

STATISTICS The "science" of processing quantitative evidence to

- make judgements or decisions. It can include
 - Design of experiments
 - Exploring data using graphics
 - Quantifying evidence for or against a hypothesis
 - Clear communication of evidence and uncertainty

Def A parametric family is a family of probability distribution indexed by a scalar or vector parameter θ in some parameter space

Ⓜ. Ex Family of $N(\mu, \sigma^2)$ distributions, $\theta = (\mu, \sigma^2) \in \mathbb{H} = \mathbb{R}^2$

- In this course we will talk about formal inference of θ , and how to test hypotheses about θ (e.g. $H: \theta = 0$)

Probability review

Let Ω be a sample space of all possible outcomes of an experiment.

An event is a measurable subset of Ω . Let \mathcal{F} be the set of all events. A probability measure $IP: \mathcal{F} \rightarrow [0, 1]$ satisfies

$$\textcircled{1} IP(\emptyset) = 0 \quad \textcircled{2} IP(\Omega) = 1 \quad \textcircled{3} IP\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} IP(A_i)$$

A random variable is a measurable function $X: \Omega \rightarrow \mathbb{R}$.

- Ex 2 coin tosses, $\Omega = \{HH, HT, TH, TT\}$

$$X(\omega) = \# \text{ heads in } \omega$$

The distribution function of X is

$$F_X(x) = IP(X \leq x)$$

We'll use f_X to denote either

ⓐ The probability mass f_X^n of X if X is discrete, $f_X(x) = IP(X=x)$

ⓑ The pdf of X if X is continuous, i.e. $F_X^+(x) = \int_{-\infty}^x f(t) dt$

The expectation of X is

$$\bullet EX = \begin{cases} \sum x f_X(x) & \text{if } X \text{ discrete} \\ \int_{-\infty}^{\infty} x f_X(x) dx & \text{if } X \text{ continuous} \end{cases}$$

Variance is $\text{Var}(X) = \mathbb{E}(X - \mathbb{E}(X))^2$

Linearity of expectation $\mathbb{E}(X+Y) = \mathbb{E}(X) + \mathbb{E}(Y)$, $\mathbb{E}(\lambda X) = \lambda \mathbb{E}(X)$

We say that X_1, \dots, X_n are independent if

$$P(X_1 \leq x_1, \dots, X_n \leq x_n) = P(X_1 \leq x_1) \cdots P(X_n \leq x_n)$$

If X_1, \dots, X_n are independent with pdfs (f_{X_1}, \dots, f_{X_n}) or pmfs then the joint pdf or pmf for the random vector $\underline{X} = (X_1, \dots, X_n)$ is

$$f_{\underline{X}}(\underline{x}) = \prod f_{X_i}(x_i)$$

For independent X_1, \dots, X_n ,

$$\text{Var}(a_1 X_1 + a_2 X_2 + \dots + a_n X_n) = a_1^2 \text{Var}(X_1) + \dots + a_n^2 \text{Var}(X_n)$$

Standardisation Let X_1, \dots, X_n be iid with $\mathbb{E}X_1 = \mu$, $\text{Var}X_1 = \sigma^2$

We usually write $S_n = \sum_{i=1}^n X_i$. By linearity $\mathbb{E}S_n = n\mu$, $\text{Var}S_n = n\sigma^2$

The empirical mean $\bar{X}_n = S_n/n$ has $\mathbb{E}\bar{X}_n = \mu$, $\text{Var}\bar{X}_n = \sigma^2/n$

A standardised statistic $Z_n = \sqrt{n}(\bar{X}_n - \mu)/\sigma$

has $\mathbb{E}Z_n = 0$, $\text{Var}Z_n = 1$.

Convergence

Weak law of large numbers: $\forall \varepsilon > 0$, $P(|\bar{X}_n - \mu| > \varepsilon) \rightarrow 0$ as $n \rightarrow \infty$

Strong law of large numbers: $P(\bar{X}_n \rightarrow \mu) = 1$

Cheese Lettuce Bacon: Z_n is approximately $N(0,1)$ when n is large

Maxima (or minima) of iid vvs If X_1, \dots, X_n are iid, $Y = \max(X_1, \dots, X_n)$

$$\begin{aligned} \text{then } F_Y(y) &= P(Y \leq y) = P(X_1 \leq y, \dots, X_n \leq y) \\ &= (F_{X_1}(y))^n \end{aligned}$$

Can derive pdf or pmf by differentiation or differencing.

Moment generating function $M_X(t) = \mathbb{E}(e^{tX})$ provided the expectation exists in some neighbourhood of 0.

We can use this to find moments, $\mathbb{E}(X^n) = M_X^{(n)}(0)$.

Also useful to find distro of a sum of vvs, because under broad conditions, $M_X = M_Y \Rightarrow F_X = F_Y$

L1.3

Ex $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\mu)$ Distribution of S_n ? A

$$\begin{aligned} M_{S_n}(t) &= \mathbb{E}(e^{tS_n}) = \mathbb{E}(e^{t(X_1 + \dots + X_n)}) \\ &= \prod_{i=1}^n \mathbb{E}(e^{tX_i}) = (\mathbb{E}e^{tX_1})^n \\ &= M_{X_1}(t)^n = e^{-n\mu(1-e^t)} \end{aligned}$$

$\Rightarrow S_n \sim \text{Poisson}(n\mu)$

Conditioning If X, Y have a joint pdf (or pmf) $f_{X,Y}$, then the marginal pdf of X is

$$\begin{aligned} f_X(x) &= \int_{-\infty}^{\infty} f_{X,Y}(x,y) dy \\ (f_X(x) &= \sum_y f_{X,Y}(x,y)) \end{aligned}$$

The conditional pdf of X given $Y=y$ is

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)} \quad \text{if } f_Y(y) \neq 0$$

Conditional expectation $\mathbb{E}(X|Y=y) = \int x f_{X|Y}(x|y) dx$
 $(\sum_x x f_{X|Y}(x|y))$

We usually write simply $\mathbb{E}(X|Y)$, a rv which is a function of Y . $\text{Var}(X|Y) = \mathbb{E}[(X - \mathbb{E}(X|Y))^2 | Y]$

Tower property $\mathbb{E}[\mathbb{E}[X|Y]] = \mathbb{E}X$

Corollary $\text{Var}(X) = \text{Var}(\mathbb{E}(X|Y)) + \mathbb{E}[\text{Var}(X|Y)]$

Probability review (cont.)

1) Change of variable Let X, Y be rvs, continuous. Let $(X, Y) \mapsto (U, V)$ be a diff. bijection. Then U, V have joint pdf

$$f_{U,V}(u,v) = f_{X,Y}(x(u,v), y(u,v)) |\det J|$$

where

$$J = \begin{pmatrix} \partial x / \partial u & \partial x / \partial v \\ \partial y / \partial u & \partial y / \partial v \end{pmatrix}$$

Important distributions

① Binomial: $X \sim \text{Bin}(n, p)$, $n \in \mathbb{N}$, $p \in [0, 1]$

$$P(X=x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, 1, \dots, n$$

$X = \#$ successes in n indep Bernoulli trials

$$EX = np, \quad \text{Var } X = np(1-p)$$

② Poisson: $X \sim \text{Poi}(\lambda)$, $\lambda > 0$

$$P(X=x) = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x = 0, 1, 2, \dots$$

$\text{Poi}(\lambda)$ is the limit of $\text{Bin}(n, p)$ as $n \rightarrow \infty$, keeping $np = \lambda$

In a Poisson process, the $\#$ of events happening in an interval of length t has a $\text{Poi}(\lambda t)$ distribution and disjoint intervals are indep

③ Negative Binomial: $X \sim \text{NB}(k, p)$, $k \in \mathbb{N}$, $p \in [0, 1]$

$$P(X=x) = \binom{x-1}{k-1} (1-p)^{x-k} p^k, \quad x = k, k+1, \dots$$

$X = \#$ trials until k^{th} success in a sequence of iid Bernoulli trials

$$EX = \frac{k}{p}, \quad \text{Var } X = \frac{k(1-p)}{p^2}$$

④ Multinomial: n indep trials with k possible outcomes, with probabilities p_1, \dots, p_k

$(N_1, \dots, N_k) \sim \text{Multin}(n, \underline{p})$ if $N_i = \#$ trials with outcome i

$$P(N_1=n_1, \dots, N_k=n_k) = \frac{n!}{n_1! \dots n_k!} p_1^{n_1} \dots p_k^{n_k}$$

for $n_1, \dots, n_k \in \{0, \dots, n\}$ with $\sum n_i = n$

L2.2

⑤ Normal: $X \sim N(\mu, \sigma^2)$, $\mu \in \mathbb{R}$, $\sigma^2 \in (0, \infty)$

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad x \in \mathbb{R}$$

We call $N(0,1)$ the standard normal. Denote its pdf ϕ , distribution function Φ .

⑥ Uniform: $X \sim \text{Unif}(a,b)$, $a, b \in \mathbb{R}$, $a < b$

$$f_X(x) = \frac{1}{b-a}, \quad x \in [a,b]$$

$$\mathbb{E}X = \frac{a+b}{2}, \quad \text{Var}X = \frac{(b-a)^2}{12}$$

⑦ Gamma: $X \sim \Gamma(\alpha, \lambda)$, $\alpha > 0$, $\lambda > 0$

$$f_X(x) = \frac{\lambda^\alpha x^{\alpha-1} e^{-\lambda x}}{\Gamma(\alpha)}, \quad x \geq 0$$

$$\mathbb{E}X = \frac{\alpha}{\lambda}, \quad \text{Var}X = \frac{\alpha}{\lambda^2}, \quad M_X(t) = \left(\frac{\lambda}{\lambda-t}\right)^\alpha$$

If $X_i \sim \Gamma(\alpha_i, \lambda)$ are indep, then

- For any $t \in \mathbb{R}$, $tX_1 \sim \Gamma(\alpha_1, \frac{\lambda}{t})$

- $\sum_{i=1}^n X_i \sim \Gamma\left(\sum_{i=1}^n \alpha_i, \lambda\right)$

Show using MGF (trivial)

⑧ Exponential: $X \sim \text{Exp}(\lambda) = \Gamma(1, \lambda)$

Memoryless, distribution of waiting times in a Poisson process with rate λ

⑨ Chi squared: $X \sim \chi_k^2$, $k \in \mathbb{N}$, degrees of freedom
 $= \Gamma\left(\frac{k}{2}, \frac{1}{2}\right)$

If $Z_1, \dots, Z_n \stackrel{\text{iid}}{\sim} N(0,1)$, then $\sum_{i=1}^n Z_i^2 \sim \chi_n^2$

⑩ Beta: $X \sim \text{Beta}(\alpha, \beta)$, $\alpha, \beta > 0$

$$f_X(x) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)} \quad \text{for } x \in [0,1]$$

$$\mathbb{E}X = \frac{\alpha}{\alpha+\beta}, \quad \text{Var}X = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$$

Estimation

A Suppose X_1, \dots, X_n are iid, each has pdf $f_X(x|\theta)$ with θ unknown.

An estimator of θ is a function of the data X_1, \dots, X_n (a statistic)

$$\hat{\theta} = T(X_1, \dots, X_n).$$

As $\hat{\theta}$ is a function of rvs, it is itself a r.v. Its distribution is called the sampling distribution of θ .

Ex $X_1, \dots, X_n \sim N(\mu, 1)$, $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i \sim N(\mu, \frac{1}{n})$

Def The bias of $\hat{\theta}$ is $\text{bias}(\hat{\theta}) = \mathbb{E}_{\theta}(\hat{\theta}) - \theta$

b

In the above example, $\mathbb{E}_{\mu}(\hat{\mu}) = \mu$.
(expectation wrt joint pdf $\prod_{i=1}^n f_{X_i}(x_i|\theta)$)

So $\text{bias}(\hat{\mu}) = 0$, we say $\hat{\mu}$ is unbiased for μ .

Def The mean squared error of $\hat{\theta}$ is $\text{mse}(\hat{\theta}) = \mathbb{E}_{\theta}((\hat{\theta} - \theta)^2)$

(a function of θ) wh

$$\begin{aligned} \mathbb{E}_{\theta}((\hat{\theta} - \theta)^2) &= \mathbb{E}_{\theta}((\hat{\theta} - \mathbb{E}_{\theta}\hat{\theta} + \mathbb{E}_{\theta}\hat{\theta} - \theta)^2) \\ &= \mathbb{E}_{\theta}((\hat{\theta} - \mathbb{E}_{\theta}\hat{\theta})^2) + \mathbb{E}_{\theta}((\mathbb{E}_{\theta}\hat{\theta} - \theta)^2) \\ &\quad + 2\mathbb{E}_{\theta}(\hat{\theta} - \mathbb{E}_{\theta}\hat{\theta})\mathbb{E}_{\theta}(\mathbb{E}_{\theta}\hat{\theta} - \theta) \\ &= \text{var}(\hat{\theta}) + \text{bias}(\hat{\theta})^2 \end{aligned}$$

NB Sometimes a biased estimator can have lower MSE than an unbiased one. There is a tradeoff between bias and variance.

Ex Two different estimators of a binomial mean

Let $X \sim \text{Bin}(n, \theta)$, θ unknown

① $T_U = \frac{X}{n}$

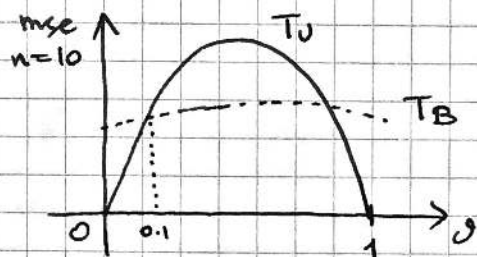
② $T_B = \frac{1}{n+2}X + \frac{1}{n+2}$

$\mathbb{E}_{\theta} T_U = \frac{n\theta}{n} = \theta$, so unbiased

$\text{Var}_{\theta} T_U = \frac{n\theta(1-\theta)}{n^2} = \frac{1}{n}\theta(1-\theta)$

$\mathbb{E}_{\theta} T_B = \frac{n}{n+2}\theta + \frac{1}{n+2}$, so biased for $\theta \neq \frac{1}{2}$

$\text{Var}_{\theta} T_B = \frac{n\theta(1-\theta)}{(n+2)^2}$, lower variance ...



L2.4

T_B has a lower mse for much of the range of θ . If we know θ is likely to be close to $\frac{1}{2}$, it makes sense to use the biased estimator.

In general, prior knowledge will change our choice of estimator.

Sufficiency

● As usual, let $X = (X_1, \dots, X_n)$ be iid, from some distribution with parameter θ .

Def A statistic T is sufficient for θ if the conditional distribution of the data X given T does not depend on θ .

NB T and θ can be vectors, not necessarily of the same dimension. In a sense, a sufficient statistic summarises all information in X relevant for estimating θ .

Ex Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\theta)$, so $f_X(x|\theta) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i}$

● This depends on ~~θ~~ only through

$T(X) = \sum_{i=1}^n X_i$. Note $T(X) \sim \text{Bin}(n, \theta)$.

$$f_{X|T(X)=t}(x|T=t) = \frac{P(X=x, T=t)}{P(T=t)} = \frac{\theta^{\sum_{i=1}^n x_i} (1-\theta)^{n-\sum_{i=1}^n x_i}}{\binom{n}{t} \theta^t (1-\theta)^{n-t}}$$

$$= \left[\binom{n}{t} \right]^{-1} \quad \text{zero if } T(x) \neq t$$

This does not depend on θ , hence $T(X) = \sum_{i=1}^n X_i$ is sufficient.

Thm (Factorisation criterion) T is sufficient for θ iff

$$f_X(x|\theta) = g(T(x), \theta) h(x)$$

● for some functions g, h .

$$\text{Ex (cont.) } f_X(x|\theta) = \underbrace{\theta^{\sum_{i=1}^n x_i} (1-\theta)^{n-\sum_{i=1}^n x_i}}_{g(T(x), \theta)} \cdot \underbrace{1}_{h(x)}$$

Hence $\sum_{i=1}^n X_i$ is sufficient (again).

Proof [Discrete case] Suppose $f_X(x|\theta) = g(T(x), \theta) h(x)$.

$$\text{Then } f_{X|T=t}(x|T=t) = \frac{P_\theta(X=x, T(X)=t)}{P_\theta(T(X)=t)} = \frac{g(t, \theta) h(x)}{\sum_{x', T(x')=t} g(t, \theta) h(x')}$$

$$= \frac{h(x)}{\sum_{x', T(x')=t} h(x')},$$

● which does not depend on θ , so $T(X)$ is sufficient.

Conversely, if $T(X)$ is sufficient, then

$$\begin{aligned}
 \mathbb{P}_\theta(X=x) &= \mathbb{P}_\theta(X=x, T(X)=T(x)) \\
 &= \underbrace{\mathbb{P}_\theta(X=x \mid T(X)=T(x))}_{\substack{\text{indep of } \theta \text{ by} \\ \text{sufficiency of } T, \text{ so} \\ \text{write as } h(x)}} \underbrace{\mathbb{P}_\theta(T(X)=T(x))}_{g(T(x), \theta)}
 \end{aligned}$$

Ex $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Unif}(0, \theta), \theta > 0$

$$\begin{aligned}
 f_X(x|\theta) &= \prod_{i=1}^n \frac{1}{\theta} \mathbb{1}_{[0, \theta]}(x_i) \\
 &= \frac{1}{\theta^n} \underbrace{\mathbb{1}_{\{\max x_i \leq \theta\}}}_{g(\max x_i, \theta)} \underbrace{\mathbb{1}_{\{\min x_i \geq 0\}}}_{h(x)}
 \end{aligned}$$

Hence $T = \max X_i$ is sufficient for θ .

Remark Sufficient statistics are not unique. In fact, $T(X) = X$ is a sufficient statistic. Consider the equivalence relation $x \sim y$ if $T(x) = T(y)$. Every (sufficient) statistic partitions the sample space into equivalence classes of the form $\{x \in \mathcal{X}^n : T(x) = t\}$. What we want is for the partition to be as coarse as possible.

Def A sufficient statistic $T(X)$ is minimal sufficient if it is a function of any other sufficient statistic, i.e. for $S(X)$ sufficient, \circledast
 $S(x) = S(y) \Rightarrow T(x) = T(y)$.

Alternatively, the partition induced by T is coarser than that of S .

Thm Suppose $T(X)$ is a statistic such that $\frac{f_X(x|\theta)}{f_X(y|\theta)}$ does not depend on θ iff $T(x) = T(y)$.

Then $T(X)$ is minimal sufficient.

Proof (sketch) First we'll prove $T(X)$ is sufficient.

For each value of $T(x)$, let x_t be a representative of $\{x \in \mathcal{X}^n : T(x) = t\}$

By hypothesis, $\frac{f_X(x|\theta)}{f_X(x_t|\theta)}$ does not depend on θ .

L3.3

Then

$$f_X(x'|\theta) = \underbrace{f_X(x_{T(x')}|\theta)}_{\text{call } g(T(x'), \theta)} \frac{f_X(x'(\theta))}{\underbrace{f_X(x_{T(x')}|\theta)}_{\text{call } h(x')}}.$$

By the factorisation criterion, T is sufficient.

Now let S be another sufficient statistic. Then \exists functions g_S, h_S such that $f_X(x|\theta) = g_S(S(x), \theta) h_S(x)$.

Say $S(x) = S(y)$. Then

$$\frac{f_X(x|\theta)}{f_X(y|\theta)} = \frac{g_S(S(x), \theta) h_S(x)}{g_S(S(y), \theta) h_S(y)} = \frac{h_S(x)}{h_S(y)},$$

which does not depend on θ .

Hence by hypothesis $T(x) = T(y)$. So T is minimal. \square

Ex $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

$$\frac{f_X(x|\mu, \sigma^2)}{f_X(y|\mu, \sigma^2)} = \frac{(2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right)}{(2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2\right)}$$

$$= \exp\left\{-\frac{1}{2\sigma^2} \left(\sum_{i=1}^n x_i^2 - \sum_{i=1}^n y_i^2\right) + \frac{\mu}{\sigma^2} \left(\sum_{i=1}^n x_i - \sum_{i=1}^n y_i\right)\right\}$$

This is indep of μ, σ^2 iff $\sum_{i=1}^n x_i = \sum_{i=1}^n y_i$ and $\sum_{i=1}^n x_i^2 = \sum_{i=1}^n y_i^2$.
Therefore $T(X) = \left(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2\right)$ is minimal sufficient. (!)

NB Any bijection of a minimal sufficient statistic is minimal sufficient (only need injection). Therefore, setting

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \text{ get } T(X) = \left(\bar{X}, \sum_{i=1}^n (X_i - \bar{X})^2\right)$$

is also minimal sufficient in (μ, σ^2) .

L4.1

Sufficiency (cont.)

- * Minimal sufficient need not exist, but do under mild conditions
- * " " are unique [If S, T both minimal sufficient, then $S(x) = S(y) \Leftrightarrow T(x) = T(y)$]

Rao-Blackwell theorem

Let T be a sufficient statistic for a parameter θ . Let $\tilde{\theta}$ be an estimator with $E \tilde{\theta}^2 < \infty$ for all θ .

Let $\hat{\theta} = E(\tilde{\theta} | T)$. Then, for all θ ,

$$E((\hat{\theta} - \theta)^2) \leq E((\tilde{\theta} - \theta)^2).$$

- The inequality is strict unless $\tilde{\theta}$ is a function of T .

$$E_{\theta}(\tilde{\theta}) = \int \tilde{\theta}(x) f_X(x|\theta) dx$$

Drop the subscripts LMAO

Proof Note $E \hat{\theta} = E(E(\tilde{\theta} | T)) = E \tilde{\theta}$, so $\hat{\theta}$ and $\tilde{\theta}$ have the same bias.

$$\text{Also, } \text{Var } \tilde{\theta} = \underbrace{E(\text{Var}(\tilde{\theta} | T))}_{> 0 \text{ with +ve prob unless } \tilde{\theta} \text{ a fn of } T} + \underbrace{\text{Var}(E(\tilde{\theta} | T))}_{\text{Var } \hat{\theta}}$$

$$\begin{aligned} \therefore \text{mse } \tilde{\theta} &= \text{bias}(\tilde{\theta})^2 + \text{var } \tilde{\theta} \\ &> \text{bias}(\hat{\theta})^2 + \text{var } \hat{\theta} = \text{mse } \hat{\theta} \end{aligned}$$

with equality iff $\tilde{\theta}$ is a function of T . □

NB Where did we use that T is sufficient?

Sufficiency implies that the conditional distribution of X given T does not depend on θ , therefore $\hat{\theta} = E(\tilde{\theta} | T)$ does not depend on θ , only on the data X . Hence it is a well defined estimator.

Ex Suppose $X_1, \dots, X_n \sim \text{Poi}(\lambda)$, let $\theta = e^{-\lambda} = P(X_1 = 0)$

$$f_X(x|\theta) = \frac{\theta^n (-\log \theta)^{\sum x_i}}{\prod_i x_i!}$$

From factorisation criterion, we see that $T = \sum_{i=1}^n X_i$ is sufficient.

Define a simple estimator of θ , $\tilde{\theta} = \mathbb{1}_{\{X_1=0\}}$. This is unbiased.

Define $\hat{\theta} = \mathbb{E}(\tilde{\theta} | T) = P(X_1=0 | \sum_{i=1}^n X_i = T)$

$$= \frac{P(X_1=0, \sum_{i=1}^n X_i = T)}{P(\sum_{i=1}^n X_i = T)} = \left(\frac{n-1}{n}\right)^T. \quad \begin{array}{l} \sim \text{ugh?} \\ t \text{ vs } T \end{array}$$

So $\hat{\theta} = \left(\frac{n-1}{n}\right)^T$. By Rao-Blackwell, $\text{mse } \hat{\theta} \leq \text{mse } \tilde{\theta} \forall \theta$.

Maximum Likelihood Estimation

Let X_1, \dots, X_n be iid rvs with joint pdf/pmf $f_X(x|\theta)$.

Def The likelihood is $L(\theta) = f_X(x|\theta)$, regarded as a function of θ , given x . The max likelihood estimator (MLE) is

$$\hat{\theta} = \arg \max_{\theta \in \Theta} L(\theta).$$

It is equivalent, and also easier, to maximise the logarithm $l(\theta)$

$$l(\theta) = \log L(\theta) = \log \prod_{i=1}^n f_{X_i}(x_i|\theta) = \sum \log f_{X_i}(x_i|\theta)$$

Ex $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$

$$l(p) = \sum x_i \log p + (n - \sum x_i) \log(1-p)$$

$$dl/dp = \sum x_i / p - (n - \sum x_i) / (1-p)$$

This is zero at $\hat{p} = \sum x_i / n$.

Since $\sum X_i \sim \text{Bin}(n, p)$, $\mathbb{E} \hat{p} = p$ is unbiased.

So the MLE is cool.

Ex 2 $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$, $\theta = (\mu, \sigma^2)$

$$l(\mu, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\partial l / \partial \mu = -\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)$$

$$\partial l / \partial \sigma^2 = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2$$

l maximised when $\partial l / \partial \mu = 0$, $\partial l / \partial \sigma^2 = 0$, this happens at

$$(\hat{\mu}, \hat{\sigma}^2) = (\bar{X}, \frac{1}{n} S_{XX})$$

where $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$, $S_{XX} = \sum_{i=1}^n (X_i - \bar{X})^2$.

Not hard to see $\hat{\mu}$ is unbiased, what about σ^2 ?

L4.3

Later, we'll show that $\frac{S_{XX}}{\sigma^2} = \frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-1}^2$

$$\Rightarrow \frac{n}{\sigma^2} \mathbb{E} \hat{\sigma}^2 = n-1 \Rightarrow \mathbb{E} \hat{\sigma}^2 = \frac{n-1}{n} \sigma^2$$

So $\hat{\sigma}^2$ is biased, but asymptotically unbiased, $\mathbb{E} \hat{\sigma}^2 \xrightarrow{n \rightarrow \infty} \sigma^2$.

It will sometimes be convenient to use the unbiased estimator

$$\tilde{\sigma}^2 = \frac{S_{XX}}{n-1}$$

Ex 3 $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Unif}(0, \theta)$. Then



$$L(\theta) = \frac{1}{\theta^n} \mathbb{1}_{\{\max_i x_i \leq \theta\}}$$

This is zero for $\theta \leq \max_i x_i$, and decreasing in θ for θ larger than $\max_i x_i$. So $\hat{\theta} = \max_i x_i$.

$$F_{\hat{\theta}}(t) = \mathbb{P}(\hat{\theta} \leq t) = \mathbb{P}(X_1 \leq t, \dots, X_n \leq t) = (t/\theta)^n$$

Differentiating, the pdf of $\hat{\theta}$ is

$$f_{\hat{\theta}}(t) = \frac{nt^{n-1}}{\theta^n}$$

for $t \in [0, \theta]$.

$$\mathbb{E} \hat{\theta} = \int_0^\theta \frac{nt^n}{\theta^n} dt = \frac{n}{n+1} \theta$$

$\therefore \hat{\theta}$ is biased, but asymptotically unbiased.

Maximum Likelihood Estimator (cont.)Properties of the MLE

- ① If T is sufficient for θ , then $L(\theta) = g(T(x), \theta) h(x)$, thus the MLE is a function of T .
- ② Invariance If h is injective, the MLE of $\phi = h(\theta)$ is $\hat{\phi} = h(\hat{\theta})$, where $\hat{\theta}$ is the MLE of θ .
- ③ Asymptotic normality (II Principles of Statistics) Under regularity conditions $\sqrt{n}(\hat{\theta} - \theta)$ is approximately $N(0, \Sigma)$, where Σ is "the smallest variance attainable"
- ④ Sometimes the MLE has no closed form expression, but it can be found numerically.

Confidence Intervals

Defⁿ A 100% confidence interval (CI) for θ is a random interval $(A(X), B(X))$ such that $P(A(X) < \theta < B(X)) = \gamma$ for all values of the parameter θ .

Remarks • It is the endpoints which are random, not θ

- We can interpret this interval in terms of "repetitions of the γ experiment", i.e. let $X^{(1)}, X^{(2)}, \dots$ be iid copies of X . Then 100% of the intervals $(A(X^{(1)}), B(X^{(1)})), \dots$ contain θ
- WRONG interpretation: there is a 100% probability that the true parameter θ is in the interval $(A(x), B(x))$

Ex 1 Let X_1, X_2 be iid $\text{Unif}(\theta - \frac{1}{2}, \theta + \frac{1}{2})$. What is a sensible confidence interval for θ ?

$$IP(\min(X_1, X_2) < \theta < \max(X_1, X_2)) = IP(X_1 < \theta < X_2) + IP(X_2 < \theta < X_1) = \frac{1}{2}$$

- Hence $(\min(X_1, X_2), \max(X_1, X_2))$ is a 50% confidence interval. What happens if $|X_1 - X_2| \geq \frac{1}{2}$? e.g. $X_1 = 0.2, X_2 = 0.9$?

L5.2

We can be sure that $\theta \in (\min(X_1, X_2), \max(X_1, X_2))$.

Hence the confidence interval is not a sensible quantifier of uncertainty.

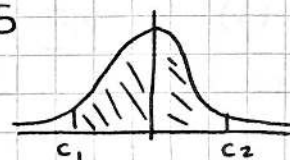
Ex 2 Let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$. What is a 95% confidence interval for θ ?

We know $\sqrt{n}(\bar{X} - \theta) \sim N(0, 1)$, no matter what θ is.

Let c_1, c_2 be such that $\Phi(c_2) - \Phi(c_1) = 0.95$

Then $IP(c_1 < \sqrt{n}(\bar{X} - \theta) < c_2) = 0.95 \quad \forall \theta$.

$$\therefore IP\left(\bar{X} - \frac{c_2}{\sqrt{n}} < \theta < \bar{X} - \frac{c_1}{\sqrt{n}}\right) = 0.95$$



So, $\left(\bar{X} - \frac{c_2}{\sqrt{n}}, \bar{X} - \frac{c_1}{\sqrt{n}}\right)$ is a 95% confidence interval.

We choose c_1, c_2 to make the interval as narrow as possible, this yields $c_2 = -c_1 \cong 1.96$.

Usual procedure ① Find a quantity $R(X, \theta)$ whose distribution does not depend on θ , e.g. $R(X, \theta) = \sqrt{n}(\bar{X} - \theta)$. This is what we call a pivot.

② Write out a statement $IP(c_1 \leq R(X, \theta) \leq c_2) = \gamma$, by choosing c_1, c_2 as appropriate quantiles of the distribution of $R(X, \theta)$.

③ Rearrange to put θ in the middle, i.e. $IP(- \leq \theta \leq -)$

Remark The asymptotic normality of the MLE is useful to find confidence intervals, because we can choose $\sqrt{n}(\hat{\theta} - \theta)$ and apply the above reasoning (when n is large)

Ex 3 $X_1, \dots, X_n \stackrel{iid}{\sim} N(0, \sigma^2)$. What is a 100 γ % CI for σ^2 ?

We know $\frac{X_i}{\sigma} \stackrel{iid}{\sim} N(0, 1)$. Then let $R(X, \sigma^2) = \frac{1}{\sigma^2} \sum_{i=1}^n X_i^2 \sim \chi_n^2$

Pick c_1, c_2 such that $F_{\chi_n^2}(c_2) - F_{\chi_n^2}(c_1) = \gamma$.

Then $IP\left(c_1 \leq \frac{1}{\sigma^2} \sum_{i=1}^n X_i^2 \leq c_2\right) = \gamma$

$$\therefore IP\left(\frac{1}{c_2} \sum_{i=1}^n X_i^2 \leq \sigma^2 \leq \frac{1}{c_1} \sum_{i=1}^n X_i^2\right) = \gamma$$

so $\left(\frac{1}{c_2} \sum_{i=1}^n X_i^2, \frac{1}{c_1} \sum_{i=1}^n X_i^2\right)$ is a 100 γ % CI for σ^2 .

L5.3

Ex 4 Let $X_1, \dots, X_n \stackrel{iid}{\sim}$ Bernoulli (p). Find an approximate (for n large) 95% CI for p .

The MLE is $\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i$. By CLT, \hat{p} is approximately $N(p, \frac{1}{n} p(1-p))$ for n large. Thus $\frac{\sqrt{n}(\hat{p}-p)}{\sqrt{p(1-p)}}$ is approximately $N(0,1)$.

$$\text{So } \mathbb{P}\left(\hat{p} - 1.96 \sqrt{\frac{p(1-p)}{n}} \leq p \leq \hat{p} + 1.96 \sqrt{\frac{p(1-p)}{n}}\right) \approx 0.95$$

The endpoints depend on p but we can plug in \hat{p} for p so

$$\left(\hat{p} - 1.96 \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + 1.96 \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right)$$

is an approx 95% is an approx CI for p (oopsie).

NB $p(1-p) \leq \frac{1}{4}$ for $p \in [0,1]$, thus a conservative 95% CI for p is $\left(\hat{p} - \frac{1}{\sqrt{n}}, \hat{p} + \frac{1}{\sqrt{n}}\right)$.

Bayesian Inference <3

● Thus far, we've worked in a Frequentist framework. We assume the parameter $\theta \in \Theta$ is fixed, observations X have pdf/pmf $f_X(\cdot | \theta)$.

Inferences on θ through confidence intervals based on sampling X repeatedly. Statements hold no matter what θ is.

In a Bayesian framework, we must specify a prior distribution $\pi(\theta)$ on the parameter, which captures our "belief" on θ before seeing any data.

Inferences are based on the posterior distribution

$$\bullet \quad \pi(\theta | x) = \frac{f_X(x | \theta) \pi(\theta)}{f_X(x)}, \quad (\text{Bayes' rule})$$

where $f_X(x) = \int f_X(x | \phi) \pi(\phi) d\phi$. The posterior $\pi(\theta | x)$ represents our belief Θ after having seen the data x .

NB · If θ, X are joint rvs with pdf $f_X(x | \theta) \pi(\theta)$, the posterior is simply the conditional pdf of θ given x .

· As $f_X(x)$ is just a normalisation constant, we usually write

$$\pi(\theta | x) \propto f_X(x | \theta) \pi(\theta).$$

● Seeing this as a function of θ , we deduce what kind of distribution this is, and $f_X(x)$.

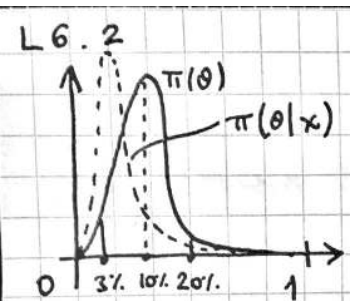
· By factorisation criterion, $\pi(\theta | x)$ only depends on x through sufficient statistics.

Ex 1 We want to estimate θ , mortality rate of new operation at Addenbrookes. In the UK, 10% of people who have the operation die, this varies from 3% to 20% in different hospitals.

Of the first 10 patients getting operation in Addenbrookes, none die.

● $X_1, X_2, \dots, X_{10} \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$.

A reasonable prior for θ is a Beta(a, b) with $a=3, b=27$.



$$\pi(3\% \leq \theta \leq 20\%) = 0.9$$

We observe $x_0 = x_1 = \dots = x_n = 0$.

$$\begin{aligned} \text{Posterior } \pi(\theta|x) &\propto \pi(\theta) f_X(x|\theta) \\ &= \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} \theta^{a-1} (1-\theta)^{b-1} \\ &= \theta^{a-1} (1-\theta)^{n+b-1}, \end{aligned}$$

which is a Beta $(a + \sum x_i, b + n - \sum x_i)$, or Beta $(3, 37)$ in our example. This has mean 0.075 (contrast this with MLE $\hat{\theta} = 0$).

Inferential statement $\pi(\theta > 0.1 | x) \approx 0.24$. It is fair to read this as "the mortality rate has probability 0.24 of exceeding 10%", but this depends on the choice of prior.

When is it best to use frequentist vs Bayesian approach? Debate!

Be a frequentist when

- There's a lot of data

- Want to make "objective" statements, which hold uniformly over the true value of θ

- Ex High-energy physics. We want an objective test of whether the Higgs Boson exists

Be a Bayesian when

- Not a lot of data, complex data

- Have useful prior information

- We wish to make decisions under uncertainty

- Ex "Nowcasting and forecasting the potential domestic & international spread of 2019 nCov..."

Bayesian estimation

Let $L(\theta, a)$ be the loss incurred when we estimate a parameter to be a when its value is θ , e.g. $L(\theta, a) = (\theta - a)^2$ or $|\theta - a|$.

The posterior expected loss is $h(a) = \int L(\theta, a) \pi(\theta|x) d\theta$.

Def A Bayes estimator under loss $L(\theta, a)$ minimises posterior expected loss, i.e. $\hat{\theta}_{\text{Bayes}} = \arg \min_{a \in \Theta} h(a)$.

Cases

① Quadratic loss, $h(a) = \int_{\Theta} (a-\theta)^2 \pi(\theta|x) d\theta$.

$$h'(a) = 0 \text{ if } a \int \pi(\theta|x) d\theta = \int \theta \pi(\theta|x) d\theta.$$

Hence Bayes estimator $\hat{\theta}$ is the posterior mean.

② Mean absolute error loss, $h(a) = \int_{\Theta} |a-\theta| \pi(\theta|x) d\theta$,

$$h(a) = \int_{-\infty}^a (a-\theta) \pi(\theta|x) d\theta + \int_a^{\infty} (\theta-a) \pi(\theta|x) d\theta,$$

$$h'(a) = 0 \text{ if } \int_{-\infty}^a \pi(\theta|x) d\theta = \int_a^{\infty} \pi(\theta|x) d\theta.$$

Hence Bayes estimator $\hat{\theta}$ is the posterior median.

Ex 2 $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, 1)$. Let $\pi(\mu)$ be a $N(0, \tau^2)$.

$$\pi(\mu|x) \propto f_x(x|\mu) \pi(\mu)$$

$$\propto \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 \right\} \cdot \exp \left\{ -\frac{\tau^2 \mu^2}{2} \right\}$$

dropping
factors
indep of μ

$$\propto \exp \left\{ -\frac{1}{2} (n+\tau^2) \left\{ \mu - \frac{\sum x_i}{n+\tau^2} \right\}^2 \right\}$$

Hence $\pi(\mu|x)$ is $N\left(\frac{\sum x_i}{n+\tau^2}, \frac{1}{n+\tau^2}\right)$.

The Bayes estimator under quadratic or absolute error loss is $\frac{\sum x_i}{n+\tau^2}$.

(contrast with MLE $\hat{\mu} = \frac{1}{n} \sum x_i$).

Bayesian Inference (cont.)

● Ex 3 $X_1, \dots, X_n \sim \text{Poi}(\lambda)$, let $\pi(\lambda)$ be an $\text{Exp}(1)$.

$$\begin{aligned} \text{Posterior } \pi(\lambda | x) &\propto f_X(x | \lambda) \cdot \pi(\lambda) \\ &\propto e^{-n\lambda} \lambda^{\sum x_i} e^{-\lambda} \\ &= \lambda^{\sum x_i} e^{-(n+1)\lambda} \end{aligned}$$

So $\pi(\lambda | x)$ is a Gamma $(\sum x_i + 1, n+1)$.

The Bayes estimator under quadratic loss is the mean $\hat{\lambda} = \frac{\sum x_i + 1}{n+1}$

Hypothesis Testing

Let $X = (X_1, \dots, X_n)$ be a vector of iid vrs taking values in \mathcal{X} ,

● with pdf / pmf f . Let H_0, H_1 be two hypotheses about f . We shall call H_0 the null hypothesis and H_1 the alternative hypothesis.

On the basis of an observation $x \in \mathcal{X}^n$ we want to choose between them.

Ex 1 $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\theta)$

$$H_0: \theta = \frac{1}{2}, \quad H_1: \theta = \frac{3}{4}$$

Ex In the same setting, we may have

$$H_0: \theta = \frac{1}{2}, \quad H_1: \theta \neq \frac{1}{2}$$

Def When a hypothesis completely determines f , we say it is simple

● (e.g. $H_0: \theta = \frac{1}{2}$), otherwise we say it is composite (e.g. $H_1: \theta \neq \frac{1}{2}$)

Ex 3 Say X_i 's take values in \mathbb{N}

H_0 : f is Poisson pmf with some unknown mean

H_1 : f is not Poisson

This is called a goodness of fit test.

Ex 4 f is in some parametric family $\{f(\cdot | \theta) : \theta \in \Theta\}$

$$H_0: \theta \in \Theta_0 \subseteq \Theta$$

$$H_1: \theta \in \Theta_1 \subseteq \Theta \quad \text{with } \Theta_0 \cap \Theta_1 = \emptyset$$

● Ex 5 $H_0: f = f_0$, $H_1: f = f_1$ where f_0, f_1 are pdfs, not necessarily in the same family

L7.2

Defⁿ We say a hypothesis test for H_0 against H_1 with critical region $C \subseteq \mathcal{X}^2$ rejects H_0 if $x \in C$, and it fails to reject H_0 when $x \notin C$.

NB You will often hear that if $x \notin C$, the test "accepts H_0 ", but this can be misleading.

A hypothesis test can make 2 types of error

Type 1 We reject H_0 when H_0 is true

Type 2 We fail to reject H_0 when H_1 is true

From now on, assume H_0, H_1 are simple hypotheses.

Let $\alpha = \text{IP}(\text{Type 1 error}) = \text{IP}(X \in C | H_0)$,

$\beta = \text{IP}(\text{Type 2 error}) = \text{IP}(X \notin C | H_1)$.

Defⁿ We say the size of a test is α , the power to detect H_1 is $1 - \beta$.

We'd like to minimise α, β , but typically they cannot be made arbitrarily small. Usually a tradeoff (decreasing α increases β).

Typical procedure for choosing test (or C) Fix Type 1 error probability α that we are happy to accept, e.g. $\alpha = 0.05$. Find test of size α which minimises β .

Given 2 simple hypotheses $H_0: f = f_0, H_1: f = f_1$, the likelihood ratio given data x is

$$\Lambda_x(H_0: H_1) = \frac{f_1(x)}{f_0(x)}$$

A likelihood ratio (LR) test is one with critical region

$$\{x \in \mathcal{X}^n : \Lambda_x(H_0, H_1) > k\} \text{ for some } k \geq 0.$$

Thm (Neyman-Pearson lemma) Suppose $H_0: f = f_0, H_1: f = f_1$, where f_0, f_1 are continuous densities which are non-zero on the same set.

Then among all tests of size $\leq \alpha$, the test with the highest power is a LR test, i.e.

$$C = \{x : \Lambda_x(H_0, H_1) > k\}$$

with k chosen such that $\alpha = \text{IP}(X \in C | H_0) = \int_C f_0(x) dx$.

L7.3

NB We require continuous f_0, f_1 such that LR of size α exists for all α with $0 \leq \alpha \leq 1$.

Proof Let $\beta = P(X \notin C | H_1) = \int_{\bar{C}} f_1(x) dx$, where \bar{C} is the complement of C in \mathcal{X}^n .

Take another test with critical region D , define $\alpha^* = \int_D f_0(x) dx$, $\beta^* = \int_{\bar{D}} f_1(x) dx$. Suppose $\alpha^* \leq \alpha$. WTS $\beta^* \geq \beta$.

$$\beta - \beta^* = \int_{\bar{C}} f_1(x) dx - \int_{\bar{D}} f_1(x) dx$$

$$= \left(\int_{\bar{C} \cap D} + \int_{\bar{C} \cap \bar{D}} - \int_{\bar{D} \cap C} - \int_{\bar{D} \cap \bar{C}} \right) f_1(x) dx$$

$$= \left(\int_{\underbrace{\bar{C} \cap D}_{\subseteq \bar{C}}} - \int_{\underbrace{\bar{D} \cap C}_{\subseteq C}} \right) f_0(x) \Delta_x(H_0, H_1) dx$$

$\therefore \lambda \leq k$ $\therefore \lambda > k$

$$\leq k \int_{\bar{C} \cap D} f_0(x) dx - k \int_{\bar{D} \cap C} f_0(x) dx$$

$$= k \left(\int_{\bar{C} \cap D} + \int_{C \cap D} - \int_{C \cap D} - \int_{\bar{D} \cap C} \right) f_0(x) dx$$

$$= k \left(\int_D - \int_C \right) f_0(x) dx$$

$$= k (\alpha^* - \alpha)$$

$$\leq 0.$$

Ex Let $X_0, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, σ^2 known. Want to find most powerful test of size α for $H_0: \mu = \mu_0$, against $H_1: \mu = \mu_1$. □

$$\begin{aligned} \Delta_x(H_0, H_1) &= \exp\left(-\frac{1}{2\sigma^2} \sum (x_i - \mu_1)^2 + \frac{1}{2\sigma^2} \sum (x_i - \mu_0)^2\right) \\ &= \exp\left(\frac{(\mu_1 - \mu_0)}{\sigma^2} n \bar{x} + \frac{n(\mu_0^2 - \mu_1^2)}{2\sigma^2}\right) \end{aligned}$$

For $\mu_0 < \mu_1$, Δ_x is increasing in \bar{x} , hence

$$\Delta_x > k \iff \bar{x} > c \text{ for some } c.$$

The critical region of LR test is $\{x: \bar{x} > c\}$ for some c .

Let $Z = \frac{(\bar{X} - \mu_0)\sqrt{n}}{\sigma}$. Then $Z \sim N(0, 1)$ under H_0 .

L6.4

And $\bar{x} > c \Leftrightarrow z > c^*$.

Hence the LR test has optimal region $C = \{x : z > c^*\}$.

Lastly, we find c^* s.t.

$$P(z > c^* | H_0) = \alpha \quad \text{i.e.} \quad c^* = z_\alpha = \Phi^{-1}(1 - \alpha).$$

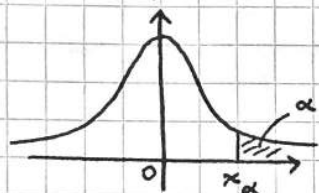
Example $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, σ^2 known

● Test $H_0: \mu = \mu_0$, $H_1: \mu = \mu_1$, $\mu_0 < \mu_1$

$\Delta_x(H_0, H_1)$ increasing function of $Z = \sqrt{n}(\bar{X} - \mu_0)/\sigma$.

Hence $\Delta_x(H_0, H_1) > k \Leftrightarrow Z > c$ for some c depending on k .

Under H_0 , $Z \sim N(0, 1)$. So a LR test of size α has critical region $\{x \in \mathcal{X}^n : z > z_\alpha\}$ where z_α is upper α point of $N(0, 1)$.



We say Z is a test statistic, the integral

$$\int_z^\infty \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = p$$

is called the p-value, i.e. the probability that the test statistic is

● more extreme than the observed value. (under H_0 ?)

Let $\sigma^2 = 1$, $n = 10$, $\mu_0 = 0$, $\mu_1 = 1$. Then $z = 6.32$. ($\bar{X} = 2$)

Say we wish to make a test of size 0.05. Then $z_\alpha = 1.69$.

Hence the test rejects H_0 , since $6.32 > 1.69$. In this case, the p-value is $p = 1.3 \times 10^{-10}$.

This tells us that we could have rejected a test of much smaller size, so this is strong evidence against H_0 .

Prop The p-value has a $\text{Unif}(0, 1)$ distribution under H_0 .

● Proof Note $p = \int_z^\infty f_Z(u) du$.

$$\begin{aligned} \mathbb{P}(P \leq t) &= \mathbb{E}[\mathbb{1}_{\{P \leq t\}}] = \mathbb{E}[\mathbb{1}_{\{Z \geq F_Z^{-1}(1-t)\}}] \\ &= 1 - \left(F_Z^{-1}(1-t) \right) = t. \end{aligned}$$

□

Composite hypotheses

Suppose X_1, \dots, X_n iid from parametric $f(x|\theta)$, $\theta \in \Theta$.

For a composite hypothesis such as $H: \theta > 0$, the probabilities of Type I and Type II error do not have a single value.

● Def The power function of a test with critical region C is

$$W(\theta) = \mathbb{P}(X \in C | \theta) = \mathbb{P}(\text{reject } H_0 | \theta).$$

L8.2

Consider $H_0: \theta \in \Theta_0$, $H_1: \theta \in \Theta_1$, $\Theta_0, \Theta_1 \in \mathcal{P}(\Theta)$.

The size of test is $\alpha = \sup_{\theta \in \Theta_0} W(\theta)$.

We say a test of H_0 against H_1 is uniformly most powerful (UMP)

of size α if (1) $\sup_{\theta \in \Theta_0} W(\theta) = \alpha$
(2) For any other critical region C^* , with power function W^* , for all $\theta \in \Theta_1$, $W(\theta) \geq W^*(\theta)$.
of size $\leq \alpha$?

NB A UMP test need not exist, but in many cases LR test is uniformly most powerful.

Ex $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, σ^2 known

Test $H_0: \mu \leq \mu_0$ vs $H_1: \mu \geq \mu_0$

Let C be the critical region $\{x \in \mathcal{X}^n: \sqrt{n}(\bar{X} - \mu_0)/\sigma > z_\alpha\}$ (i.e. the LR test for $H_0: \mu = \mu_0$, $H_1: \mu = \mu_1$, $\mu_0 < \mu_1$).

The power function

$$\begin{aligned} W(\mu) &= \mathbb{P}_\mu(\text{reject } H_0) = \mathbb{P}_\mu(\sqrt{n}(\bar{X} - \mu_0)/\sigma > z_\alpha) \\ &= \mathbb{P}_\mu\left(\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} > z_\alpha + \frac{\sqrt{n}(\mu_0 - \mu)}{\sigma}\right) \\ &= 1 - \Phi\left(z_\alpha + \frac{\sqrt{n}(\mu_0 - \mu)}{\sigma}\right). \end{aligned}$$

So $W(\mu)$ is increasing in μ . So

$$\sup_{\mu \leq \mu_0} W(\mu) = W(\mu_0) = \alpha \text{ by previous example.}$$

So (1) is satisfied. Now, take C^* to be any other critical with size $\leq \alpha$. Then by Neyman-Pearson lemma, $W(\mu_1) \geq W^*(\mu_1)$ for any $\mu_1 > \mu_0$, i.e. any μ_1 in H_1 . Hence (2) is also satisfied and the LR test is UMP.

NB In many examples of the form $H_0: \theta \leq \theta_0$, $H_1: \theta > \theta_0$, we can apply the same reasoning.

(1) Determine LR test for $H_0: \theta = \theta_0$, $H_1: \theta = \theta_1$, $\theta_0 < \theta_1$

(2) Observe that critical region C does not depend on θ_1

(3) Hence by NP lemma, the test is UMP

L8.3

Suppose $\Theta_0 \subseteq \Theta_1$. (nested hypotheses). The generalised likelihood

● ratio is
$$\Lambda_x(H_0, H_1) = \frac{\sup_{\theta \in \Theta_1} f(x|\theta)}{\sup_{\theta \in \Theta_0} f(x|\theta)} \geq 1.$$

Ex $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$, σ^2 known

Test $H_0: \theta = \theta_0$, $H_1: \theta \in \mathbb{R}$ ← meant μ 's

The generalised likelihood ratio is

$$\Lambda_x(H_0, H_1) = \frac{(2\pi\sigma)^{-n/2} \exp(-\frac{1}{2\sigma^2} \sum (x_i - \bar{x})^2)}{(2\pi\sigma)^{-n/2} \exp(-\frac{1}{2\sigma^2} \sum (x_i - \mu_0)^2)}$$

Note $2 \log \Lambda_x(H_0, H_1) = \frac{n}{\sigma^2} (\bar{X} - \mu_0)^2$.

● Under H_0 , this is a χ_1^2 v.v. The GLR test rejects when $\Lambda_x(H_0, H_1)$ is large.

Hence for a size α test, we reject H_0 when $\frac{n}{\sigma^2} (\bar{X} - \mu_0)^2 > \chi_1^2(\alpha)$, the upper α point of a χ_1^2 dist.

Note $Z^2 = \frac{(\bar{X} - \mu_0)^2 n}{\sigma^2}$. Thus this test rejects when Z is very large or very small (called a "2-tailed test").

L9.1 Chi-squared tests

Recall example from last time: $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, σ^2 known

● $H_0: \mu = \mu_0$, $H_1: \mu \neq \mu_0$

The generalised likelihood ratio Λ had

$$2 \log \Lambda = \frac{(\bar{X} - \mu_0)^2 n}{\sigma^2} \sim \chi_1^2$$

~~Thm~~ Suppose X_1, \dots, X_n are iid from a parametric model $f(\cdot | \theta)$.

Let $H_0: \theta \in \mathcal{H}_0$, $H_1: \theta \in \mathcal{H}_1$, with $\mathcal{H}_0 \subseteq \mathcal{H}_1$.

Let $\dim \mathcal{H}_0$ denote the dimension of \mathcal{H}_0 . Usually $\mathcal{H} \subseteq \mathbb{R}^k$ and \mathcal{H}_0 is a linear subspace of \mathbb{R}^k , $\{\theta \in \mathbb{R}^k: A\theta = b\}$ for some matrix $A \in \mathbb{R}^{p \times k}$, $b \in \mathbb{R}^p$; then $\dim \mathcal{H}_0 = k - p$.

● Similarly define $\dim \mathcal{H}_1$.

Thm (Wilks') Suppose $\mathcal{H}_0 \subseteq \mathcal{H}_1$, $\dim \mathcal{H}_1 - \dim \mathcal{H}_0 = p$. Then under regularity conditions, if H_0 is true, i.e. X_1, \dots, X_n are iid with pdf $f(\cdot | \theta)$ for $\theta \in \mathcal{H}_0$, then as $n \rightarrow \infty$, the distribution of $2 \log \Lambda$ converges to a χ_p^2 distribution.

[More precisely, if F_n is the distribution of $2 \log \Lambda$, and F is that of a χ_p^2 distribution, then $F_n \rightarrow F$ pointwise]

If H_1 is true but H_0 is not, then $2 \log \Lambda$ will tend to be

● bigger. Hence a test which rejects H_0 when $2 \log \Lambda > \chi_p^{-1}(\alpha)$ has size approximately α , by the theorem, and some power to detect H_1 .

Goodness-of-fit test

Birth month of everyone admitted to Cambridge in 2012

Month	Sep	Oct	Nov	Dec	...	Aug
n_i	470	515	470	457	...	399

Is this compatible with the proportions $\tilde{p}_1, \dots, \tilde{p}_{12}$ of UK births each month?

● NB \tilde{p} can be determined from Census data, is not uniform
"Christmas effect" more births around September

Model Let N_i be the # births in month i

$k=12$

$(N_1, \dots, N_{12}) \sim \text{Multinomial}(n, (p_1, \dots, p_{12}))$, n fixed

$H_0: p_i = \tilde{p}_i$, $H_1: p_i$'s unrestricted, $p \in [0, 1]^{12}$, $\sum_{i=1}^{12} p_i = 1$

Here $\dim(\Theta_0) = 0$, $\dim(\Theta_1) = 11$

$$\Lambda = \frac{\sup_{p \in \Theta_1} \binom{n}{n_1 \dots n_k} p_1^{n_1} \dots p_k^{n_k}}{\sup_{p \in \Theta_0} \binom{n}{n_1 \dots n_k} p_1^{n_1} \dots p_k^{n_k}} = \frac{\hat{p}_1^{n_1} \dots \hat{p}_k^{n_k}}{\tilde{p}_1^{n_1} \dots \tilde{p}_k^{n_k}} \leftarrow \hat{p} \text{ is MLE}$$

The log likelihood is

$$l(p) = \text{const.} + \sum_{i=1}^k n_i \log p_i$$

Using the Lagrangian $L(p, \lambda) = \sum_{i=1}^k n_i \log p_i - \lambda (\sum_{i=1}^k p_i - 1)$ we can minimise $l(p)$ subject to $\sum p_i = 1$.

We obtain $\hat{p}_i = \frac{n_i}{n}$ for $i=1, \dots, k$.

Hence

$$2 \log \Lambda = 2 \log \left[\left(\frac{n_i}{n \tilde{p}_i} \right)^{n_i} \dots \left(\frac{n_k}{n \tilde{p}_k} \right)^{n_k} \right]$$

$$= 2 \sum_{i=1}^k n_i \log \left(\frac{n_i}{n \tilde{p}_i} \right).$$

A chi-squared test of size approx α rejects the null H_0 when

$$2 \log \Lambda > \chi_{k-1}^2{}^{-1}(\alpha).$$

In this example we have $2 \log \Lambda = 44.9$. Say $\alpha = 0.01$, we have p-value $IP(\chi_{k-1}^2 > 44.9) \approx 3 \times 10^{-9}$. As this is < 0.01 we can reject H_0 at level 1%.

Modify H_0 to be $H_0: p_i = p_i(\theta)$ for some scalar parameter θ ,

e.g. $p_1(\theta) = \frac{1}{12} + \theta$, $p_2(\theta) = \dots = p_{12}(\theta) = \frac{1}{12} - \frac{\theta}{11}$, $\theta \in \mathbb{R}$ represents the Christmas effect.

We have still $H_1: p_i$'s unrestricted.

In this case $2 \log \Lambda = \sum n_i \log \left(\frac{n_i}{n p_i(\hat{\theta})} \right)$ where $\hat{\theta}$ is the MLE of θ under H_0 .

L 9.3

Pearson's χ^2 statistic Let $o_i = n_i$, the "observed # samples of type i ".

Let e_i be the "expected # samples of type i " under H_0 .

In first example, $e_i = n \times \tilde{p}_i$.

In our second example, $e_i = n \times p_i(\hat{\theta})$.

We can write

$$\begin{aligned} 2 \log \bar{\Lambda} &= 2 \sum o_i \log \left(\frac{o_i}{e_i} \right), \quad \text{set } \delta_i = o_i - e_i \\ &= 2 \sum (e_i + \delta_i) \log \left(1 + \frac{\delta_i}{e_i} \right) \\ &\approx 2 \sum \left(\delta_i + \frac{\delta_i^2}{e_i} - \frac{\delta_i^3}{2e_i} \right) \quad [\text{Taylor}] \\ &= \sum \frac{\delta_i^2}{e_i} \quad \left[\sum \delta_i = \sum o_i - \sum e_i = n - n \right] \\ &= \sum \frac{(o_i - e_i)^2}{e_i} \end{aligned}$$

We call this the Pearson χ^2 statistic. This is referred to a χ_p^2 distribution, where $p = \dim(\mathbb{N}_1) - \dim(\mathbb{N}_0) = k - 1 - \dim(\mathbb{N}_0)$

Ex Mendel crossed 556 SY peas with WG peas.

From progeny, N_1 were SY

N_2 were SG

N_3 were WY

N_4 were WG

He observed $(n_1, n_2, n_3, n_4) = (315, 108, 103, 31)$

Mendel's hypothesis $(p_1, \dots, p_4) = \left(\frac{9}{16}, \frac{3}{16}, \frac{3}{16}, \frac{1}{16} \right)$

Here, $2 \log \bar{\Lambda} = 0.618$, $\sum \frac{(o_i - e_i)^2}{e_i} = 0.604$.

The p-value is $P(\chi_3^2 > 0.604) \approx 0.96$.

So no evidence against Mendel's hypothesis.

It wouldn't be rejected by a test of size 10%.

Actually, the high p-value (small test statistic) would suggest that the $\hat{\Lambda}$ numbers are too close to the theoretical proportions.

observed

Chi-squared tests (cont.)Contingency tables

Ex Survey 500 people

	Prefer Kate	Prefer Meghan	No preference
Leave	56	52	42
Remain	50	83	67
No vote	18	51	81

 N_{ij} : # of samples in row i column j Observations are $(N_{ij}, i=1, \dots, r, j=1, \dots, c)$ Denote $N_{i+} = \sum_{j=1}^c N_{ij}$, $N_{+j} = \sum_{i=1}^r N_{ij}$ Model 1 $(N_{ij}, i=1, \dots, r, j=1, \dots, c) \sim \text{Multinomial}(n, (p_{ij}, \substack{i=1, \dots, r \\ j=1, \dots, c}))$ n fixed, in this example $n=500$ p_{ij} probability that a sample falls in cell (i, j) Independence test H_0 : Preferred royal indep of Brexit votei.e. $p_{ij} = p_{i+} p_{+j}$ where $p_{i+} = \sum_{j=1}^c p_{ij}$, $p_{+j} = \sum_{i=1}^r p_{ij}$ $(p_{1+}, \dots, p_{r+}), (p_{+1}, \dots, p_{+c})$ are unconstrained probability vectors H_1 : $(p_{ij}: i=1, \dots, r, j=1, \dots, c)$ is unconstrained probability vectorMLE of p under H_1 $\hat{p}_{ij} = \frac{N_{ij}}{n}$ (same as multinomial model!
in previous lecture)MLE of p under H_0

$$\mathcal{L}(p) = \frac{n!}{\prod_{i,j} N_{ij}!} \prod_{i,j} (p_{i+} p_{+j})^{N_{ij}}$$

Take logarithm, maximise subject to $\sum_{i=1}^r p_{i+} = \sum_{j=1}^c p_{+j} = 1$ Use Lagrangian method to obtain $\tilde{p}_{ij} = \tilde{p}_{i+} \tilde{p}_{+j} = \frac{N_{i+} N_{+j}}{n^2}$ Recall $2 \log \Lambda = 2 \sum_{i,j} o_{ij} \log \left(\frac{o_{ij}}{e_{ij}} \right) \approx \sum_{i,j} \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$ In this case, $o_{ij} = N_{ij}$, "observed number" $e_{ij} = \tilde{p}_{ij} n$ "expected number under H_0 "

$$= N_{i+} N_{+j} / n$$

With these numbers we can compute $2 \log \Lambda$ or the Pearson χ^2 stat.By Wilkes' theorem, this has distribution approx $\chi^2_{\dim \Theta_1 - \dim \Theta_0}$

Here $\dim \mathbb{H}_1 = rc - 1$ [1 equality condition $\sum_{ij} p_{ij} = 1$]

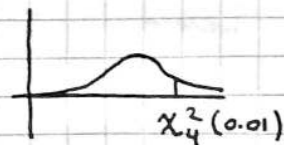
$\dim \mathbb{H}_0 = (r-1)c + (c-1)$ [(p_{1+}, \dots, p_{r+}) has 1 equality constraint
 (p_{+1}, \dots, p_{+c}) "]

Thus, the degrees of freedom are $rc - 1 - (r-1)c - (c-1) = (r-1)(c-1)$

In this example $\sum_{ij} \frac{(o_{ij} - e_{ij})^2}{e_{ij}} = 36.20$

This is referred to a $\chi^2_{(3-1)(3-1)} = \chi^2_4$.

Observe that upper 1% point $\chi^2_4(0.01) = 13.28$



As $36.20 > 13.28$, a test of size 1% rejects H_0 .

Strong evidence against independence.

Ex 50 patients are assigned to each of 3 treatments. Outcomes are recorded in a contingency table.

	Improvement	No Improvement	Worsening	
Placebo	18	17	15	50
1/2 dose	20	20	20	50
Full dose	25	13	12	50

Model 2 n_{i+} fixed for $i=1, \dots, r$

$(N_{i1}, \dots, N_{ic}) \sim \text{Multinomial}(n_{i+}, (p_{i1}, \dots, p_{ic}))$

independently for each row $i=1, \dots, r$

In the example, $n_{1+} = n_{2+} = n_{3+} = 50$

Homogeneity test H_0 : Treatment has no effect on the outcome

i.e. the probabilities of each outcome are the same for different treatments

$p_{1j} = p_{2j} = \dots = p_{rj} =: q_j$ for $j=1, \dots, c$

for some vector of probabilities (q_1, \dots, q_c)

H_1 : (p_{i1}, \dots, p_{ic}) is unconstrained probability vector for each row

MLE under H_1 $\mathcal{L}(p) = \prod_{i=1}^r \binom{n_{i+}}{N_{i1} \dots N_{ic}} p_{i1}^{N_{i1}} \dots p_{ic}^{N_{ic}}$

Take a log, maximise subject to $\sum_j p_{ij} = 1$ for each i

$\Rightarrow \hat{p}_{ij} = \frac{N_{ij}}{n_{i+}}$ [NB This amounts to maximising each factor separately]

MLE of p under H_0 $\mathcal{L}(q) = \prod_{i=1}^r \binom{n_{i+}}{N_{i1} \dots N_{ic}} q_1^{N_{i1}} \dots q_c^{N_{ic}}$

Maximise subject to $\sum_j q_j = 1$

$\Rightarrow \hat{q}_j = \frac{N_{+j}}{n}$ where $n = \sum_i n_{i+}$

people getting each treatment is fixed

L10.3

$$o_{ij} = N_{ij}, \quad e_{ij} = n_{i+} \tilde{q}_j = n_{i+} N_{+j} / n$$

● NB e_{ij}, o_{ij} are identical to in the independence test

The only difference is that previously n_{i+} was random, whereas now it is fixed.

$$\text{Therefore } 2 \log \Lambda = 2 \sum_{i,j} o_{ij} \log \left(\frac{o_{ij}}{e_{ij}} \right)$$

$$\text{Pearson's } \chi^2 = \sum_{i,j} \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$$

are identical in both tests.

The reference distribution is $\chi^2_{\dim \mathbb{H}_1 - \dim \mathbb{H}_0}$.

Here, $\dim \mathbb{H}_1 = r(c-1)$ [r probability vectors of length c]

$$\dim \mathbb{H}_0 = c-1$$

So degrees of freedom are $(r-1)(c-1)$.

Conclude Homogeneity test has the same conclusions as the independence test despite the differences in their models.

L11.1 Relationship between confidence sets and hypothesis tests

Suppose $(X_1, \dots, X_n) = X$ have joint pdf/pmf $f_X(x|\theta)$ for some $\theta \in \Theta$.

● Thm (i) Suppose for any $\theta_0 \in \Theta$ there is a size α test for $H_0: \theta = \theta_0$ with critical region $C(\theta_0)$. Then the set $I(X) = \{\theta: X \notin C(\theta)\}$ is a $100(1-\alpha)\%$ confidence set for θ ,

i.e. we build a confidence set from all θ_0 such that $H_0: \theta = \theta_0$ is not rejected.

(ii) Suppose that $I(X)$ is a $100(1-\alpha)\%$ confidence set for θ .

Then $C(\theta_0) = \{X: \theta_0 \notin I(X)\}$ is the critical region of a size α test for $H_0: \theta = \theta_0$,

● i.e. the test rejects H_0 when θ_0 is not in $I(X)$

Proof Note that in each case $\theta_0 \in I(X) \Leftrightarrow X \notin C(\theta_0)$.

Observe

$$\begin{aligned} \text{IP}(\text{reject } H_0 \mid H_0 \text{ is true}) &= \text{IP}(X \in C(\theta_0) \mid H_0 \text{ is true}) \\ &= 1 - \text{IP}(\theta_0 \in I(X) \mid \theta = \theta_0). \end{aligned}$$

For part (i), we assume LHS = α , thus

$$\text{IP}(\theta_0 \in I(X) \mid \theta = \theta_0) = 1 - \alpha$$

which implies $I(X)$ is a $100(1-\alpha)\%$ confidence set.

● For part (ii), we assume $\text{IP}(\theta_0 \in I(X) \mid \theta = \theta_0) = 1 - \alpha$ which implies $\text{IP}(\text{reject } H_0 \mid H_0 \text{ is true}) = 1 - (1 - \alpha) = \alpha$ and so $C(\theta_0)$ is the critical region of a test of size α for $H_0: \theta = \theta_0$. \square

Ex $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, 1)$. Suppose we want a 95% confidence set for μ . Recall that a size 0.05 test for $H_0: \mu = \mu_0$ rejects when the test statistic $Z = \sqrt{n}(\bar{X} - \mu_0)$ has $|Z| > 1.96$ [1.96 is the upper 2.5% point of a $N(0,1)$ distribution]

By part (i) of the theorem, a 95% conf. set for μ is given by

$$\begin{aligned} I(X) &= \{\mu_0: X \notin C(\mu_0)\} = \{\mu_0: |Z| \leq 1.96\} \\ &= \{\mu_0: |(\bar{X} - \mu_0)\sqrt{n}| \leq 1.96\} = \left(\bar{X} - \frac{1.96}{\sqrt{n}}, \bar{X} + \frac{1.96}{\sqrt{n}}\right). \end{aligned}$$

L11.2

This is the CI we found previously.

Orthogonal projections

Let V be a subspace of \mathbb{R}^n . Let $V^\perp = \{w \in \mathbb{R}^n : w^T v = 0 \forall v \in V\}$, be its orthogonal complement. Recall any vector $x \in \mathbb{R}^n$ can be written uniquely as $x = v + w$ with $v \in V, w \in V^\perp$.

Def An orthogonal projection matrix $\Pi \in \mathbb{R}^{n \times n}$ onto V acts on any $x = v + w, v \in V, w \in V^\perp$ by $\Pi(x) = v$.

Properties (1) The column space of Π is V .

(2) $I - \Pi$ is an orthogonal projection onto V^\perp . Indeed, if $x = v + w, v \in V, w \in V^\perp$, then

$$(I - \Pi)x = v + w - v = w.$$

(3) $\Pi^2 = \Pi = \Pi^T$
 ↑ ↑
 idempotent self-adjoint

$\Pi^2 x = \Pi v = v$, so idempotency is clear from defⁿ

For self-adjoint-ness, let $x_1, x_2 \in \mathbb{R}^n$. Then

$$(\Pi x_1)^T ((I - \Pi)x_2) = 0$$

$$\therefore x_1^T (\Pi^T - \Pi^T \Pi) x_2 = 0$$

This holds for all $x_1, x_2 \in \mathbb{R}^n$, so $\Pi^T = \Pi^T \Pi$.

Since $\Pi^T \Pi$ is symmetric, so are Π^T and Π .

Conversely, any Π which satisfies $\Pi^2 = \Pi = \Pi^T$ is an orthogonal projection onto $\text{col}(\Pi)$. Proof left as an exercise.

(4) An orthonormal basis for V , and V^\perp is made up of eigenvectors of Π with eigenvalues 1 and 0 respectively. So we can write $\Pi = UDU^T$ for U orthogonal, D diagonal.

Moreover, D has p 1s and $n-p$ 0s on its diagonal, where

$$p = \dim V.$$

(5) $\text{Tr}(\Pi) = \text{Tr}(UDU^T) = \text{Tr}(D) = p = \dim V \leftarrow = \text{Rank}(\Pi)$

Multivariate Normal Distributions

● If $X = (X_1, \dots, X_n)$ is a random vector, we define

$$\mathbb{E}X = (\mathbb{E}X_1, \dots, \mathbb{E}X_n).$$

Similarly for a random matrix X , $(\mathbb{E}X)_{ij} = \mathbb{E}(X_{ij})$.

For 2 ^{col.} vectors V, W define $\text{cov}(V, W) = \mathbb{E}[(V - \mathbb{E}V)(W - \mathbb{E}W)^T]$

We'll also let $\text{cov}(X) = \text{cov}(X, X)$.

Linearity of expectation implies

① For any fixed $A \in \mathbb{R}^{m \times n}$, $\mathbb{E}(AX) = A\mathbb{E}(X)$.

② For any matrices A, B ,

$$\begin{aligned} \bullet \text{cov}(AX, BX) &= \mathbb{E}[(AX - \mathbb{E}(AX))(BX - \mathbb{E}(BX))^T] \\ &= \mathbb{E}[A(X - \mathbb{E}X)(X - \mathbb{E}X)^T B^T] \\ &= A \text{cov}(X) B^T. \end{aligned}$$

Def We say a random vector X is multivariate normal if $\forall t \in \mathbb{R}^n$, $t^T X$ has a normal distribution.

Lemma The MVN distribution is completely determined by

$$\mu = \mathbb{E}X \text{ and } \Sigma = \text{cov}(X).$$

Hence we write $X \sim N(\mu, \Sigma)$.

● Proof The multivariate moment generating function is

$$M_X(t) = \mathbb{E}e^{t^T X} \text{ for } t \in \mathbb{R}^n.$$

When this exists for all $t \in \mathbb{R}^n$, it fully determines the distribution of X . Now if $X \sim N(\mu, \Sigma)$, then $t^T X \sim N(t^T \mu, t^T \Sigma t)$, by the properties above.

Hence MGF is $M_X(t) = \mathbb{E}e^{t^T X} = M_{t^T X}(1) = \exp(t^T \mu + \frac{1}{2} t^T \Sigma t)$, which only depends on μ and Σ . □

Lemma Let $\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{pmatrix}\right)$ where $X_i, \mu_i \in \mathbb{R}^{n_i}$, and

• $\Sigma_{ij} \in \mathbb{R}^{n_i \times n_j}$. Then

(i) $X_i \sim N(\mu_i, \Sigma_{ii})$,

(ii) $X_1 \perp X_2 \iff \Sigma_{12} = 0$.

Proof (i) WLOG $i=1$. Note X_1 is multivariate normal, because $\forall t \in \mathbb{R}^{n_1}$, $t^T X_1 = (t \ 0)^T \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ is normal by assumption.

And $\mathbb{E} X_1 = \mu_1$, $\text{Cov} X_1 = \Sigma_{11}$.

(ii) (\implies) If X_1, X_2 are indep, then

$$\Sigma_{12} = \text{cov}(X_1, X_2) = \mathbb{E}((X_1 - \mu_1)(X_2 - \mu_2)^T) = 0.$$

• (\impliedby) As the MVN distⁿ is characterised by its mean and covariance, if $\Sigma_{12} = 0$, then $X_1 \perp X_2$. □

NB (ii) does not necessarily hold if X_1, X_2 are not jointly normal.

Lemma Let π be an orthogonal projection of rank p . Let $X \sim N(0, I_n)$. Then $\|\pi X\|_2^2 \sim \chi_p^2$.

Proof Let $\pi = UDU^T$ be an eigendecomposition of π .

$$\|\pi X\|_2^2 = \|UDU^T X\|_2^2 \stackrel{U \text{ orthog}}{=} \|DU^T X\|_2^2 \stackrel{(d)}{=} \|DX\|_2^2$$

The last identity follows from the fact that $U^T X \sim N(0, I_n)$.

• Indeed, $U^T X$ is MVN (check $t^T U^T X = (Ut)^T X$ is normal), and $\mathbb{E} U^T X = U^T \mathbb{E} X = 0$, $\text{Cov}(U^T X) = U^T \text{Cov} X U = U^T U = I_n$.

We note $\|DX\|_2^2 = \sum_{i=1}^p X_i^2 \sim \chi_p^2$.

Indeed, the X_i are iid $N(0,1)$ variables. □

Normal random sample

Suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$. We don't assume σ^2 is known.

We want to estimate μ .

Let $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$, $S_{XX} = \sum_{i=1}^n (X_i - \bar{X})^2$

• Claim (i) $\bar{X} \sim N(\mu, \sigma^2/n)$ [this was shown previously]

(ii) $S_{XX}/\sigma^2 \sim \chi_{n-1}^2$ (iii) $\bar{X} \perp S_{XX}$

L12.3

Hence $\bar{X} \pm t_{n-1}(\frac{\alpha}{2}) \frac{\tilde{\sigma}}{\sqrt{n}}$ is a $100(1-\alpha)\%$ CI for μ .

● NB When σ^2 is known, we found a CI $\bar{X} \pm z(\frac{\alpha}{2}) \frac{\sigma}{\sqrt{n}}$. So in the new procedure, we replace σ by an estimator $\tilde{\sigma}$ and the normal quantile by a t_{n-1} quantile. This makes a big difference when n is small.

Linear models

Observations or responses Y_1, \dots, Y_n modelled as

$$Y_i = x_{i1}\beta_1 + \dots + x_{ip}\beta_p + \varepsilon_i \quad \text{for } i=1, \dots, n$$

- β_1, \dots, β_p are coefficients, the parameters of the model
- x_{i1}, \dots, x_{ip} are fixed covariates or predictors or inputs
- $\varepsilon_1, \dots, \varepsilon_n$ are noise rvs

Ex Y_i # covid-19 cases in region i detected Jan 23-29

x_{i1} # daily flights to i from Wuhan

x_{i2} # tests conducted in i , Jan 23-29

Matrix formulation Let $Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$, $\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$, $\varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$

and $X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & \dots & \dots & x_{np} \end{pmatrix} \in \mathbb{R}^{n \times p}$, the design matrix. We will assume

that X has full column rank, i.e. no covariate is a linear combination of the others. The model is $Y = X\beta + \varepsilon$.

$\underbrace{\quad}_{\text{random}} \quad \underbrace{\quad}_{\text{fixed}} \quad \underbrace{\quad}_{\text{random}}$

Moment assumptions for ε

① $\mathbb{E} \varepsilon = 0$

② $\text{Cov}(\varepsilon) = \sigma^2 I_n$ for some $\sigma^2 \in (0, \infty)$ *homoskedasticity*

i.e. the noise variables are uncorrelated with the same variance

Ordinary least-squares (OLS) estimator

Note $\mathbb{E} Y = \mathbb{E}(X\beta + \varepsilon) = X\beta$.

A reasonable way to estimate β is to minimize the residual sum of squares $S(\beta) = \sum_i (Y_i - (X\beta)_i)^2 = \|Y - X\beta\|_2^2$

The OLS estimator is the minimiser of $S(\beta)$ over \mathbb{R}^p , written $\hat{\beta}$.

$\frac{\partial S(\beta)}{\partial \beta} = 2X^T(Y - X\beta)$ ^{sign?} so $2X^T(Y - X\hat{\beta}) = 0$

$\therefore \hat{\beta} = (X^T X)^{-1} X^T Y$ NB $X^T X$ is invertible because X has rank p

117
L13.2

Lemma $E\hat{\beta} = \beta$ [$\hat{\beta}$ unbiased for β], $\text{cov}(\hat{\beta}) = \sigma^2(X^T X)^{-1}$

Proof $E\hat{\beta} = E((X^T X)^{-1} X^T Y) = E((X^T X)^{-1} X^T (X\beta + \varepsilon))$
 $= \beta + \underbrace{(X^T X)^{-1} X^T E(\varepsilon)}_{\text{zero}}$

$$\text{cov}\hat{\beta} = \text{cov}\left(\underbrace{\beta}_{\text{const.}} + (X^T X)^{-1} X^T \varepsilon\right) = \text{cov}\left((X^T X)^{-1} X^T \varepsilon\right)$$

$$= (X^T X)^{-1} X^T \text{cov}(\varepsilon) X (X^T X)^{-1}$$

$$= \sigma^2 (X^T X)^{-1} \quad \square$$

Theorem (Gauss-Markov) Let $\hat{\beta}$ be the OLS estimator. Let $\beta^* = AY$ be a linear unbiased estimator of β . Then for all $t \in \mathbb{R}^p$,

$$\text{Var}(t^T \hat{\beta}) \leq \text{var}(t^T \beta^*)$$

The prediction of the linear model at any input t has the smallest variance with the OLS estimator. We say $\hat{\beta}$ is the Best Linear Unbiased Estimator (BLUE) of β .

Proof $E(AY) = E(A(X\beta + \varepsilon)) = AX\beta = \beta$ as β^* unbiased

Define $B = A - (X^T X)^{-1} X^T$, observe $BX = 0$.

$$\text{Then } \text{cov}(\beta^*) = E[(\beta^* - \beta)(\beta^* - \beta)^T]$$

$$= E[(\beta^* - \hat{\beta} + \hat{\beta} - \beta)(\beta^* - \hat{\beta} + \hat{\beta} - \beta)^T]$$

$$= \text{cov}(\hat{\beta}) + \underbrace{\text{cov}(BY)}_{BY} + \underbrace{E[BY(\hat{\beta} - \beta)^T]}_M + \underbrace{E[(\hat{\beta} - \beta)(BY)^T]}_{M^T}$$

$$\text{cov}(BY) = \text{cov}(BX\beta + B\varepsilon) = \text{cov}(B\varepsilon) = \sigma^2 BB^T$$

$$M = E[B\varepsilon \underbrace{((X^T X)^{-1} X^T (X\beta + \varepsilon) - \beta)^T}_{\text{zero}}]$$

$$= E[B\varepsilon \underbrace{((X^T X)^{-1} X^T \varepsilon)^T}_{\text{zero}}] = B \text{cov}(\varepsilon) X (X^T X)^{-1}$$

$$= \sigma^2 \underbrace{BX}_{\text{zero}} (X^T X)^{-1} = 0.$$

Hence $\text{cov}(\beta^*) \neq \text{cov}(\hat{\beta}) + \sigma^2 BB^T$.

Now for any $t \in \mathbb{R}^p$,

$$\text{var}(t^T \beta^*) = t^T \text{cov}(\beta^*) t = \text{var}(t^T \hat{\beta}) + \underbrace{\sigma^2 t^T BB^T t}_{\geq 0} \quad \square$$

L13.3

Def The fitted values in OLS estimator are the predictions of the response, $\hat{Y} = X\hat{\beta} = \underbrace{X(X^T X)^{-1} X^T}_P Y$.

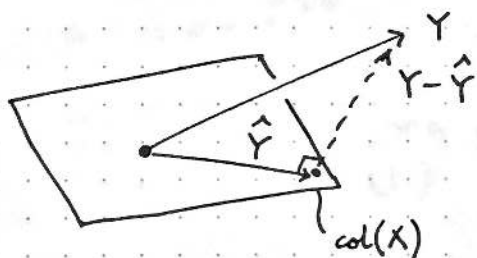
We call P the hat matrix. The residuals are $Y - \hat{Y} = (I - P)Y$.

Lemma P is an orthogonal projection matrix onto $\text{col}(X)$.

By inspection, P is idempotent, symmetric. So an orthog projection.

P acts as the identity on $\text{col}(X)$, maps onto $\text{col}(X)$. \square

\mathbb{R}^n :



Normal Linear model

From now on, we'll assume that

$$E \sim N(0, \sigma^2 I).$$

This will make statistical inference much easier.

L14.1 Normal linear model

We shall assume $\varepsilon \sim N(0, \sigma^2 I_p)$ [ε_i independent]

The log likelihood of this model is

$$\begin{aligned} \ell(\beta, \sigma^2) &= \sum_{i=1}^n \left(-\frac{\log 2\pi}{2} - \frac{\log \sigma^2}{2} - \frac{1}{2\sigma^2} (Y_i - (X\beta)_i)^2 \right) \\ &= \text{constants} - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} S(\beta) \end{aligned}$$

$\leftarrow \|Y - X\beta\|_2^2$

For any value of σ^2 , $\ell(\beta, \sigma^2)$ is maximised as a function of β at the minimiser of $S(\beta)$.

So the MLE for β is $\hat{\beta} = (X^T X)^{-1} X^T Y$.

MLE of σ^2 ?

$$\frac{\partial \ell(\sigma^2, \beta)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} S(\beta)$$

So at the MLE $(\hat{\sigma}^2, \hat{\beta})$,

$$\begin{aligned} -\frac{n}{2\hat{\sigma}^2} + \frac{1}{2\hat{\sigma}^4} S(\hat{\beta}) &= 0 \Rightarrow \hat{\sigma}^2 = \frac{1}{n} S(\hat{\beta}) = \frac{\|Y - X\hat{\beta}\|_2^2}{n} \\ &= \frac{\|(I-P)Y\|_2^2}{n} \end{aligned}$$

Where $P = X(X^T X)^{-1} X^T$, the "hat matrix".

Theorem In the normal linear model,

(i) The fitted values PY , and the residuals $(I-P)Y$ are indpt

(ii) $\hat{\beta} \sim N(\beta, \sigma^2 (X^T X)^{-1})$, $\frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p}^2$, $\hat{\beta} \perp \hat{\sigma}^2$

Proof The vector $\begin{pmatrix} PY \\ (I-P)Y \end{pmatrix} \in \mathbb{R}^{2n}$ is an affine function of ε , and

hence is multivariate normal [check that for $t \in \mathbb{R}^{2n}$, $t^T \begin{pmatrix} PY \\ (I-P)Y \end{pmatrix} \text{ norm}^L$]

Furthermore, $\text{cov} \begin{pmatrix} PY \\ (I-P)Y \end{pmatrix} = \text{cov} \begin{pmatrix} P\varepsilon \\ (I-P)\varepsilon \end{pmatrix} = \sigma^2 \begin{pmatrix} P P^T & P(I-P)^T \\ (I-P)P^T & (I-P)(I-P)^T \end{pmatrix}$

Hence $PY = \hat{Y}$ is \perp of $(I-P)Y$.

Note that $\hat{\sigma}^2$ is a function of $(I-P)Y$. Also, $\hat{\beta}$ is a function of PY . Indeed, $\hat{\beta} = (X^T X)^{-1} X^T Y = (X^T X)^{-1} X^T P Y$.

Therefore, $\hat{\beta} \perp \hat{\sigma}^2$.

Then $\hat{\beta}$ is a linear function of Y , therefore it is MVN.

L14.2

Using the mean and covariance found in the last lecture,

conclude $\hat{\beta} \sim N(\beta, \sigma^2(X^T X)^{-1})$

$$\begin{aligned} \text{Finally, } \frac{n\hat{\sigma}^2}{\sigma^2} &= \frac{1}{\sigma^2} \|(I-P)Y\|_2^2 = \frac{1}{\sigma^2} \|(I-P)(X\beta + \varepsilon)\|_2^2 \\ &= \frac{1}{\sigma^2} \|(I-P)\varepsilon\|_2^2 = \|(I-P) \frac{\varepsilon}{\sigma}\|_2^2 \\ &\sim \chi_{n-p}^2 \cdot 0. \end{aligned}$$

\uparrow rank $n-p$ \uparrow $N(0, I_p)$

Corollary $E\left[\frac{n\hat{\sigma}^2}{\sigma^2}\right] = n-p = \frac{n}{\sigma^2} E[\hat{\sigma}^2]$

$\therefore E[\hat{\sigma}^2] = \frac{n-p}{n} \sigma^2$, biased for small n

We'll define an unbiased estimator of σ^2 , $\tilde{\sigma}^2 = \frac{n}{n-p} \hat{\sigma}^2 = \frac{\|(I-P)Y\|_2^2}{n-p}$

The F distribution

Def If $U \sim \chi_m^2$, $V \sim \chi_n^2$, $U \perp V$, then $X = \frac{U/m}{V/n}$ has an F-distribution with m & n degrees of freedom, $X \sim F_{m,n}$.

NB If $X \sim F_{1,n}$, then $\sqrt{X} \sim t_n$.

As a consequence of the theorem, the following is a pivot:

$$\frac{\frac{1}{p} (\hat{\beta} - \beta)^T (X^T X) (\hat{\beta} - \beta)}{\tilde{\sigma}^2} = \frac{\frac{1}{p\sigma^2} (\hat{\beta} - \beta)^T (X^T X) (\hat{\beta} - \beta)}{\tilde{\sigma}^2 / \sigma^2}$$

- Numerator (function of $\hat{\beta}$) \perp denominator (function of $\tilde{\sigma}^2$)

- The denominator $\tilde{\sigma}^2 / \sigma^2 \sim \chi_{n-p}^2 / (n-p)$

- $(\hat{\beta} - \beta) \sim N(0, (X^T X)^{-1} \sigma^2)$ Let $X^T X = UDU^T$, eigendecomp.

define $(X^T X)^{1/2} = UD^{1/2}U^T$. Then $(\hat{\beta} - \beta)^T (X^T X)^{1/2} \sim N(0, \sigma^2 I_p)$.

Hence $\frac{1}{p\sigma^2} (\hat{\beta} - \beta)^T (X^T X) (\hat{\beta} - \beta)$

$$= \frac{1}{p\sigma^2} [(\hat{\beta} - \beta)^T (X^T X)^{1/2}] [(X^T X)^{1/2} (\hat{\beta} - \beta)] \sim \chi_p^2 / p$$

Putting these together, the pivot has $F_{p, n-p}$ distribution.

Confidence sets for β

$$IP \left(\frac{\frac{1}{p} (\hat{\beta} - \beta)^T (X^T X) (\hat{\beta} - \beta)}{\tilde{\sigma}^2} \leq F_{p, n-p}(\alpha) \right) = 1 - \alpha$$

\uparrow upper α point



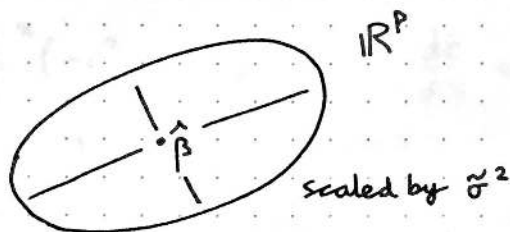
L14.3

The set of β which satisfy the inequality

$$C = \{ b \in \mathbb{R}^p : (\hat{\beta} - b)^T X^T X (\hat{\beta} - b) \leq F_{p, n-p}(\alpha) \hat{\sigma}^2_p \}$$

is an ellipsoid.

$\Rightarrow P_{\beta}(\beta \in C) = 1 - \alpha$ hence C is $100(1 - \alpha)\%$ confidence set for β . \square



Inference for a single coeff. β_j

By the same argument, $(\hat{\beta}_j - \beta_j)^2 / [(X^T X)^{-1}]_{jj} \hat{\sigma}^2 \sim F_{1, n-p}$

Or, taking the square root,

$$\frac{\hat{\beta}_j - \beta_j}{[(X^T X)^{-1}]_{jj}^{1/2} \hat{\sigma}} \sim t_{n-p}$$

Hence, $\hat{\beta}_j \pm t_{n-p}(\frac{\alpha}{2}) \hat{\sigma} [(X^T X)^{-1}]_{jj}^{1/2}$ is an exact $100(1 - \alpha)\%$ confidence interval for β_j .

L15.1 Hypothesis Testing in the Normal Linear Model

$$Y = X\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n)$$

Suppose $\underbrace{X}_{n \times p} = \left[\underbrace{X_0}_{n \times p_0} \mid \underbrace{X_1}_{n \times (p-p_0)} \right]$, and $\beta = \begin{pmatrix} \beta_0 \uparrow \text{in } \mathbb{R}^{p_0} \\ \beta_1 \uparrow \text{in } \mathbb{R}^{p-p_0} \end{pmatrix}$,

where $\text{rank}(X) = p$ and $\text{rank}(X_0) = p_0$.

Want to test $H_0: \beta_1 = 0$ [predictors in X_1 are not associated] to response

$$H_1: \beta_1 \neq 0$$

Under H_0 , the linear model is $Y = X_0\beta_0 + \varepsilon$, so the MLEs are

$$\hat{\beta}_0 = (X_0^T X_0)^{-1} X_0^T Y, \quad \hat{\beta}_1 = 0, \quad \hat{\sigma}^2 = \frac{\|Y - X_0 \hat{\beta}_0\|_2^2}{n} = \frac{1}{n} \|(I - P_0)Y\|_2^2 =: \frac{RSS_0}{n}$$

where $P_0 = X_0(X_0^T X_0)^{-1} X_0^T$ is the orthogonal projection onto $\text{col}(X_0)$.

The generalised likelihood ratio is

$$\begin{aligned} \Lambda_Y(H_0, H_1) &= \frac{(2\pi \hat{\sigma}^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\hat{\sigma}^2} \|Y - X\hat{\beta}\|_2^2\right)}{(2\pi \hat{\sigma}_0^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\hat{\sigma}_0^2} \|Y - X_0 \hat{\beta}_0\|_2^2\right)} = \left(\frac{\hat{\sigma}_0^2}{\hat{\sigma}^2}\right)^{n/2} \\ &= \left(\frac{RSS_0}{RSS}\right)^{n/2} = \left(1 + \frac{RSS_0 - RSS}{RSS}\right)^{n/2} \end{aligned}$$

So generalised LR test rejects when $\Lambda_Y(H_0, H_1)$ is large.

A version of Wilks' theorem says that if p, p_0 fixed, under H_0 ,

as $n \rightarrow \infty$, $2 \log \Lambda_Y(H_0, H_1)$ is approximately $\chi_{p-p_0}^2$.

$\Lambda_Y(H_0, H_1)$ is monotone in the statistic

$$F =: \frac{(RSS_0 - RSS)/(p-p_0)}{RSS/(n-p)}$$

and we can find the exact distribution of F under H_0 .

Thm Under H_0 , $F \sim F_{p-p_0, n-p}$

Proof It's enough to show $RSS \perp RSS_0 - RSS$, $RSS \sim \chi_{n-p}^2$,

$$RSS_0 - RSS \sim \chi_{p-p_0}^2 \leftarrow \sigma^2$$

Note $RSS_0 - RSS = \|(I - P_0)Y\|_2^2 - \|(I - P)Y\|_2^2$
 $= Y^T(I - P_0)Y - Y^T(I - P)Y$
 $= Y^T(P - P_0)Y = \|(P - P_0)Y\|_2^2$

L15.2

where in the last identity we used that $P - P_0$ is idempotent symmetric.

Symmetry is clear from that of P, P_0 .

$$\text{Also, } (P - P_0)(P - P_0) = P^2 - P_0 P - P P_0 + P_0^2 = P + P_0 \overset{\# \text{ epic}}{\leftarrow} - P_0 \overset{\downarrow}{P_0} = P - P_0.$$

Thus $P - P_0$ is an orthogonal projection matrix and

$$\begin{aligned} \text{rank}(P - P_0) &= \text{tr}(P - P_0) = \text{tr} P - \text{tr} P_0 = \text{rank} P - \text{rank} P_0 \\ &= \text{rank}(X) - \text{rank}(X_0) = p - p_0 \end{aligned}$$

Under H_0 , $Y = X_0 \beta_0 + \varepsilon$

$$\begin{aligned} \text{RSS}_0 - \text{RSS} &= \|(P - P_0)Y\|_2^2 = \|(P - P_0)(X_0 \beta_0 + \varepsilon)\|_2^2 \quad \leftarrow \text{dies} \\ &= \|(P - P_0)\varepsilon\|_2^2 \end{aligned}$$

$$\begin{aligned} \text{RSS} &= \|(I - P)Y\|_2^2 = \|(I - P)(X_0 \beta_0 + \varepsilon)\|_2^2 \\ &= \|(I - P)\varepsilon\|_2^2 \end{aligned}$$

Thus by a lemma in L13, $\text{RSS}_0 - \text{RSS} \sim \chi_{p-p_0}^2 \sigma^2$

$$\text{RSS} \sim \chi_{n-p}^2 \sigma^2$$

Finally, note $\begin{pmatrix} (P - P_0)\varepsilon \\ (I - P)\varepsilon \end{pmatrix} \sim N \left(0, \begin{pmatrix} (P - P_0)(P - P_0)^T & (P - P_0)(I - P)^T \\ (I - P)(P - P_0)^T & (I - P)(I - P)^T \end{pmatrix} \right)$

where $(I - P)(P - P_0) = (P - P_0)(I - P) = 0$ because they project onto orthogonal subspaces of \mathbb{R}^n .

So $(P - P_0)\varepsilon \perp (I - P)\varepsilon$ and the result follows. □

So the generalised LRT of size α rejects H_0 when

$$F \geq F_{p-p_0, n-p}(\alpha).$$

How to interpret F ?

RSS is the squared error of the full model

$\text{RSS}_0 - \text{RSS}$ is the reduction in squared error when we include X_1

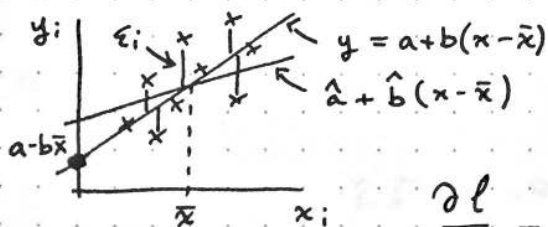
The statistic $R^2 = \frac{\text{RSS}_0 - \text{RSS}}{\text{RSS}}$ is "the proportion of the variance in the responses explained by the predictors in X_1 "

The F statistic is a version of R^2 scaled by degrees of freedom

Simple linear regression

② $Y_i = a + b(x_i - \bar{x}) + \varepsilon_i$ where $\bar{x} = \frac{1}{n} \sum x_i$, $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$
 2 parameters a, b and also σ^2

Ex Sheet 3 How to define the design matrix X for this model.



The log-likelihood is

$$l(a, b, \sigma^2) = \sum_{i=1}^n \left\{ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y_i - a - b(x_i - \bar{x}))^2 \right\}$$

$$\frac{\partial l}{\partial a} = \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - a - b(x_i - \bar{x}))$$

$$= \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - a)$$

③ this is zero at $\hat{a} = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$.

$$\frac{\partial l}{\partial b} = \sum_{i=1}^n \frac{1}{2\sigma^2} (y_i - a - b(x_i - \bar{x})) (x_i - \bar{x})$$

Plug in \hat{a} for a , set to zero, solve for b

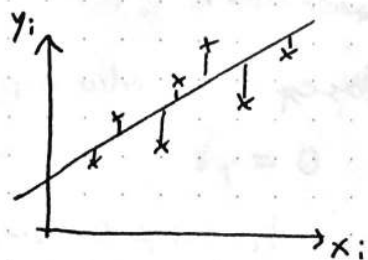
$$\hat{b} = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Applications of the general normal linear model

- ① Simple linear model $Y_i = a + b(x_i - \bar{x}) + \varepsilon_i$, $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$

Last time derived the MLEs $\hat{a} = \bar{Y}$

$$\hat{b} = \frac{\sum (Y_i - \bar{Y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$



We wish to test $H_0: b=0$ vs $H_1: b \neq 0$.

Under H_0 , the MLE is $\hat{a} = \bar{Y}$, $\hat{b} = 0$.

The generalised LR test rejects when F is large, where $F = \frac{(RSS_0 - RSS) / (p - p_0)}{RSS / (n - p)}$.

- In this case $RSS_0 = \sum_i (Y_i - \bar{Y})^2$, $RSS = \sum_i (Y_i - \bar{Y} - \hat{b}(x_i - \bar{x}))^2$

$$R^2 = \frac{RSS_0 - RSS}{RSS_0} = \frac{[\sum (Y_i - \bar{Y})(x_i - \bar{x})]^2}{\sum (Y_i - \bar{Y})^2 \sum (x_i - \bar{x})^2} \leftarrow \text{which?}$$

oh ok he made a mistake

this is the square of Pearson's correlation coefficient.

The F statistic $F = R^2 \frac{n-p}{p-p_0} = R^2 \cdot \frac{n-2}{2-1}$ has, under H_0 , an

$F_{p-p_0, n-p} = F_{1, n-2}$ distribution.

② Analysis of Variance (ANOVA)

Example Patients in a clinical trial are randomly assigned to groups

- $1, \dots, I$ each of size J (balanced design)

Y_{ij} response of patient j in group i

$$Y_{ij} = \alpha + \mu_i + \varepsilon_{ij}, \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$$

We'll assume (WLOG) that group 1 is a "control" group, let $\mu_1 = 0$.

How to interpret parameters:

- α is the mean response in the control group,
- μ_i is the difference in mean response between group i and the control group.

- How to write this as a linear model $Y = X\beta + \varepsilon$?

$$Y = \begin{pmatrix} Y_{11} \\ Y_{12} \\ \vdots \\ Y_{1j} \\ Y_{21} \\ \vdots \\ Y_{Ij} \end{pmatrix} = \begin{pmatrix} 1 & 0 & \dots & \dots & 0 & 0 \\ 1 & 0 & \dots & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots & \vdots \\ 1 & 0 & \dots & \dots & 0 & 0 \\ 1 & 1 & \dots & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots & \vdots \\ 1 & 0 & \dots & \dots & 0 & 1 \end{pmatrix} \begin{pmatrix} \alpha \\ \mu_2 \\ \mu_3 \\ \vdots \\ \mu_I \end{pmatrix} + \begin{pmatrix} \epsilon_{11} \\ \epsilon_{12} \\ \vdots \\ \epsilon_{1j} \\ \epsilon_{21} \\ \vdots \\ \epsilon_{Ij} \end{pmatrix}$$

The matrix has one column with all 1s, the subsequent columns are indicators (0,1) for groups 2,3,..., I.

Equality of means test $H_0: \mu_2 = \mu_3 = \dots = \mu_I = 0$ vs $H_1: \mu_i \neq 0$ for $i \neq 1$

experimental treatments have no effect

Let $X_0 = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$, $X = (X_0 \ X_1)$, $P = X(X^T X)^{-1} X^T$,

$P_0 = X_0(X_0^T X_0)^{-1} X_0^T$

$PY = \begin{pmatrix} \bar{Y}_1 \\ \vdots \\ \bar{Y}_2 \\ \vdots \\ \bar{Y}_I \end{pmatrix}$, $\bar{Y}_i = \sum_j \frac{Y_{ij}}{J}$ "the average response for i"

$P_0 Y = \begin{pmatrix} \bar{Y} \\ \vdots \\ \bar{Y} \end{pmatrix}$, $\bar{Y} = \sum_i \frac{\bar{Y}_i}{I} = \sum_{ij} \frac{Y_{ij}}{IJ}$

$F = \frac{(RSS_0 - RSS)/(p - p_0)}{RSS/(n - p)} = \frac{\frac{1}{p - p_0} \|(P_0 - P)Y\|_2^2}{\frac{1}{n - p} \|(I - P)Y\|_2^2}$

$n = IJ$, $p = I$, $p_0 = 1$

$F = \frac{I(J-1)}{I-1} \cdot \frac{\sum_{ij} (\bar{Y} - \bar{Y}_i)^2 \rightarrow \text{"between group" variance}}{\sum_{ij} (Y_{ij} - \bar{Y}_i)^2 \rightarrow \text{"within group" variance}}$

The statistic has, under H_0 , an $F_{p-p_0, n-p} = F_{I-1, I(J-1)}$ distⁿ

Choice of control group If we chose another control group, $col(X)$, $col(X_0)$ would be the same. Thus, P , P_0 would be the same, the fitted values PY , $P_0 Y$ would be the same (under H_0, H_1), and the F statistic would be the same.

The interpretation of the parameters α, μ changes.

L16.3

③ Paired observations Consider an ANOVA model with $I=2$,

● but the people in the trial come in pairs and one member of each pair is assigned to each group, (at random).

Y_{ij} = response for person in j^{th} pair, group i

$$Y_{ij} = \alpha + \mu_{X_i} + \delta_j + \varepsilon_{ij}, \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$\mu_1 = 0, \delta_1 = 0$$

We care about the estimates parameter μ_2 , not about $\delta_2, \dots, \delta_J$.

We can test $H_0: \mu_2 = 0$ vs $H_1: \mu_2 \neq 0$ using the F statistic