

# Statistics 2

Hadrien Montanelli

## 1 Estimation

**Delta method** Suppose  $X_1, \dots, X_n$  are iid with mean  $\mu$  and variance  $\sigma^2$ . The CLT yields

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \stackrel{d}{\approx} N(0, 1).$$

Suppose we need the asymptotic distribution of  $g(\bar{X})$ , for some function  $g$ . Expanding  $g(\bar{X})$  yields

$$g(\bar{X}) \approx g(\mu) + (\bar{X} - \mu)g'(\mu),$$

*asymptotic mean and asymptotic variance*

$$E(g(\bar{X})) \approx g(\mu) \quad \text{and} \quad \text{var}(g(\bar{X})) \approx \frac{g'(\mu)^2 \sigma^2}{n},$$

and *asymptotic distribution*

$$g(\bar{X}) \stackrel{d}{\approx} N\left(g(\mu), \frac{g'(\mu)^2 \sigma^2}{n}\right).$$

**Order statistics (data)** The *order statistics* of data  $x_1, \dots, x_n$  are their values in increasing order, which we denote  $x_{(1)} \leq \dots \leq x_{(n)}$ .

**Sample median** The *sample median* is  $m = x_{([n+1]/2)}$  if  $n$  is odd, or  $m = 1/2(x_{(n/2)} + x_{(n/2+1)})$  if  $n$  is even.

**Lower and upper quartile** The *lower quartile* has 1/4 of the sample that is less than it, and the *upper quartile* has 3/4 of the sample that is less than it.

**Inter-quartile range (IQR)** The *inter-quartile range* is defined by  $\text{IQR} = \text{upper quartile} - \text{lower quartile}$ .

**Order statistic (random sample)** The  $r$ th *order statistic* of the random sample  $X_1, \dots, X_n$  is the random variable  $X_{(r)}$ , where  $X_{(1)} \leq \dots \leq X_{(n)}$ .

**Pdf of order statistic** Suppose  $X_1, \dots, X_n$  are iid and continuous, each having cdf  $F$  and pdf  $f$ . Then, the pdf of  $X_{(r)}$  is

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

**Quantile** For a distribution with cdf  $F$  and pdf  $f$ , the  $p$ th *quantile* is the value  $x_p$  such that

$$F(x_p) = \int_{-\infty}^{x_p} f(u) du = p, \quad 0 \leq p \leq 1.$$

**Probability integral transform** Suppose  $X$  is a continuous random variable taking values in  $(a, b)$ , with strictly increasing cdf  $F$ . Then,  $F(X) \sim U(0, 1)$  is called the *probability integral transform* of  $X$ .

**Lemma** If  $U_{(1)}, \dots, U_{(n)}$  are the order statistics of a random sample of size  $n$  from a  $U(0, 1)$  distribution, then

$$E(U_{(r)}) = \frac{r}{n+1},$$
$$\text{var}(U_{(r)}) = \frac{r}{(n+1)(n+2)} \left(1 - \frac{r}{n+1}\right).$$

**Q-Q plot** If data  $x_1, \dots, x_n$  are from a distribution with cdf  $F$ , then Q-Q plots use the approximation

$$F(x_{(k)}) \approx \frac{k}{n+1}.$$

**Normal Q-Q plot** If data  $x_1, \dots, x_n$  are from a  $N(\mu, \sigma^2)$  distribution, for some unknown  $\mu$  and  $\sigma^2$ , then

$$x_{(k)} \approx \sigma \phi^{-1}\left(\frac{k}{n+1}\right) + \mu.$$

**Exponential Q-Q plot** If data  $x_1, \dots, x_n$  are from a  $\text{Exp}(\mu)$  distribution, for some unknown  $\mu$ , then

$$x_{(k)} \approx -\mu \log\left(1 - \frac{k}{n+1}\right).$$

**Pareto Q-Q plot** If data  $x_1, \dots, x_n$  are from a  $\text{Par}(\alpha, \theta)$  distribution, for some unknown  $\alpha$  and  $\theta$ , then

$$\log x_{(k)} \approx \log \alpha - \frac{1}{\theta} \log\left(1 - \frac{k}{n+1}\right).$$

**Observed/Fisher information** In a model with scalar  $\theta$  and log-likelihood  $\ell(\theta; \mathbf{X})$ , the *observed information*  $J(\theta)$  and the *Fisher information*  $I(\theta)$  are

$$J(\theta) = -\frac{d^2 \ell}{d\theta^2}, \quad I(\theta) = E\left(-\frac{d^2 \ell}{d\theta^2}\right).$$

When  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)$ , the *observed information matrix* and the *Fisher information matrix* are  $p \times p$  symmetric matrices  $J(\boldsymbol{\theta})$  and  $I(\boldsymbol{\theta})$  whose  $(j, k)$  elements are

$$J(\boldsymbol{\theta})_{j,k} = -\frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k}, \quad I(\boldsymbol{\theta})_{j,k} = E\left(-\frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k}\right).$$

Note that expectations are taken over  $\mathbf{X}$ .

**Fisher information (iid)** If  $X_1, \dots, X_n$  are iid from  $f(\mathbf{x}; \boldsymbol{\theta})$ , then  $I(\boldsymbol{\theta}) = ni(\boldsymbol{\theta})$ , where  $i(\boldsymbol{\theta})$  is the Fisher information in a sample of size 1; that is, in the scalar case,

$$I(\theta) = ni(\theta), \quad i(\theta) = E\left(-\frac{d^2 \log f(X_1; \theta)}{d\theta^2}\right).$$

**Properties of MLEs (scalar case)**

- $\sqrt{I(\hat{\theta})}(\hat{\theta} - \theta) \xrightarrow{d} N(0, 1)$ , that is,  $\hat{\theta} \stackrel{d}{\approx} N(\theta, I(\theta)^{-1})$ ;
- $\hat{\theta} \xrightarrow{p} \theta$  as  $n \rightarrow \infty$ ;
- $\hat{\theta}$  is asymptotically unbiased;
- the variance of  $\hat{\theta}$  is  $\sim 1/n$ , and as small as possible;
- if  $\psi = g(\theta)$  then  $\hat{\psi} = g(\hat{\theta})$ .

## 2 Confidence intervals

**CI using MLE** Using  $\sqrt{I(\hat{\theta})}(\hat{\theta} - \theta) \stackrel{d}{\approx} N(0, 1)$ ,

$$P\left(-z_{\alpha/2} < \sqrt{I(\hat{\theta})}(\hat{\theta} - \theta) < z_{\alpha/2}\right) \approx 1 - \alpha.$$

One can either solve for  $\theta$ , or use  $I(\theta) \approx I(\hat{\theta})$ .

**Student  $t$ -distribution** Let  $Z \sim N(0, 1)$  and  $Y \sim \chi_n^2$  be independent. We say that  $T = Z/(\sqrt{Y/n})$  has a *student  $t$ -distribution*, and write  $T \sim t_n$ .

**Independence of  $\bar{X}$  and  $S$**  Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ . Then,  $\bar{X}$  and  $S$  are independent, and

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right), \quad \frac{(n-1)S}{\sigma^2} \sim \chi_{n-1}^2.$$

**Pivot** A *pivot* is a random variable, function of both  $\mathbf{X}$  and  $\theta$ , whose distribution does not depend on  $\theta$ .

**Pivot (normal)** Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ . Examples of pivot include

$$\frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}, \quad \frac{(n-1)S}{\sigma^2} \sim \chi_{n-1}^2.$$

**CI using pivot** Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ . We have

$$P\left(-t_{n-1}(\alpha/2) < \frac{\bar{X} - \mu}{S/\sqrt{n}} < t_{n-1}(\alpha/2)\right) = 1 - \alpha,$$

where  $t_{n-1}(\alpha/2)$  satisfies  $P(t_{n-1} > t_{n-1}(\alpha/2)) = \alpha/2$ .

## 3 Hypothesis testing

**General setup** Let  $X_1, \dots, X_n$  be a random sample from  $f(x; \theta)$  where  $\theta \in \Theta$  is a scalar or vector parameter. Suppose we are interested in testing, for  $\Theta_0 \cap \Theta_1 = \emptyset$ ,

- the null hypothesis  $H_0: \theta \in \Theta_0$ ;
- against the alternative hypothesis  $H_1: \theta \in \Theta_1$ .

Consider a statistic  $t(\mathbf{X})$  such that large values of  $t(\mathbf{X})$  cast doubt on  $H_0$ , and let  $t_{obs} = t(\mathbf{x})$  be the value observed. Then, the *p-value* is  $p = P(t(\mathbf{X}) \geq t_{obs} | H_0)$ .

**Critical region** The *critical region*  $C$  is such that if we have  $\mathbf{x} \in C \subset \mathbb{R}^n$ , we reject  $H_0$ , and keep it otherwise.

### Errors in hypothesis testing

- type I error: rejecting  $H_0$  when  $H_0$  is true;
- type II error: not rejecting  $H_0$  when  $H_0$  is false.

**Size and power** The *type I error probability*, also called the *size*, is

$$\alpha = P(\text{reject } H_0 | H_0 \text{ true}) = \sup_{\theta \in \Theta_0} P(\mathbf{X} \in C | \theta),$$

while the *type II error probability* is

$$\beta(\theta) = P(\text{don't reject } H_0 | \text{true value is } \theta) = P(\mathbf{X} \notin C | \theta),$$

and  $w(\theta) = 1 - \beta(\theta) = P(\mathbf{X} \in C | \theta)$  is called the *power*. We want  $w(\theta) \approx 1$  for  $\theta \in \Theta_1$ , and  $w(\theta) \approx 0$  for  $\theta \in \Theta_0$ .

**Neyman–Pearson lemma** Consider testing  $H_0: \theta = \theta_0$  against  $H_1: \theta = \theta_1$ . Define the critical region  $C$  by

$$C = \left\{ \mathbf{x} : \frac{L(\theta_0; \mathbf{x})}{L(\theta_1; \mathbf{x})} \leq k \right\},$$

and suppose  $k$  and  $\alpha$  are such that  $P(\mathbf{X} \in C | H_0) = \alpha$ . Then, among all tests of size  $\leq \alpha$ , the test with critical region  $C$  has maximum power.

**Uniformly most powerful** Consider testing  $H_0: \theta = \theta_0$  against  $H_1: \theta \in \Theta_1$ . If the critical region is the same for all  $\theta_1 \in \Theta_1$ , then  $C$  is said to be *uniformly most powerful*.

**Likelihood ratio** Consider testing  $H_0: \theta \in \Theta_0$  against  $H_1: \theta \in \Theta \setminus \Theta_0$ . The *likelihood ratio* is defined by

$$\lambda(\mathbf{x}) = \frac{\sup_{\theta \in \Theta_0} L(\theta; \mathbf{x})}{\sup_{\theta \in \Theta} L(\theta; \mathbf{x})}.$$

A *likelihood ratio test* (LRT) of  $H_0$  against  $H_1$  has critical region  $C = \{\mathbf{x} : \lambda(\mathbf{x}) \leq k\}$ . For a test of size  $\alpha$ , we must choose  $k$  such that  $\sup_{\theta \in \Theta_0} P(\lambda(\mathbf{X}) \leq k | \theta) = \alpha$ .

**Likelihood ratio statistic** The *likelihood ratio statistic* is defined by  $\Lambda(\mathbf{X}) = -2 \log \lambda(\mathbf{X})$ . For  $\Lambda$ , the LRT has critical region  $C = \{\mathbf{x} : \Lambda(\mathbf{x}) \geq k\}$ . If  $H_0$  is true, then

$$\Lambda(\mathbf{x}) \xrightarrow{d} \chi_p^2 \text{ as } n \rightarrow \infty, \quad p = \dim \Theta - \dim \Theta_0.$$

For a size- $\alpha$  test, we choose  $k$  such that  $P(\chi_p^2 \geq k) = \alpha$ .

**Goodness of fit tests** Consider  $n$  independent observations of categories  $i = 1, \dots, k$ . Let  $n_i$  be the number of observations in category  $i$ , with  $\sum_i n_i = n$ , and  $\pi_i$  the probability of being category  $i$ , with  $\sum_i \pi_i = 1$ . We test

- the null hypothesis  $H_0: \pi_i = \pi_i(\theta)$ , where  $\theta \in \Theta$ , and  $\dim H_0 = q < k - 1$ ;
- against the general alternative  $H_1$ : the  $\pi_i$  are unrestricted except for  $\sum_i \pi_i = 1$ , and  $\dim H_1 = k - 1$ .

In this case, the likelihood ratio statistic is given by

$$\Lambda = -2 \log \frac{\sup_{H_0} L(\theta; \mathbf{x})}{\sup_{H_1} L(\theta; \mathbf{x})} = 2 \sum_{i=1}^k n_i \log \left( \frac{n_i}{n_i \pi_i(\hat{\theta})} \right),$$

and  $\Lambda \approx \chi_{k-1-q}^2$  when  $H_0$  is true.

## 4 Bayesian inference

**Notation** We write  $f(\mathbf{x}|\theta)$  instead of  $f(\mathbf{x};\theta)$  to emphasize that we have a model for data  $\mathbf{x}$  given the value  $\theta$ .

**Prior** We summarize our beliefs about  $\theta$  in a *prior* pmf/pdf  $\pi(\theta)$ ; we treat  $\theta$  as a random variable.

**Posterior** Using the continuous Bayes' theorem,

$$f_{Z|Y}(z|y) = \frac{f_{Z|Y}(y|z)f_Z(z)}{f_Y(y)},$$

we define the *posterior* pmf/pdf as

$$\begin{aligned}\pi(\theta|\mathbf{x}) &\propto f(\mathbf{x}|\theta) \times \pi(\theta), \\ \text{posterior} &\propto \text{likelihood} \times \text{prior}.\end{aligned}$$

### Posterior summaries

- the *posterior mode* ( $\theta$  that maximizes  $\pi(\theta|\mathbf{x})$ );
- the *posterior mean* (expectation over  $\theta$ );
- the *posterior variance* (variance over  $\theta$ );
- the *posterior median* (satisfies  $\int_{-\infty}^m \pi(\theta|\mathbf{x})d\theta = 1/2$ );
- other quantiles of  $\pi(\theta|\mathbf{x})$ .

**Credible interval** A  $100(1 - \alpha)\%$  *credible interval* for  $\theta$  is an interval  $(\theta_a, \theta_b)$  such that

$$P(\theta_a \leq \theta \leq \theta_b|\mathbf{x}) = \int_{\theta_a}^{\theta_b} \pi(\theta|\mathbf{x})d\theta = 1 - \alpha.$$

If  $P(\theta \leq \theta_a|\mathbf{x}) = P(\theta \geq \theta_b|\mathbf{x}) = \alpha/2$ , then the interval is called *equal-tailed*.

**Highest posterior density** A credible interval  $I$  is called a *highest posterior density* interval if  $\pi(\theta|\mathbf{x}) \geq \pi(\theta'|\mathbf{x})$ , for all  $\theta \in I$  and  $\theta' \notin I$ . (For the normal, every equal-tailed credible interval is a highest posterior density interval.)

**Posterior predictive density** Let  $X_{n+1}$  represent a future observation, independent of  $X_1, \dots, X_n$ , and let  $\mathbf{x} = (x_1, \dots, x_n)$  denote the observed data. The pdf of  $X_{n+1}$ , the *posterior predictive density*, is defined by

$$f(x_{n+1}|\mathbf{x}) = \int_{\Theta} f(x_{n+1}, \theta|\mathbf{x})d\theta = \int_{\Theta} f(x_{n+1}|\theta)\pi(\theta|\mathbf{x})d\theta,$$

where we used  $P(A \cap B|C) = P(A|B \cap C)P(B|C)$ .

**Proper prior** A prior is *proper* if  $\int \pi = 1$ , and is *improper* if the integral cannot be normalized.

**Jeffreys prior** The *Jeffreys prior* is defined by

$$\pi(\theta) \propto I(\theta)^{1/2},$$

where  $I(\theta)$  is the expected information.

**Prior/posterior odds** Suppose we want to compare two hypotheses  $H_0$  and  $H_1$ , exactly one of which is true. The *prior* and *posterior odds* of  $H_0$  relative to  $H_1$  are

$$\text{prior odds} = \frac{P(H_0)}{P(H_1)}, \quad \text{posterior odds} = \frac{P(H_0|\mathbf{x})}{P(H_1|\mathbf{x})}.$$

**Bayes factor** The *Bayes factor*  $B_{01}$  is defined via

$$\frac{P(H_0|\mathbf{x})}{P(H_1|\mathbf{x})} = \frac{P(\mathbf{x}|H_0)}{P(\mathbf{x}|H_1)} \times \frac{P(H_0)}{P(H_1)},$$

$$\text{posterior odds} = \text{Bayes factor} \times \text{prior odds}.$$

**General setup** We are assuming we have:

- prior probabilities  $P(H_0)$  and  $P(H_1)$ , which satisfy  $P(H_0) + P(H_1) = 1$ ;
- prior distributions for  $\theta_0 \in \Theta_0$  and  $\theta_1 \in \Theta_1$  under  $H_0$  and  $H_1$ , which we write as  $\pi(\theta_0|H_0)$  and  $\pi(\theta_1|H_1)$ ;
- models for data  $\mathbf{x}$  under  $H_0$  and  $H_1$ , which we write as  $f(\mathbf{x}|\theta_0, H_0)$  and  $f(\mathbf{x}|\theta_1, H_1)$ .

**Computing the Bayes factor** Use

$$P(H_i) = \int_{\Theta_i} \pi(\theta_i)d\theta_i \quad \text{and} \quad P(H_i|\mathbf{x}) = \int_{\Theta_i} \pi(\theta_i|\mathbf{x})d\theta_i.$$

Alternatively, use

$$\pi(\theta_i|H_i) = \pi(\theta_i) \Big/ \int_{\Theta_i} \pi(\theta_i)d\theta_i,$$

and

$$P(\mathbf{x}|H_i) = \int_{\Theta_i} f(\mathbf{x}|\theta_i, H_i)\pi(\theta_i|H_i)d\theta_i.$$

Note that if  $H_i : \theta = \theta_i$ , then  $P(\mathbf{x}|H_i) = f(\mathbf{x}|\theta_i)$ .

**Assessing evidence** The quantity  $2 \log B_{01}$  is used to summarize the evidence for  $H_0$  compared to  $H_1$ :

$B_{01}$	$2 \log B_{01}$	Evidence for $H_0$
$< 1$	$< 0$	Negative
$1 - 3$	$0 - 2$	Hardly worth a mention
$3 - 20$	$2 - 6$	Positive
$20 - 150$	$6 - 10$	Strong
$> 150$	$> 10$	Very strong

**Asymptotic normality** Let  $\tilde{\ell}(\theta) = \log \pi(\theta|\mathbf{x})$  and  $\tilde{\theta}$  be the posterior mode. Expanding  $\tilde{\ell}(\theta)$  yields

$$\tilde{\ell}(\theta) \approx \tilde{\ell}(\tilde{\theta}) + (\theta - \tilde{\theta})\tilde{\ell}'(\tilde{\theta}) + \frac{1}{2}(\theta - \tilde{\theta})^2\tilde{\ell}''(\tilde{\theta}),$$

with  $\tilde{\ell}'(\tilde{\theta}) = 0$ . Therefore,

$$\pi(\theta|\mathbf{x}) = \exp(\tilde{\ell}(\theta)) \propto \exp\left(-\frac{1}{2}\tilde{J}(\tilde{\theta})(\theta - \tilde{\theta})^2\right),$$

i.e.,  $\theta|\mathbf{x} \stackrel{d}{\approx} N(\tilde{\theta}, \tilde{J}(\tilde{\theta})^{-1})$ . (Similarly,  $\theta|\mathbf{x} \stackrel{d}{\approx} N(\hat{\theta}, J(\hat{\theta})^{-1})$ .)