

Today's outline

Midterm

Hypothesis , classification

Presentation

Q 2

(b) 2 r.v.s X, Y

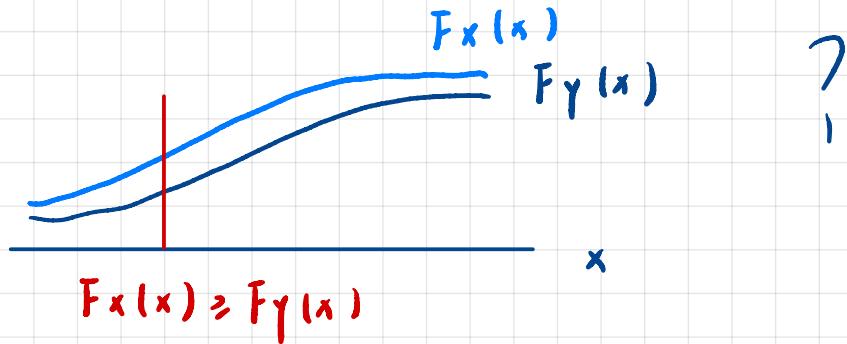
$$P(X \leq Y) = 1$$

$$F_X(x) \geq F_Y(x) \quad \forall x \quad \text{True or False?}$$

$$F_Y(x) = P(Y \leq x)$$

$$\leq P(X \leq x)$$

$$\{Y \leq x\} \subset \{X \leq x\} = F_X(x) \Rightarrow \text{True}$$



(c) $X, Y, \max\{X, Y\}$

$$F_{\max\{X, Y\}}(x) = F_X(x) \cdot F_Y(x), \text{ given } X \perp\!\!\!\perp Y$$

$$F_{\max\{X, Y\}}(x)$$

$$= P(\max\{X, Y\} \leq x)$$

$$= P(\{X \leq x\} \cap \{Y \leq x\})$$

$$= P(X \leq x) \cdot P(Y \leq x)$$

$$F_X(x) \quad F_Y(x)$$

Q3 (c)

$x \perp\!\!\! \perp y$

x, y jointly normal and $\text{cov}(x, y) = 0$ True

$$\Rightarrow E(y|x) = E(y)$$

(E) $x \perp\!\!\! \perp y$, $\text{cov}(x^2, y) = 0$ True

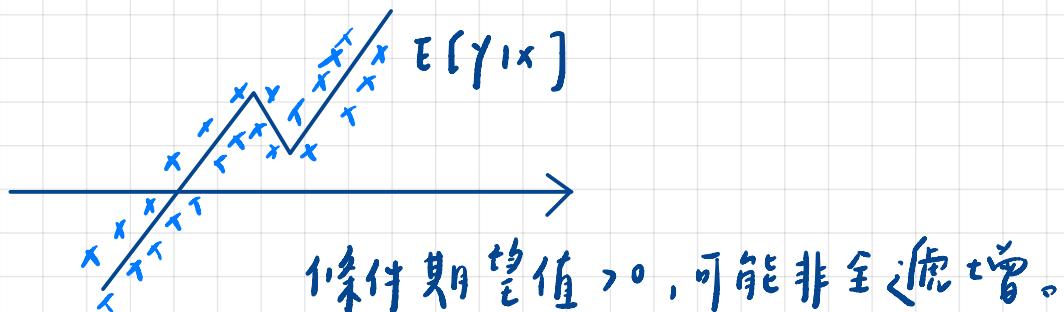
Q4 (B)

$$E(y|x) = ax + b \text{ obs} \rightarrow E(y|x) \text{ True}$$

Ans: obs : best linear predictor

$E(y|x)$ best predictor

(E) $\text{Cov}(x, y) > 0 \Rightarrow E(y|x=x) \text{ increases in } x$ False



Q5

$$\Rightarrow \text{Var}(x_1) > \text{Var}(x_2)$$

$$(x_1, x_2) \quad a_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad (x_2, 2x_1) \Rightarrow \tilde{a}_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix} ?$$

a_1 = solution

$$\Rightarrow \text{Var}(2x_1) > \text{Var}(x_2)$$

$$\max_a \text{Var}(a^T x) \quad \text{s.t. } a^T a = 1$$

Q6

$$(x_1, 2x_2) \quad \tilde{a}_2 = \begin{pmatrix} ? \end{pmatrix}$$

$$\tilde{a}_2^T \tilde{a}_1 = 0 \Rightarrow \tilde{a}_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

Q8

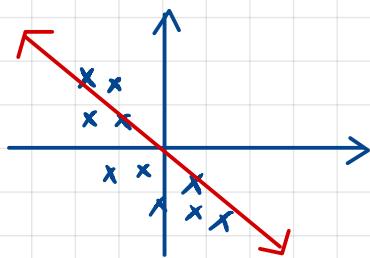
$$x_1, x_2, \quad a_1 = \begin{pmatrix} 3/5 \\ 4/5 \end{pmatrix} \quad \text{Cov}(x_1, x_2) = ?$$

>
|
<

$$\text{Var}(a_1^T x) = \text{Var}(a_{11}x_1 + a_{12}x_2)$$

$$= \text{Var}(a_{11}x_1) + \text{Var}(a_{12}x_2) + 2a_{11}a_{12}\text{Cov}(x_1, x_2)$$

if $\text{Cov}(x_1, x_2) < 0 \Rightarrow \text{var}(a_{12}x_2) < 0 \nmid \max \text{Var}$



Q9. 10. 11

Find the pair of closest point for the $d(x, y)$ given

13, 14, 15

Find out $\hat{\beta}_1, \hat{\beta}_2$

for $y = \beta_1 x_1 + \beta_2 x_2 + \epsilon$ given data

Sol(1)

$\hat{\beta}_1, \hat{\beta}_2$ are sample means for $x_1 = 1$ & $x_2 = 1$

Sol(2), (3)

$\hat{\beta}_1, \hat{\beta}_2$ are the solutions

$$\min_{\beta_1, \beta_2} \frac{1}{n} \sum (y_i - \beta_1 x_{i1} - \beta_2 x_{i2})^2$$

solve by F.O.C or 配方法

↳ 抛物線最小值

Sol(4)

out has closed form

$$\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = (x^T x)^{-1} x^T Y$$

Hypothesis Testing

H_0 : Statement 1

H_1 : Statement 2

真項只有一個，which one?

We have data $X = X_1, X_2, \dots, X_n$

A test is data-driven decision rule $D(x)$

$$D(x) \in \{ \begin{array}{l} \text{Reject } H_0 \\ \text{Accept } H_0 \end{array} \}$$

A good test = $\begin{cases} \text{Reject } H_0 \text{ when } H_0 \text{ false} \\ \text{Accept } H_0 \text{ when } H_0 \text{ true} \end{cases}$

Example: a fair coin

x = outcome of coin toss

$x \sim \text{Ber}(p)$

$$p \in [0, 1]$$

$$H_0: p = 0.5$$

$$H_1: p \neq 0.5$$

How to construct a test $D(x)$?

Suppose we have $x_1, x_2, \dots, x_n \sim \text{Ber}(p)$

An intuitive test

$$D(x) = \begin{cases} \text{reject } |x_n - 0.5| > c \\ \text{accept } H_0 \text{ o.w.} \end{cases}$$

$$\Rightarrow \text{what } c = ?$$

a type I error

if c too small, likely to reject H_0 when H_0 is true.

If c too large, likely to accept H_1 even if H_0 false.

y

type II error

It's impossible to get rid of type-I & type-II errors

but we can calculate the probability of error.

Recall:

$$\text{size} = p(\text{type-I error})$$

$$\text{power} = 1 - p(\text{type-II error})$$

Next Step

$$n = 10, c = 0.15$$

$b(x) = \begin{cases} \text{reject if } |\bar{x}_n - 0.5| > 0.15 \\ \text{accept } \alpha. \omega \end{cases}$

$$p(\text{Type I error}) \quad p = 0.5$$

$$p(b(x) = \text{rejects})$$

$$\bar{x}_n \in \{0, 0.1, 0.2, \dots, 1\}$$

$$p(\text{reject}) = p(|\bar{x}_n - 0.5| > 0.15)$$

$$= p(\bar{x}_n \in \{0, 0.1, 0.2, 0.3, 0.7, 0.8, 0.9, 1\})$$

$$= 1 - p(\bar{x}_n \in \{0.4, 0.5, 0.6\})$$

$$p(\text{Type I error})$$

$$= 1 - p(\bar{x}_n = 0.4) - p(\bar{x}_n = 0.5) - p(\bar{x}_n = 0.6)$$

$$= 1 - C_4^{10} 0.5^4 0.5^6 - C_5^{10} 0.5^5 0.5^5 - C_6^{10} 0.5^6 0.5^4$$

Remark:

$p(\text{Type-I error})$ depends on c

e.g. $c = 0.08$ ($n = 10$)

$|\bar{x}_n - 0.5| > 0.08$ rejects

when $n=10$ $\{|\bar{x}_n - 0.5| < 0.08\} = \{\bar{x}_n = 0.5\}$

$$\text{size} = 1 - P(\bar{x}_n = 0.5)$$

c 越小 \rightarrow size 越大

$P(\text{type II error})$

$$p \neq 0.5 \quad n=10, c=0.15$$

$P(D(x) \text{ Accepts})$

$$= P(\bar{x}_n = 0.6) + P(\bar{x}_n = 0.5) + P(\bar{x}_n = 0.4)$$

$$C_6^{10} p^6 (1-p)^4 + C_5^{10} p^5 (1-p)^5 + C_4^{10} p^4 (1-p)^6$$

Remark:

1. $P(\text{type II error})$ depends on c

$\Rightarrow c$ 越大 $P(\text{type II error})$ 越大

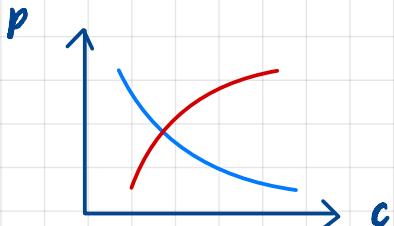
2. $P(\text{type II error})$ depends on p

① $p = 0.49$

② $p = 0.01$

3. There is a trade-off between

$P(\text{type-I})$ & $P(\text{type-II})$



Problem of exact calculations when $n = 1000$

$$\frac{c^{1000}}{500} 0.5^{500} 0.5^{500}$$

Such terms will appear . . .

Error probability approximation with CLT

CLT

$$\sqrt{n}(\bar{x}_n - p) \xrightarrow{d} N(0, p(1-p))$$

$$\Rightarrow \frac{\sqrt{n}(\bar{x}_n - p)}{\sqrt{p(1-p)}} \stackrel{\sim}{=} N(0,1) \text{ by CLT (} n \rightarrow \infty \text{)}$$

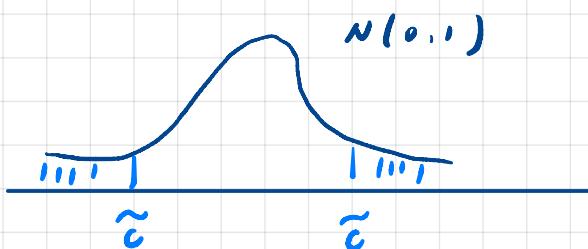
$P(\text{type I error}) = P(b|x) = \text{reject}$

$$P = 0.5$$

$$= P(|\bar{x}_n - 0.5| > c)$$

$$= P\left(\frac{\sqrt{n}(\bar{x}_n - 0.5)}{\sqrt{0.5(1-0.5)}} > \frac{\sqrt{n}c}{0.5}\right)$$

$$= P\left(|N(0,1)| > \frac{\sqrt{n}c}{0.5}\right) = 2\Phi\left(-\frac{\sqrt{n}c}{0.5}\right)$$



$$P(|N(0,1)| > \tilde{c}) = 2P(N(0,1) < -\tilde{c}) = 2\Phi(-\tilde{c})$$

\downarrow
 $N(0,1) \text{ CDF}$

(when n large)

Given n, c

$p(\text{type I error})$

$$= 2 \Phi \left(-\frac{\sqrt{n}c}{0.5} \right) \stackrel{\text{want}}{=} 0.05$$

when $n = 1000$ want $p(\text{type I error}) = 0.05$

solve $\Phi(-2\sqrt{1000}c) = 0.05$

take $\Phi^{-1}(\dots)$

$$-2\sqrt{1000}c = \Phi^{-1}(-0.025)$$

$$c = \frac{\Phi^{-1}(-0.025)}{-2\sqrt{1000}}$$

$$\therefore \frac{-1.96}{-2\sqrt{1000}} \approx \frac{-1.96}{-64} \div 0.03 = 3\%$$

case study : 藍白合

Assume our goal is to find the best pair

HK v.s KH

侯柯 v.s 柯侯

Q1 HK v.s 朝

$$X_{HK} = 1, X_{HK} \sim \text{Ber}(p_{HK})$$

p_{HK} : 侯柯支持率

Q2 KH v.s 朝

$$X_{KH} \sim \text{Ber}(p_{KH})$$

$$\left\{ \begin{array}{l} H_0: p_{HK} > p_{KH} \quad (\Rightarrow \text{讓柯侯強除非強力證據說柯侯強}) \\ H_1: p_{HK} < p_{KH} \end{array} \right.$$

Recall: we always prioritize β (type II error) = 0.5

侯柯強 but test says 柯侯強

$$H_0: p_{HK} > p_{KH}$$

A test if $\hat{p}_{HK} - \hat{p}_{KH} < c, c > 0$ we reject

$$P(\text{type I error}) = P(|\hat{p}_{HK} - p_{KH}| < -c) \stackrel{\text{Want}}{=} 0.05$$

朱立倫： if $|\hat{p}_{HK} + 0.03 - (p_{KH} - 0.03)| < 0$

$$\Leftrightarrow p_{HK} - p_{KH} < -0.06$$

$$P(|\hat{p}_{HK} - p_{HK}| > 0.03) = 0.05 \quad \left. \right\} \text{True}$$

$$P(|\hat{p}_{KH} - p_{KH}| > 0.03) = 0.05$$

~~X~~

$$\text{"So" } P(\hat{p}_{HK} - \hat{p}_{KH} < 0.06) = 0.05$$

$$0.03 = 0.1 \text{ (柯言謬太多)}$$

柯太哲

$$= 0.01 \text{ (柯言謬太多)}$$

協議：

誤差內 \Rightarrow 候柯 +1

is ambiguous

1. 什麼民調？

2. 誤差幾 %？

3. 40

$$H_0: p_{HK} \geq p_{KH}$$

$$H_0: p_{HK} \geq p_{LKH}$$

$$p_{KH} > p_{LKH}$$

Error calculation with concentration inequality

Concentration inequality

"It's unlikely that X deviate from μ a lot"

e.g. Chebyshev inequality

$$P(|X - E[X]| > M) \leq \text{Var}(X)/M$$

(上界)

X "concentrates" on $E[X]$

Hoeffding inequality for Bernoulli r.v.

Thm. Let $X_1, X_2, \dots, X_n \sim \text{Ber}(p)$

$$P(|\bar{X}_n - p| > \epsilon) \leq 2e^{-2n\epsilon^2} \quad \forall p$$

$$(\epsilon = 0.1) \quad n = 10 \quad n = 100 \quad n = 200 \quad n = 500$$

$$2e^{-2n\epsilon^2} \quad 1.63 \quad 0.27 \quad 0.039 \quad \leq 10^{-6}$$

沒用

機率 < 1

n 越大 ... ?

$$P(\text{type I error}) = P(|\bar{X}_n - 0.5| > c) \leq 2e^{-2nc^2} \stackrel{\text{Want}}{=} 0.05$$

by Hoeffding

$$e^{-2nc^2} = 0.05/2$$

$$-2nc^2 = \ln(0.025)$$

$$c^2 > \ln(0.025) / -2n$$

when $n = 1000$, $c \approx 0.03$

n 夠大，樣本不會離真實值太遠

		Actual		$TP = \text{true positive}$
		cat	dog	
predict	cat	TP	FP	
	dog	FN	TN	

$$\text{precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$