

STATISTICAL MODELLING

Statistical model is collection of probability distributions

$\{P_\theta : \theta \in \Theta\}$ on a given sample space.

↳ statistical model called IDENTIFIABLE if $\theta \mapsto P_\theta$ one-to-one,

i.e. $P_{\theta_1} = P_{\theta_2} \Rightarrow \theta_1 = \theta_2$, so distinct parameter values give rise to distinct distributions. (no 2 different θ lead to same observed distribution)

BIAS: given estimator T , $\text{bias}_\theta(T) = E_\theta(T) - \theta$

↳ if T estimator for $g(\theta)$, then $\text{bias}(T) = E(T) - g(\theta)$.

STANDARD ERROR: $SE_\theta(T) = \sqrt{\text{Var}_\theta(T)} = \sqrt{E_\theta(T^2) - E_\theta(T)^2}$

MEAN SQUARE ERROR: $MSE_\theta(T) = E_\theta[(T - \theta)^2]$.

$$\Rightarrow MSE_\theta(T) = \text{Var}_\theta(T) + \text{bias}_\theta(T)^2.$$

CRAMER-RAO LOWER BOUND: suppose $T = T(X)$ unbiased:

$$\text{Var}_\theta(T) \geq \frac{1}{I(\theta)}, \text{ where } I(\theta) = E_\theta \left[\left(\frac{\partial}{\partial \theta} \log f_\theta(X) \right)^2 \right] \\ = - E_\theta \left[\frac{\partial^2}{\partial \theta^2} \log f_\theta(X) \right].$$

NB: $f_\theta(x)$ joint p.d.f. of $x = (x_1, \dots, x_n)$.

For x_1, \dots, x_n and $f_\theta^{(1)}$ pdf of single observation,

$$\Rightarrow \underline{I_f(\theta) = n \cdot I_{f_\theta^{(1)}}(\theta)}. \Rightarrow I \propto \text{sample size for random sample.}$$

JENSEN'S: for convex functⁿ g and r.v. X ,
 $g(E(X)) \leq E(g(X))$.

- $X_n \xrightarrow{\text{a.s.}} X$ if $P(\lim_{n \rightarrow \infty} X_n = X) = 1$.
- $X_n \xrightarrow{P} X$ if $\forall \varepsilon > 0, P(|X_n - X| > \varepsilon) = 0$.
- $X_n \xrightarrow{D} X$ if $\lim_{n \rightarrow \infty} P(X_n \leq x) = F_X(x) = P(X \leq x)$ where X has c.d.f. F_X .

(a.s. \Rightarrow P \Rightarrow D.)

CONSISTENT: sequence of estimators T_n consistent if

$$\forall \theta \in \Theta, T_n \xrightarrow{P_\theta} g(\theta) \rightarrow T_n \text{ converging in probability.}$$

PORTMANTEAU LEMMA:

$$X_n \xrightarrow{D} X \iff E(f(X_n)) \rightarrow E(f(X)).$$

for all bounded + continuous $f: \mathbb{R} \rightarrow \mathbb{R}$.

ASYMPTOTICALLY UNBIASED: T_n for $g(\theta)$ if $E(T_n) \rightarrow g(\theta)$.

If T_n asymptotically unbiased and $\text{Var}_\theta(T_n) \rightarrow 0$,
 $\Rightarrow T_n$ consistent for $g(\theta)$.

proof by Markov's inequality:
 $P(|X| \geq a) \leq \frac{E(|X|)}{a}$.

ASYMPTOTICALLY NORMAL: T_n for θ if $\sqrt{n}(T_n - \theta) \xrightarrow{D} N(0, \sigma^2(\theta))$

CENTRAL LIMIT THM: Y_1, \dots, Y_n i.i.d., $E(Y_i) = \mu$, $\text{Var}(Y_i) = \sigma^2$

$$\Rightarrow \sqrt{n}(\bar{Y} - \mu) \xrightarrow{D} N(0, \sigma^2).$$

\hookrightarrow i.e. sample averages are asymptotically normal.

SLUTSKY'S: if $X_n \xrightarrow{D} X$ and $Y_n \xrightarrow{P} c$ for constant c ,

$$\bullet X_n + Y_n \xrightarrow{D} X + c$$

$$\bullet Y_n X_n \xrightarrow{D} cX, \quad \bullet Y_n^{-1} X_n \xrightarrow{D} c^{-1}X.$$

DELTA METHOD: T_n asymptotically normal, so $\sqrt{n}(T_n - \theta) \xrightarrow{D} N(0, \sigma^2(\theta))$,

$$\Rightarrow \sqrt{n}(g(T_n) - g(\theta)) \xrightarrow{D} N(0, g'(\theta)^2 \cdot \sigma^2(\theta)) \quad (g' \neq 0)$$

(note: odds of event A happening defined as $P(A) / (1 - P(A))$.)

CONTINUOUS MAPPING THM: convergence in $d, P, a.s.$ all preserved under continuous mappings.

MLE: parameter value $\theta \in \Theta$ for which observed data is most likely.

$$L(\theta) = \prod_{i=1}^n f(y_i; \theta) \rightarrow \text{product of pdfs} \rightarrow \text{likelihood fun}^n.$$

HYPOTHESIS TEST: H_0 : for which values of sample X_1, \dots, X_n to accept H_0 , or reject otherwise + accept H_1 .

↳ subset of sample space to reject H_0 is called CRITICAL REGION.

TYPE I ERROR: false +ve \rightarrow reject H_0 when H_0 true.

TYPE II ERROR: false -ve \rightarrow accept / don't reject H_0 when H_0 false.

α -level test where $P_\theta(\text{reject } H_0) \leq \alpha$. (small α .)

POWER FUNCTION: $\beta(\theta) = P_\theta(\text{reject } H_0)$

↳ if $\theta \in \Theta_0$ then want $\beta(\theta)$ be small.

↳ if $\theta \in \Theta_1$ then want $\beta(\theta)$ be large.

p-value: reject H_0 iff $p \leq \alpha$ (p-value reported from observed values)

construct test for confidence region: for region $A(Y)$ of $1-\alpha$

and test $H_0: \theta \in \Theta_0$, $H_1: \theta \notin \Theta_0$, reject H_0 if:

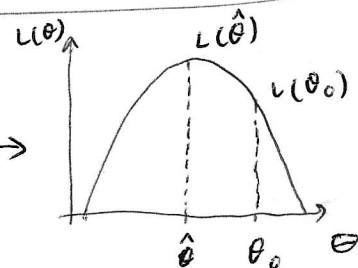
$\Theta_0 \cap A(Y) = \emptyset$ i.e. if none of elements of H_0 are in confidence region.

LIKELIHOOD RATIO TEST: $t(\underline{y}) = \frac{\sup_{\theta \in \Theta} L(\theta; \underline{y})}{\sup_{\theta \in \Theta_0} L(\theta; \underline{y})} = \frac{\text{max lik. under } H_0 + H_1}{\text{max lik. under } H_0}$

for observed data \underline{y} .

↳ reject H_0 if $t(\underline{y})$ large i.e. $t(\underline{y}) \geq k$.

compare $L(\hat{\theta})$ to $L(\theta_0)$ and if $L(\theta_0) < L(\hat{\theta})$, θ_0 likely one.



DEFⁿ of χ^2 : For $X_1, X_2, \dots, X_n \sim N(0, 1)$ i.i.d. $\Rightarrow \sum_{i=1}^n X_i^2 \sim \chi_n^2$.

LRT DISTRIBUTION: $2 \log t(\underline{y}) \xrightarrow{D} \chi_r^2$ (LRT statistic asymptotically χ^2)

where $r = \#$ of indep params in full model - $\#$ of indep params under H_0 .

(pf: ① Taylor's expansion of $\log L(\theta)$. ② Slutsky's + continuous mapping thm + MLE asymptotically normal + WLLN.)

MAXIMUM LIKELIHOOD ESTIMATOR: $\hat{\theta}$ is $\hat{\theta}$ such that

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

e.g. for $\text{Bern}(\theta)$ and $\text{Pois}(\theta)$, MLE is \bar{X} and for $\text{exp}(\theta)$, MLE is $\frac{1}{\bar{X}}$.

MLE not necessarily unbiased.

MLE functionally invariant: bijective g then $\hat{\phi} = g(\hat{\theta})$ MLE of $\phi = g(\theta)$.

if $\hat{\theta}_n$ MLE of X_1, \dots, X_n sequence, then $\hat{\theta}_n$ asymptotically normal

$$\Rightarrow \sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D} N\left(0, \frac{1}{I_f(\theta)}\right) \quad \text{for } I_f(\theta) = E\left[\left(\frac{\partial}{\partial \theta} \log f_{\theta}(X)\right)^2\right].$$

\hookrightarrow for $\hat{\theta}_n = \frac{1}{n} \sum X_i \Rightarrow \sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D} N(0, \frac{\sigma^2}{n})$ where $\sigma^2 = \text{Var}(X_i)$
 \hookrightarrow CLT.

CONFIDENCE REGION: random interval containing true parameter with probability of $1 - \alpha$.

$$\hookrightarrow P(\theta \in I) \geq 1 - \alpha.$$

construct confidence interval by pivotal quantity.

PIVOTAL QUANTITY: $t(\underline{Y}, \theta)$ s.t. distribution of t completely known + doesn't depend on any unknown params.

• by asymptotic normality of MLE, $\sqrt{n} \frac{(T_n - \theta)}{\sigma(\theta)} \xrightarrow{D} N(0, 1)$

$\Rightarrow \sqrt{n} \frac{(T_n - \theta)}{\sigma(\theta)} \sim N(0, 1)$ approx. can be used as pivotal quantity.

BONFERRONI CORRECTION: for confidence intervals $[L_i, U_i]$ of θ_i ($i=1, \dots, k$) of $1 - \alpha/k$. $\Rightarrow (L_1, U_1) \times \dots \times (L_k, U_k)$ $1 - \alpha$ CI for $(\theta_1, \dots, \theta_k)^T$.

(e.g. if $[L_1, U_1]$ is $0.99 = 1 - 0.01$ and $[L_2, U_2]$ is $0.97 = 1 - 0.03$ then $[L_1, U_1] \times [L_2, U_2]$ is $1 - 0.01 - 0.03 = 0.96$. ($1 - \alpha_1 - \alpha_2$))

(simple)

LINEAR MODEL: $Y_i = \beta_0 + \alpha_1 \beta_1 + \epsilon_i$, $i=1, \dots, n$

and goal how to "estimate" linear parameters β_1, β_2 (and also $\text{Var}(\epsilon_i) = \sigma^2$)?

STATISTICAL ERROR ϵ_i : amount by which observation differs from its expected value $\rightarrow E(Y_i) = \beta_0 + \alpha_1 \beta_1$.

LINEAR ALGEBRA lemma: $X_{n \times p}$ matrix, $\text{rank}(X^T X) = \text{rank}(X)$.

for $\underline{X} = (X_1, \dots, X_n)^T$, $E(\underline{X}) = (E(X_1), \dots, E(X_n))^T$.

for deterministic $\underline{A}, \underline{B}$: $E(\underline{A}\underline{X}) = \underline{A} E(\underline{X})$ and $E(\underline{X}^T \underline{B}) = E(\underline{X})^T \underline{B}$.

Cov($\underline{X}, \underline{Y}$) = $(\text{Cov}(X_i, Y_j))_{i,j} = E(\underline{X}\underline{Y}^T) - E(\underline{X}) \cdot E(\underline{Y})^T$

$(\text{Cov}(\underline{X}) = \text{Cov}(\underline{X}, \underline{X}))$
 $(\text{Cov}(\underline{X}) = \text{Var}(\underline{X}))$

• $\text{Cov}(X, Y) = \text{Cov}(Y, X)^T$ • $\text{Cov}(aX + bY, Z) = a \text{Cov}(X, Z) + b \text{Cov}(Y, Z)$

• $\text{Cov}(AX, BY) = A \text{Cov}(X, Y) B^T$ • $\text{Cov}(\underline{X})$ +ve definite + symmetric.

if $\text{Cov}(X, Y) = 0$, X, Y same distribution $\Rightarrow X, Y$ independent.

X, Y independent $\Rightarrow \text{Cov}(X, Y) = 0$.

(general)

LINEAR MODEL: $\underline{Y} = \underline{X} \underline{\beta} + \underline{\epsilon}$ $\rightarrow \underline{X} \in \mathbb{R}^{n \times p}$, $\underline{\beta} \in \mathbb{R}^p$, $\underline{Y} \in \mathbb{R}^n$.

$E(\underline{\epsilon}) = \underline{0} \Rightarrow E(\underline{Y}) = \underline{X} \underline{\beta}$ (\underline{X} design matrix)

SECOND ORDER ASSUMPTION: $\text{Cov}(\underline{\epsilon}) = \sigma^2 I_n$ \rightarrow errors of 2 different observat^{ns} independent + variance of all error identical.

NORMAL THEORY ASSUMPTION: $\underline{\epsilon} \sim N(\underline{0}, \sigma^2 I_n)$. \rightarrow used to construct tests.

FULL RANK: X has full rank of p (since assuming $n > p$).

LEAST SQUARES: $S(\underline{\beta}) = (\underline{Y} - \underline{X}\underline{\beta})^T (\underline{Y} - \underline{X}\underline{\beta})$ \leftarrow want to minimise w.r.t. $\underline{\beta}$.

\Rightarrow LS eqⁿ: $X^T X \hat{\underline{\beta}} = X^T \underline{Y}$

Solution exists iff $(X^T X)^{-1}$ exists iff $\text{rank}(X^T X) = p = \text{rank}(X)$.

$\hookrightarrow \text{Cov}(\hat{\underline{\beta}}) = \sigma^2 (X^T X)^{-1}$.

GAUSS-MARKOV THM: under FR, so A , $\forall \underline{c} \in \mathbb{R}^p$:

estimator $\underline{c}^T \hat{\underline{\beta}}$ has smallest variance out of all unbiased estimator $\underline{c}^T \underline{\beta}$.

P projection matrix iff $\underline{P}^T = \underline{P}$ and $\underline{P}^2 = \underline{P}$.

\bullet P projection onto L given by: $\underline{P} = \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T$

with $\underline{X} = (\underline{x}_1 \dots \underline{x}_r)$ for $\underline{x}_1, \dots, \underline{x}_r$ basis of L .

\hookrightarrow $I - P$ projects to L^\perp .

LEMMA: for $n \times n$ proj. matrix \underline{P} w/ rank r
 \hookrightarrow \underline{P} has r evales of 1 and $n-r$ evales of 0.
 \hookrightarrow rank $\underline{P} = \text{trace } \underline{P}$.

$\hat{\underline{Y}} = \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y} = \underline{P} \underline{Y}$; $\underline{e} = \underline{Y} - \hat{\underline{Y}} = \text{vector of residuals.}$

RESIDUAL SUM of SQUARES: $\text{RSS} = \underline{e}^T \underline{e} = \sum_{i=1}^n e_i^2 = \underline{Y}^T \underline{Q} \underline{Y}$

\hookrightarrow it is minimum value of $S(\underline{\beta})$. for $\underline{Q} = \underline{I} - \underline{P}$.

$\hat{\sigma}^2 = \frac{\text{RSS}}{n - p}$ is unbiased estimator for σ^2 , where $\text{cov}(\underline{\varepsilon}) = \sigma^2 \underline{I}_n$
 $n - p \leftarrow p = \text{rank } \underline{X}$.

COEFFICIENT of DETERMINATION: $R^2 = 1 - \frac{\text{RSS}}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{\text{RSS}}{\text{RSS in interest only model}}$

\hookrightarrow more parameter, lower RSS,
 smaller RSS \Rightarrow larger R^2 , $0 \leq R^2 \leq 1$ shows "fit" of model.

MULTIVARIATE NORMAL: for $X_1, \dots, X_r \sim N(0, 1)$ i.i.d. and $\underline{\mu} \in \mathbb{R}^n$, $\underline{A} \in \mathbb{R}^{n \times r}$

$\Rightarrow \underline{z} = \underline{A} \underline{X} + \underline{\mu} \sim N(\underline{\mu}, \underline{A} \underline{A}^T)$.

if $\underline{z} \sim N(\underline{\mu}, \underline{\Sigma})$, then $\underline{A} \underline{z} + \underline{b} \sim N(\underline{A} \underline{\mu} + \underline{b}, \underline{A} \underline{\Sigma} \underline{A}^T)$

LEMMA: if $\underline{z} = \begin{pmatrix} z_1 \\ \vdots \\ z_k \end{pmatrix} \sim N(\underline{\mu}, \underline{\Sigma})$ and $\underline{\Sigma} = \begin{pmatrix} A_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & A_k \end{pmatrix}$
 $\Rightarrow z_1, \dots, z_k$ independent.

uncorrelated
 jointly normal
 \Downarrow
 independent

NON-CENTRAL χ^2 : $\underline{z} \sim N(\underline{\mu}, I_n) \Rightarrow \underline{U} = \underline{z}^T \underline{z} = \sum z_i^2 \sim \chi_n^2(\delta)$

where $\delta = \sqrt{\underline{\mu}^T \underline{\mu}}$ and $\underline{\mu} \in \mathbb{R}^n$.

if $U \sim \chi_n^2(\delta)$, then $E(U) = n + \delta^2$ and $\text{Var}(U) = 2n + 4\delta^2$

NON-CENTRAL t : $X \sim N(\delta, 1)$ and $U \sim \chi_n^2 \Rightarrow \frac{X}{\sqrt{U/n}} \sim t_n(\delta)$

F-DISTRIBUTION: $W_1 \sim \chi_{n_1}^2$ and $W_2 \sim \chi_{n_2}^2$
 $\Rightarrow F = \frac{W_1/n_1}{W_2/n_2} \sim F_{n_1, n_2}(\delta)$ (ratio of 2 chi-squared)
 \downarrow
cannot take -ve values.

for +ve semidefinite, symmetric $\underline{A} \in \mathbb{R}^{n \times n}$, $\underline{A} = \underline{L}\underline{L}^T$ and $\underline{L}^T \underline{L} = \text{diag}(\text{non-zero values of } \underline{A})$.

if $\underline{X} \sim N(\underline{\mu}, I)$, \underline{A} +ve semidefinite, symmetric and \underline{B} s.t. $\underline{B}\underline{A} = 0$
 $\Rightarrow \underline{X}^T \underline{A} \underline{X}$ and $\underline{B} \underline{X}$ are independent.

Pf: since if r.v. X indep of Y then $g(X)$ indep of $f(Y)$.

if $\underline{z} \sim N(\underline{\mu}, I_n)$ and $\underline{A}_1, \underline{A}_2$ proj matrices with $\underline{A}_1 \underline{A}_2 = 0$
 $\Rightarrow \underline{z}^T \underline{A}_1 \underline{z}$ and $\underline{z}^T \underline{A}_2 \underline{z}$ independent.

FISHER-COCHRAN THM: $\underline{A}_1, \dots, \underline{A}_k$ proj matrices s.t. $\sum \underline{A}_i = I_n$
and $\underline{z} \sim N(\underline{\mu}, I_n)$, then $\underline{z}^T \underline{A}_1 \underline{z}, \dots, \underline{z}^T \underline{A}_k \underline{z}$ all independent
AND $\underline{z}^T \underline{A}_i \underline{z} \sim \chi_{r_i}^2(\delta_i)$ where $r_i = \text{rank } \underline{A}_i$ and $\delta_i^2 = \underline{\mu}^T \underline{A}_i \underline{\mu}$
($\delta_i = \sqrt{\underline{\mu}^T \underline{A}_i \underline{\mu}}$)

• under NTA, $\underline{Y} \sim N(\underline{X}\beta, \sigma^2 I_n)$.

• MLE for σ^2 is $\hat{\sigma}^2 = \frac{RSS}{n} \rightarrow$ biased since $\frac{RSS}{n-p}$ unbiased.

$\frac{RSS}{\sigma^2} \sim \chi_{n-r}^2$ where $r = \text{rank } X$ ← PIVOTAL QUANTITY for σ^2

TEST for all component of $\underline{\beta}$: $\frac{\underline{c}^T \hat{\underline{\beta}} - \underline{c}^T \underline{\beta}}{\sqrt{\underline{c}^T (X^T X)^{-1} \underline{c} \cdot \frac{RSS}{n-p}}} \sim t_{n-p}$
 (under FR, NTA.) \rightarrow where $p = \text{rank } X$

e.g. if we want to test for $\beta_3 = 0$, then $\underline{c} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}$.

equivalently, $\frac{\underline{c}^T \hat{\underline{\beta}} - \underline{c}^T \underline{\beta}}{\sqrt{\underline{c}^T (X^T X)^{-1} \underline{c} \sigma^2}} \sim N(0, 1)$ if σ^2 known.

F-TEST: testing for more than one component of $\underline{\beta}$.

$$F = \frac{RSS_0 - RSS}{RSS} \cdot \frac{(n-r)}{(r-s)} \sim F_{r-s, n-r} \quad \text{where } r = \text{rank } X \\ s = \text{rank } X_0.$$

\hookrightarrow i.e. RSS of reduced model - RSS of full model / RSS of full model.

we use $RSS = \underline{y}^T \underline{Q} \underline{y}$ and $RSS_0 = \underline{y}^T \underline{Q}_0 \underline{y}$ ($\underline{Q}_0 = \underline{I} - \underline{P}_0$)
 onto $\text{span}(X_0)^\perp$

(\hookrightarrow reject if $F > c$ at α -sig level where $P(X > c) = \alpha$ for $X \sim F_{r-s, n-r}$.)

OUTLIERS: look for residuals that are too large!

$\frac{e_i}{\sqrt{(1-p_{ii})\sigma^2}} \sim N(0, 1)$ where \underline{P} is proj onto $\text{span } X$
 \leftarrow if dk σ^2 , then plug in unbiased $\sigma^2 = \frac{RSS}{n-p}$.

$\text{Cov}(\underline{e}) = \sigma^2(\underline{I}_n - \underline{P})$ and $\text{Var}(e_i) = \sigma^2(1 - p_{ii})$
 \leftarrow known or leverage.

COOK'S DISTANCE: $D_i = \frac{(\hat{\beta}_{(i)} - \hat{\beta})^T X^T X (\hat{\beta}_{(i)} - \hat{\beta})}{p \cdot RSS / (n-p)}$ where $\hat{\beta}_{(i)}$ is OLS with ith observation removed.

\hookrightarrow how much $\hat{\beta}$ changes if I remove observation i .

OR $D_i = r_i^2 \cdot \frac{p_{ii}}{(1-p_{ii})r}$ $r = \text{rank } X$
 $r_i^2 = \frac{e_i^2}{(1-p_{ii})\sigma^2}$ standardized residual.