$\hookrightarrow$  Most classic example when have  $a_n(\hat{\theta}_n - \theta) \xrightarrow{d} \mathcal{N}(0, \sigma^2(\theta))$ .



What about higher dimensional parameters?

### Definition (Confidence Region)

Let  $\mathbf{Y} = (Y_1, \dots, Y_n)$  be random variables with joint distribution depending on  $\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^p$ . A random subset  $R(\mathbf{Y})$  of  $\Theta$  depending on  $\mathbf{Y}$  is called a  $100(1 - \alpha)\%$  confidence region for  $\boldsymbol{\theta}$  if

$$\mathbb{P}_{\boldsymbol{\theta}}[R(\mathbf{Y}) \ni \boldsymbol{\theta}] \geq 1 - \alpha$$

for all  $\boldsymbol{\theta} \in \Theta$ , with equality for at least one value of  $\boldsymbol{\theta}$ .

- No restriction requiring  $R(\mathbf{Y})$  to be convex or even connected
  - ↪ So when  $p = 1$  get more general notion than CI
- Nevertheless, many notions extend immediately to CR case
  - ↪ e.g. notion of a pivotal quantity

Let  $g : \mathcal{Y}^n \times \Theta \rightarrow \mathbb{R}$  be a function such that  $\text{dist}[g(\mathbf{Y}, \boldsymbol{\theta})]$  independent of  $\boldsymbol{\theta}$   
 $\hookrightarrow$  Since image space is the real line, can find  $a < b$  s.t.

$$\mathbb{P}_{\boldsymbol{\theta}}[a \leq g(\mathbf{Y}, \boldsymbol{\theta}) \leq b] = 1 - \alpha$$

$$\implies \mathbb{P}_{\boldsymbol{\theta}}[R(\mathbf{Y}) \ni \boldsymbol{\theta}] = 1 - \alpha$$

where  $R(\mathbf{y}) = \{\boldsymbol{\theta} \in \Theta : g(\mathbf{y}, \boldsymbol{\theta}) \in [a, b]\}$

Notice that region can be “wild” since it is a random level set of  $g$

## Example

Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n \stackrel{iid}{\sim} \mathcal{N}_k(\boldsymbol{\mu}, \Sigma)$ . Two unbiased estimators of  $\boldsymbol{\mu}$  and  $\Sigma$  are

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{Y}_i$$

$$\hat{\Sigma} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{Y}_i - \hat{\boldsymbol{\mu}})(\mathbf{Y}_i - \hat{\boldsymbol{\mu}})^T$$

## Example (cont'd)

Consider the random variable

$$g(\{\mathbf{Y}\}_{i=1}^n, \boldsymbol{\mu}) := \frac{n(n-k)}{k(n-1)} (\hat{\boldsymbol{\mu}} - \boldsymbol{\mu})^T \hat{\boldsymbol{\Sigma}}^{-1} (\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}) \sim F\text{-dist with } k \text{ and } n-k \text{ d.f.}$$

A pivot!

$\hookrightarrow$  If  $f_q$  is  $q$ -quantile of this distribution, then get  $100q\%$  CR as

$$R(\{\mathbf{Y}\}_{i=1}^n) = \left\{ \boldsymbol{\theta} \in \mathbb{R}^n : \frac{n(n-k)}{k(n-1)} (\hat{\boldsymbol{\mu}} - \boldsymbol{\mu})^T \hat{\boldsymbol{\Sigma}}^{-1} (\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}) \leq f_q \right\}$$

- An ellipsoid in  $\mathbb{R}^n$
- Ellipsoid centred at  $\hat{\boldsymbol{\mu}}$
- Principle axis lengths given by eigenvalues of  $\hat{\boldsymbol{\Sigma}}^{-1}$
- Orientation given by eigenvectors of  $\hat{\boldsymbol{\Sigma}}^{-1}$

Visualisation of high-dimensional CR's can be hard

- When these are ellipsoids, spectral decomposition helps
- But more generally?

Things are especially easy when dealing with rectangles - **but they rarely occur!**

→ What if we construct a CR as Cartesian product of CI's?

Let  $[L_i(\mathbf{Y}), U_i(\mathbf{Y})]$  be  $100q_i\%$  CI's for  $\theta_i$ ,  $i = 1, \dots, p$ , and define

$$R(\mathbf{Y}) = [L_1(\mathbf{Y}), U_1(\mathbf{Y})] \times \dots \times [L_p(\mathbf{Y}), U_p(\mathbf{Y})]$$

Bonferroni's inequality implies that

$$\mathbb{P}_{\theta}[R(\mathbf{Y}) \ni \theta] \geq 1 - \sum_{i=1}^p \mathbb{P}[\theta_i \notin [L_i(\mathbf{Y}), U_i(\mathbf{Y})]] = 1 - \sum_{i=1}^p (1 - q_i)$$

→ So pick  $q_i$  such that  $\sum_{i=1}^p (1 - q_i) = \alpha$  (can be conservative...)

Discussion on CR's  $\rightarrow$  provides no guidance on choosing “good” regions

**But:**  $\exists$  close relationship between CR's and hypothesis tests!

$\hookrightarrow$  exploit this to transform good testing properties into good CR properties

Suppose  $R(\mathbf{Y})$  is an exact  $100q\%=100(1-\alpha)\%$  CR for  $\theta$ . Consider

$$H_0 : \theta = \theta_0 \quad \text{vs} \quad H_1 : \theta \neq \theta_0$$

Define test function:

$$\delta(\mathbf{Y}) = \begin{cases} 1 & \text{if } \theta_0 \notin R(\mathbf{Y}), \\ 0 & \text{if } \theta_0 \in R(\mathbf{Y}). \end{cases}$$

Then,  $\mathbb{E}_{\theta_0}[\delta(\mathbf{Y})] = 1 - \mathbb{P}_{\theta_0}[\theta_0 \in R(\mathbf{Y})] \leq \alpha$

Can use a CR to construct test with significance level  $\alpha$ !

Going the other way around, *can invert* tests to get CR's:

Suppose we have tests at level  $\alpha$  for any choice of simple null,  $\theta_0 \in \Theta$ .

→ Say that  $\delta(\mathbf{Y}; \theta_0)$  is appropriate test function for  $H_0 : \theta = \theta_0$

Define 
$$R^*(\mathbf{Y}) = \{\theta_0 : \delta(\mathbf{Y}; \theta_0) = 0\}$$

Coverage probability of  $R^*(\mathbf{Y})$  is

$$\mathbb{P}_\theta[\theta \in R^*(\mathbf{Y})] = \mathbb{P}_\theta[\delta(\mathbf{Y}; \theta) = 0] \geq 1 - \alpha$$

Obtain a  $100(1 - \alpha)\%$  confidence region by choosing all the  $\theta$  for which the null would not be rejected given our data  $\mathbf{Y}$ .

→ If test inverted is powerful, then get “small” region for given  $1 - \alpha$ .

Modern example: looking for signals in noise

- Interested in detecting presence of a signal  $\mu(x_t)$ ,  $t = 1, \dots, T$  over a discretised domain,  $\{x_1, \dots, x_t\}$ , on the basis of noisy measurements
- This is to be detected against some known background, say 0.
- May or may not be specifically interested in detecting the presence of the signal in some particular location  $x_t$ , but in detecting whether the a signal is present anywhere in the domain.

Formally:

Does there exist a  $t \in \{1, \dots, T\}$  such that  $\mu(x_t) \neq 0$ ?

or

for which  $t$ 's is  $\mu(x_t) \neq 0$ ?

More generally:

- Observe

$$Y_t = \mu(x_t) + \varepsilon_t, \quad t = 1, \dots, T.$$

- Wish to test, at some significance level  $\alpha$ :

$$\begin{cases} H_0 : \mu(x_t) = 0 & \text{for all } t \in \{1, \dots, T\}, \\ H_A : \mu(x_t) \neq 0 & \text{for some } t \in \{1, \dots, T\}. \end{cases}$$

- May also be interested in which specific locations signal deviates from zero
- More generally: May have  $T$  hypotheses to test simultaneously at level  $\alpha$  (they may be related or totally unrelated)
- Suppose we have a test statistic for each individual hypothesis  $H_{0,t}$  yielding a  $p$ -value  $p_t$ .



## Bonferroni Method.

If we test each hypothesis individually, we will not maintain the level!

Can we maintain the level  $\alpha$ ?

Idea: use the same trick as for confidence regions!

### Bonferroni

- 1 Test individual hypotheses separately at level  $\alpha_t = \alpha/T$
- 2 Reject  $H_0$  if at least one of the  $\{H_{0,t}\}_{t=1}^T$  is rejected

Global level is bounded as follows:

$$\mathbb{P}[\text{not } H_0 | H_0] = \mathbb{P}\left[\bigcup_{t=1}^T \{\text{not } H_{0,t}\} \middle| H_0\right] \leq \sum_{t=1}^T \mathbb{P}[\text{not } H_{0,t} | H_0] = T \frac{\alpha}{T} = \alpha$$

## Holm-Bonferroni Method.

- Advantage: Works for any (discrete domain) setup!
- Disadvantage: Too conservative when  $T$  large

Holm's modification increases average # of hypotheses rejected at level  $\alpha$  (but does not increase power for overall rejection of  $H_0 = \cap_{t \in T} H_{0,t}$ )

### Holm's Procedure

- 1 We reject  $H_{0,t}$  for small values of a corresponding  $p$ -value,  $p_t$
- 2 Order  $p$ -values from most to least significant:  $p_{(1)} \leq \dots \leq p_{(T)}$
- 3 Starting from  $t = 1$  and going up, reject all  $H_{0,(t)}$  such that  $p_{(t)}$  significant at level  $\alpha/(T - t + 1)$ . Stop rejecting at first insignificant  $p_{(t)}$ .

Genuine improvement over Bonferroni if want to detect as many signals as possible, not just existence of some signal

Both Holm and Bonferroni reject the global  $H_0$  if and only if  $\inf_t p_t$  significant at level  $\alpha/T$ .

## Taking Advantage of Structure: Independence.

In the (special) case where individual test statistics are independent, one may use Sime's (in)equality,

$$\mathbb{P} \left[ p_{(j)} \geq \frac{j\alpha}{T}, \text{ for all } j = 1, \dots, T \mid H_0 \right] \geq 1 - \alpha$$

(strict equality requires continuous test statistics, otherwise  $\leq \alpha$ )

### Yields Sime's procedure (assuming independence)

- 1 Suppose we reject  $H_{0,j}$  for small values of  $p_j$
- 2 Order  $p$ -values from most to least significant:  $p_{(1)} \leq \dots \leq p_{(T)}$
- 3 If, for some  $j = 1, \dots, T$  the  $p$ -value  $p_{(j)}$  is significant at level  $\frac{j\alpha}{T}$ , then reject the global  $H_0$ .

Provides a test for the global hypothesis  $H_0$ , but does not “localise” the signal at a particular  $x_t$

One can, however, devise a sequential procedure to “localise” Sime’s procedure, at the expense of lower power for the global hypothesis  $H_0$ :

### Hochberg’s procedure (assuming independence)

- 1 Suppose we reject  $H_{0,j}$  for small values of  $p_j$
- 2 Order  $p$ -values from most to least significant:  $p_{(1)} \leq \dots \leq p_{(T)}$
- 3 Starting from  $j = T, T - 1, \dots$  and down, accept all  $H_{0,(j)}$  such that  $p_{(j)}$  insignificant at level  $\alpha/(T - j + 1)$ .
- 4 Stop accepting for the first  $j$  such that  $p_{(j)}$  is significant at level  $\alpha/j$ , and reject all the remaining ordered hypotheses past that  $j$  going down.

Genuine improvement over Holm-Bonferroni both overall ( $H_0$ ) and in terms of signal localisation:

- 1 Rejects “more” individual hypotheses than Holm-Bonferroni
- 2 Power for overall  $H_0$  “weaker” than Sime’s (for  $T > 2$ ), much “stronger” than Holm (for  $T > 1$ ).

# Bonferroni, Hochberg, Simes

