

10.12
(四) W3

Review
如何用 CLT

理論層次
假設 $H_0: \mu \leq 2.5 \text{ ppm}$ 對立 $H_1: \mu > 2.5 \text{ ppm}$ 互斥
描述母體的特徵
 H_0 null 無異 於 ≤ 2.5

假設 H_0 是對的
則從母體逐次抽取 成千上萬次 抽取 from 母體 Population
每次抽取 n ($n \geq 30$)
求得的樣本平均數為 \bar{x}
將 \bar{x} 拿來作圖

樣本平均數 (\bar{x}) sample

CLT
critical value
 $\bar{x} = 2.5$
拒絕區
若 H_0 是 true

高工字: 理論層次: Reference
可以不倚模樣 (Repeated)

一般而言, 是要拒絕 H_0 , 要支持 $H_1: \mu > 2.5 \text{ ppm}$ → 收集資料去拒絕

$\bar{x} \sim N(2.5, \frac{\sigma^2}{n})$ 事實上, 只有一個樣本 n .
一個樣本平均數 \bar{x} 成千上萬次的其中一次
假設我重複抽 m 次
 $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_m = \bar{x}$
 $\mu = 2.5 \text{ ppm}$
 $\frac{\sigma^2}{n} = \frac{(\bar{x}_1 - \mu)^2 + (\bar{x}_2 - \mu)^2 + \dots + (\bar{x}_m - \mu)^2}{m}$

我們看看 \bar{x} 落在 \bar{x} 的抽樣分配列表
如果 \bar{x} 反決母體 μ , $\bar{x} = 2.7$ 不在 Reject 區 → 不拒絕 H_0 (無異於 2.5)
則 \bar{x} 會落在 \bar{x} 抽樣分配的中部地區
如果 \bar{x} 不是對 H_0 , → 黑色虛線區 → H_0 不成立
則 \bar{x} 應該會落在 \bar{x} 抽樣分配的拒絕區

→ 代入 data, 看其落在何處
估計是反證法, 回推

不拒絕 H_0 (無異於 2.5)

拒絕區

(follow)

四個抽樣分配

$X \sim N(\mu, \sigma^2)$ 常態分布的 2 個參數

i.i.d. 獨立同質分配
z 分配
縮變

$f(x) | \mu, \sigma^2 = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ 常態 (高斯分配)

反曲夫

標準化 (Repeated)
 $Z = \frac{x - \mu}{\sigma}$
 $f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$

從 $X \sim N(\mu, \sigma^2)$ 隨機抽樣 n , 計算 \bar{x}

1 st	$x_1, x_2, \dots, x_n \rightarrow \bar{x}_n$	$\frac{\bar{x}_n - \mu}{\frac{\sigma}{\sqrt{n}}}$
2 nd	$x_1, x_2, \dots, x_n \rightarrow \bar{x}_n$	$\frac{\bar{x}_n - \mu}{\frac{\sigma}{\sqrt{n}}}$
...
m th	$x_1, x_2, \dots, x_n \rightarrow \bar{x}_n$	$\frac{\bar{x}_n - \mu}{\frac{\sigma}{\sqrt{n}}}$

卡方分配 χ^2

很重要
 $Z^2 = \text{卡方}$

從 $X \sim N(\mu, \sigma^2)$ 隨機抽樣, 離均差平方和

1st $x_1, x_2, \dots, x_{1+n} \rightarrow \sum_{i=1}^n (x_i - \bar{x}_n)^2 / \sigma^2$

2nd ...

...

Mth $\rightarrow \sum_{i=1}^n (x_{mi} - \bar{x}_n)^2 / \sigma^2$

卡方 自由度 $df = n-1$

$S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$
過去所學之變異數公式

$\frac{(n-1)S^2}{\sigma^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sigma^2}$ 卡方公式

z 分配平方是卡方

χ^2 < p.p >

$Z_i = \frac{x_i - M}{\sigma}$ $Z^2 = (\frac{x_i - M}{\sigma})^2 = \chi^2_{(n)}$

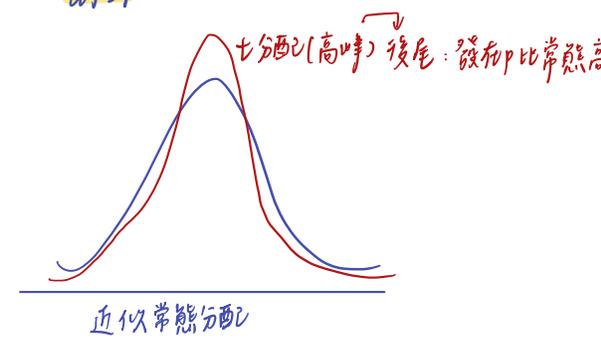
$\sum_{i=1}^n Z_i^2 = \sum_{i=1}^n (\frac{x_i - M}{\sigma})^2 = \chi^2_{(n)}$
 $= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - M)^2 = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - M)^2 = \frac{1}{\sigma^2} [\sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - M)^2]$
 $= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \bar{x})^2 + (\frac{\bar{x} - M}{\frac{\sigma}{\sqrt{n}}})^2$
 統計量
 自由度 = $k(x)$
 $df = 1$

* 補充 回家查
布拉格分配

$Z^2 = \frac{(O-E)^2}{?}$

從 $X \sim N(\mu, \sigma^2)$ 隨機抽樣

$t = \frac{\bar{x} - M}{\frac{s}{\sqrt{n}}}$ (以前學的) \times
 $t = \frac{Z}{\sqrt{\frac{\chi^2}{df}}} = \frac{\frac{\bar{x} - M}{\frac{\sigma}{\sqrt{n}}}}{\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}}$ (Z)
 平方
 平方



用於 迴歸分析

F 分配

用於 ANOVA, 迴歸分析

來自 2 常態分配

單因子變異數分析
2 個變異數相除

SV	SS	MS	F
組間	SSB	MSB	$\left[\begin{matrix} = SS^* \\ + \\ = SW^* \end{matrix} \right]$
組內	SSW	MSW	
total			Mean Square = S^2 (變異數)

$\sum (y_i - \bar{y})^2 = \sum (y_i - \bar{y})^2 + \sum \sum (\bar{y}_i - \bar{y})^2$
 互斥

迴歸分析中的
變異數分析

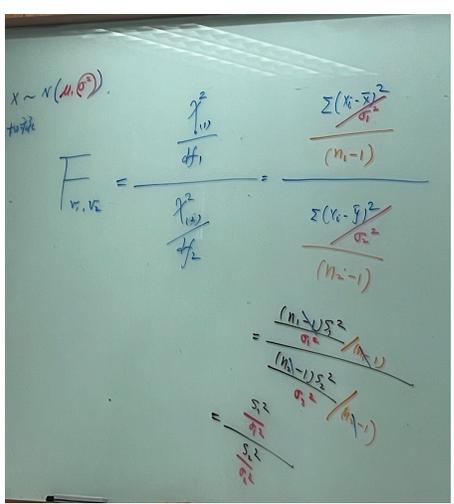
SV	df	MS	F
迴歸 $\sum (\hat{y}_i - \bar{y})^2 = SSR$	#	$\frac{SSR}{\#}$	$\left[\begin{matrix} \\ + \\ \text{误差} \end{matrix} \right]$
誤差 $\sum (y_i - \hat{y}_i)^2 = SSE$	$n-1-\#$	$\frac{SSE}{n-1-\#}$	

四個抽樣分配使用時機

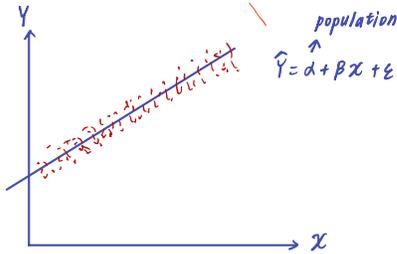
Z 分配: $n \geq 30$. 標準常態分配

t 分配: $n \leq 30$. p 隨 df 不同而改變

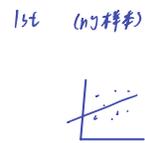
$F_{v_1, v_2} =$



Sample



population
 $\hat{Y} = \alpha + \beta X + \epsilon$
 $(\epsilon \sim N(0, \sigma^2))$

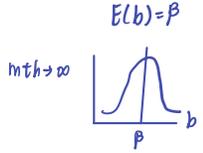


$\hat{Y} = a + bx$
 $y = a + b_1x + e$

$x \rightarrow y$ 沒有效果
 \rightarrow 身高的解釋

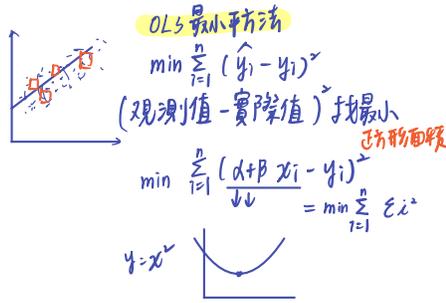
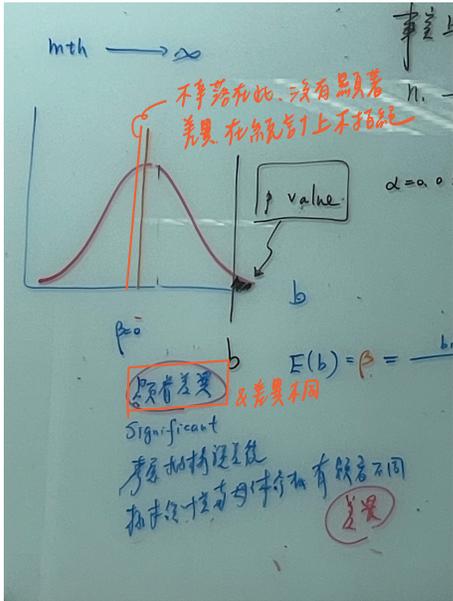


$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$



$$\frac{b_1 + b_2 + b_3 + \dots + b_m}{m} = \beta$$

且為常態分配



MLE / OLS
 Ordinary
 Least
 Square / Estimate

$$\hat{\beta} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$= \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

$$= \frac{\text{長方形面積}(xy)}{\text{正方形面積}(x)}$$

MLE = OLS
 求解尋找 α, β
 一定有分配的假設

- Today:
- 分配的正態
 - 四種分配
 - 抽樣分配

Q3, Q4 \rightarrow 理想狀態

Q5 \rightarrow TRUELY

Q5

$H_0: \mu \geq 100$ $n = 16$ $\bar{x} = 98.9375$

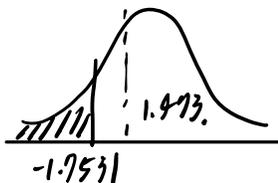
$H_1: \mu < 100$ $s = 2.794$
 2.886

$$t = \frac{\bar{x} - \mu}{s / \sqrt{n}} = \frac{98.9375 - 100}{2.886 / 4} = \frac{-1.0625}{0.7215}$$

查七表 \rightarrow $= 1.473$

$t_{0.05, 15} = 1.7531$

$1.473 < 1.7531$, 不拒絕



$$\frac{70}{69} = 1.31$$

10.26 統計一定會有「誤差項」

4.3 Measure of fit
使每個點到相減平方和為 min

好的參數估計值包括 μ .

① 不偏性: $E(s^2) = \sigma^2$
 $E\left(\frac{\sum(x_i - \bar{x})^2}{n-1}\right) = \frac{\sum(x_i - \mu)^2}{n}$

② 有效性
(有最小的變異誤有 min 的標準誤)

③ 充分性: 用到 all 樣本資訊

④ 一致性: $n \uparrow, B \rightarrow \beta$ 趨近到條件 β

$s^2 = \frac{SSY}{df_y}$ 自由度
 $= \frac{\sum(y_i - \bar{y})^2}{n-1}$
 = mean square (均方)

11.2
 $y\text{-score} = \frac{\text{read} + \text{math}}{2}$
 $x = \text{生師比 STR}$) 連續變項

$y_i = a + bx_i + e_i = \hat{y}_i + e_i$

$\hat{y}_i = a + b x_i$ if $b=0$ ($= a = \bar{y}$)

Sum of Square (y 的變異)

$SS_T = SS_Y = SS_R + SS_E$

$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2$

$= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$

$df = 2, 1 \text{ for } a \& b \quad df = 1 \text{ from } \bar{x}$

$= SS_{\text{ERROR}} + SS_{\text{Reg}}$

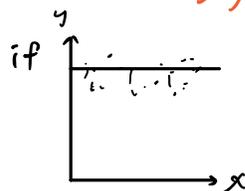
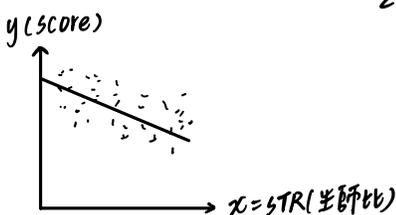
連續

$x = \begin{cases} 1 & \text{if } STR < 20 \\ 0 & \text{o.w. 當參考組} \end{cases}$

無論連續 or dummy

$b = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}$ <不變>

= dummy
 = 兩個類別 (0,1)

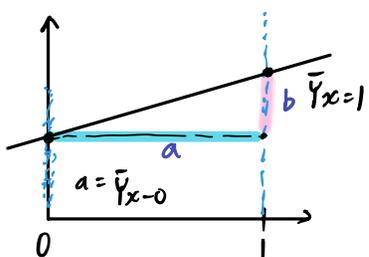


x 對 y 沒有影響.
 No pattern

R^2 (可以解釋 y 的變異程度, 決定係數)

$= \frac{SSR}{SS_T}$ or $1 - \frac{SSE}{SS_T}$

Dummy



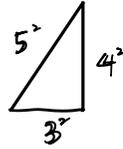
$= \bar{y}_{x=1} - \bar{y}_{x=0}$
 = 2個獨立樣本平均數的差
 = t-test

帶去解平單其變異

OVERALL TEST

(卡方)

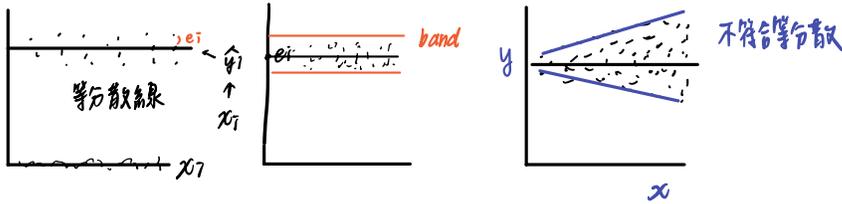
$$F_{\frac{df_R}{df_E}} = \frac{\frac{SSR}{df_R}}{\frac{SSE}{df_E}} = \frac{\frac{(n-1) \cdot s_1^2}{df_1(p)}}{\frac{(n-1) \cdot s_2^2}{df_2(n-p-1)}}$$



變異數分析 → 2個 s^2 相除

適配度檢定 (F-Test)
以圖卡方, 獨立

sum of square



Lesson 2
信賴區間

If H_0 是對的話.

Ch 6 多元迴歸

$X_1 = STR$

$X_2 = lunch$

$X_3 = English$

$$y = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3 + b_4 X_4 + e$$

$H_0: \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$ $H_1: \text{至少一個 } \beta \neq 0$

null = 無異於 0

到台有一變項, 斜率 $\neq 0$
但不知道是哪的.

(Overall)

F 用來檢定 H_0

$$SS_T = SS_Y = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$= \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2$$

$$= SSE + SSR$$

迴歸係數 → partial regression

解釋 "Marginal"

X_2, X_3, X_4 在 model 中, 當 2 個人有相同

設其他條件不變的時候, X 增加一單位, y 可以增加 b_1 個單位

多元共線性: 指的是 X_1 之間資訊變化的程度 or X_1 之間彼此被 Copy 的程度

Another

$$\hat{y}_1 = \alpha_0 + \alpha_1 X_1 + \alpha_2 X_2 + \alpha_3 X_3 + \alpha_4 X_4$$

①

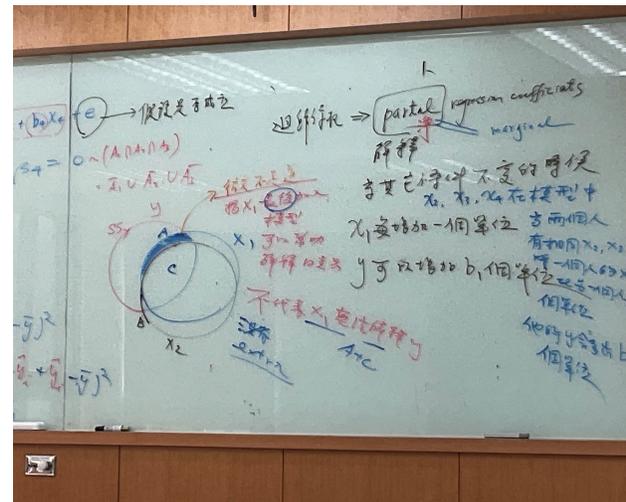
$$VIF_{X_1} = \frac{1}{1 - R_{X_1}^2}$$

有幾個解釋變數就有幾個 VIF.

(Variance Inflation factor)
變異數 膨脹

②

CN C Condition Index > 30



$$R^2_{adjusted} = R^2 - \text{pseudo } R^2$$

$$R^2 = 1 - \frac{SSE}{SST} \rightarrow R^2_{adjusted} = 1 - \frac{n-1}{n-p-1} (1 - R^2) = \frac{S_y^2 - s_e^2}{S_y^2} = 1 - \frac{s_e^2}{S_y^2}$$

多元迴歸 > 1.6 箭, type "text" 才有報表.