

巨量資料管理學院碩士在職專班

統計分析

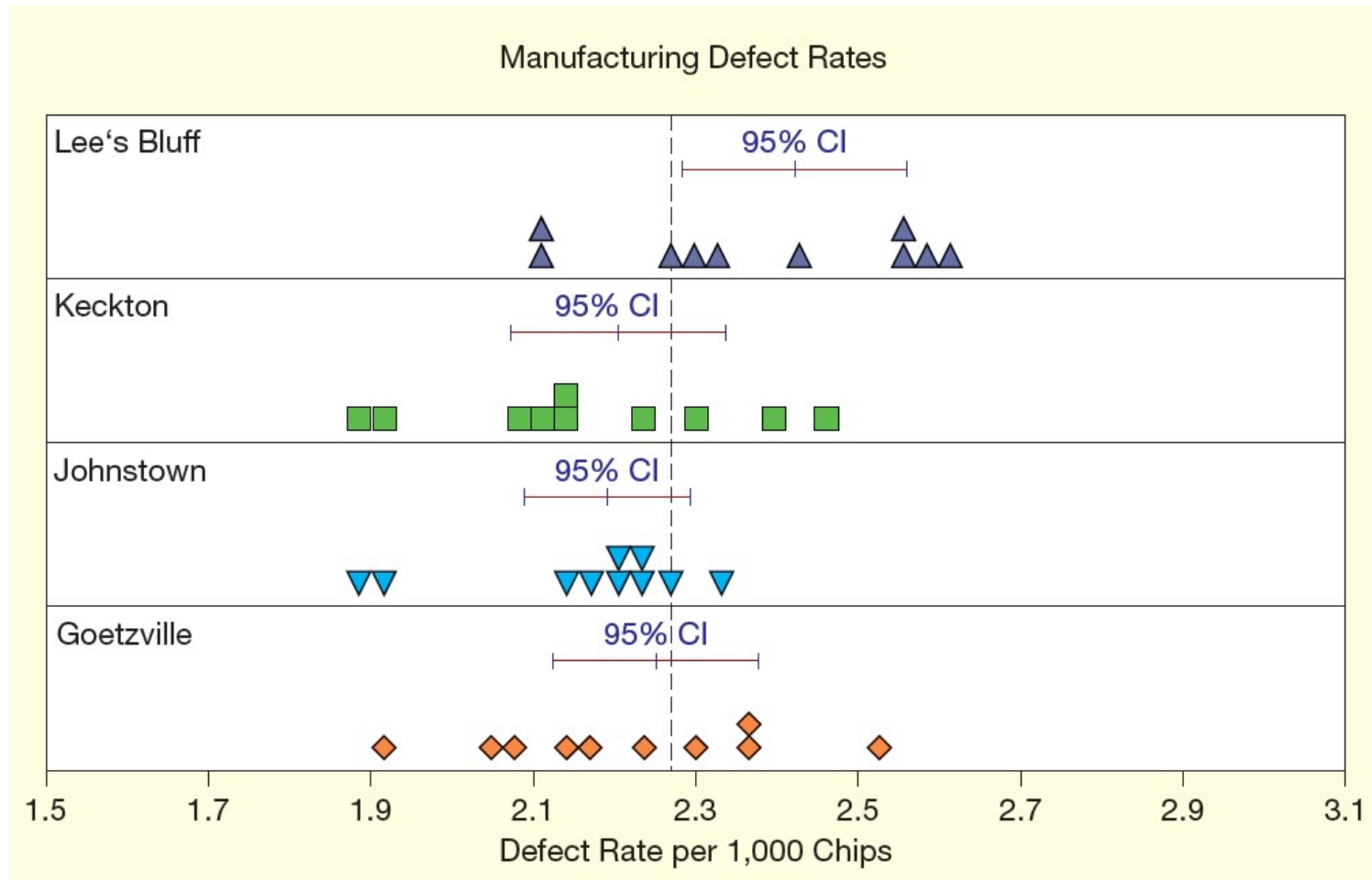
2022/11/18

陳光宏

變異數分析 (Analysis of variance)

多個獨立樣本的檢定

Analysis of variance (ANOVA)

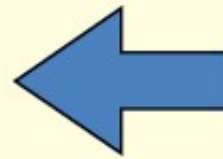


涉及兩個變項

*Dependent variable
(numerical)*

$Y = \text{defect rate}$

*may be
affected by*



*Independent variable
(categorical)*

Treatment (plant location)

$T_1 = \text{Lee's Bluff}$

$T_2 = \text{Keckton}$

$T_3 = \text{Johnstown}$

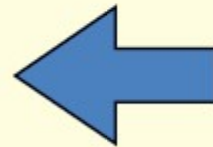
$T_4 = \text{Goetzville}$

*Dependent variable
(numerical)*

*may be
affected by*

*Independent variable
(categorical)*

$Y = \text{length of stay}$



Treatment (fracture type)

$T_1 = \text{facial}$

$T_2 = \text{radius or ulna}$

$T_3 = \text{hip or femur}$

$T_4 = \text{other lower extremity}$

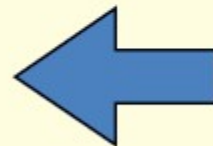
$T_5 = \text{all other}$

*Dependent variable
(numerical)*

*may be
affected by*

*Independent variable
(categorical)*

$Y = \text{paint viscosity}$



Treatment (temperature)

$T_1 = \text{low (15°C)}$

$T_2 = \text{medium (20°C)}$

$T_3 = \text{high (25°C)}$

多個獨立樣本的檢定

Analysis of variance (ANOVA)

- 比較三組或三組以上的平均值
- 能否沿用兩個獨立樣本檢定的概念？
- 要解決什麼問題？
 - 假說怎麼建立？
 - 檢定統計量怎麼計算？

$$H_0: \mu_1 - \mu_2 = 0$$

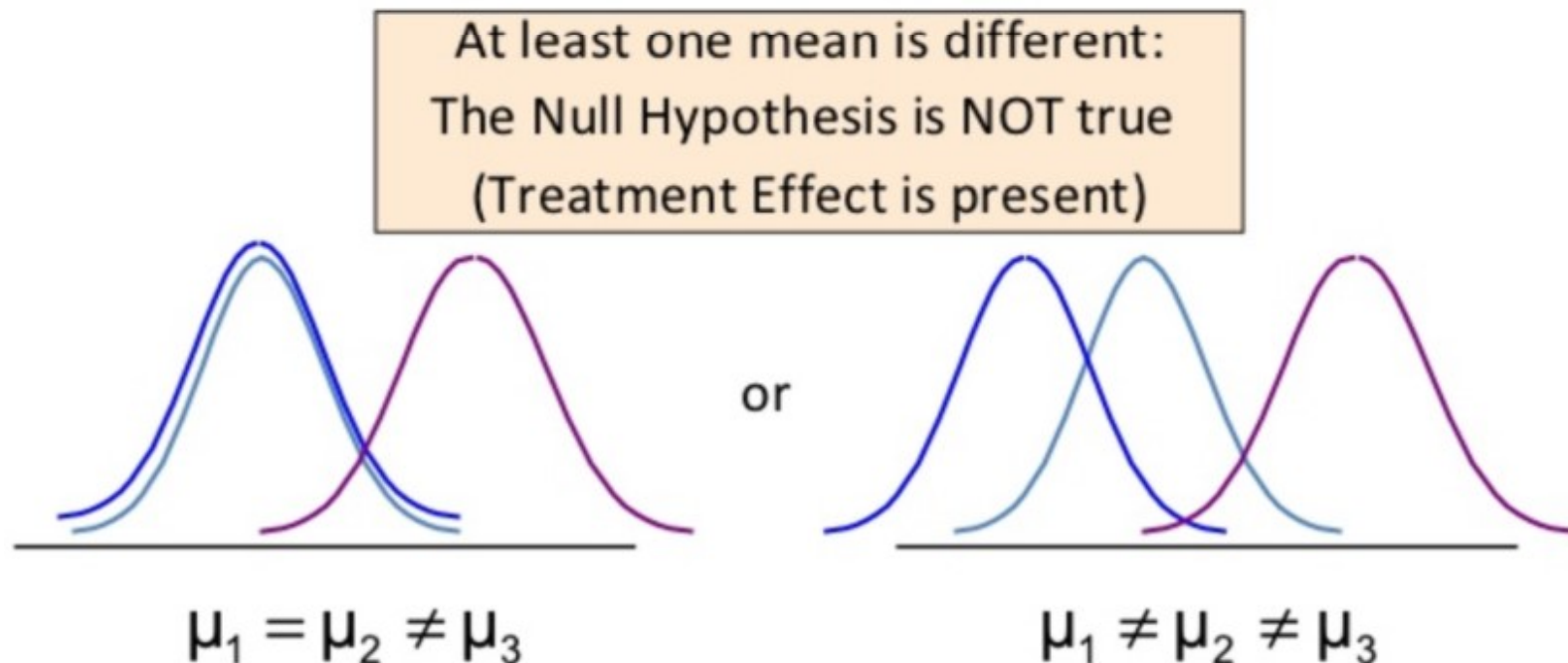
$$H_1: \mu_1 - \mu_2 \neq 0$$

假說建立

- 虛無假說：三個平均值均相等
- 對立假說：至少有兩個平均值不相等

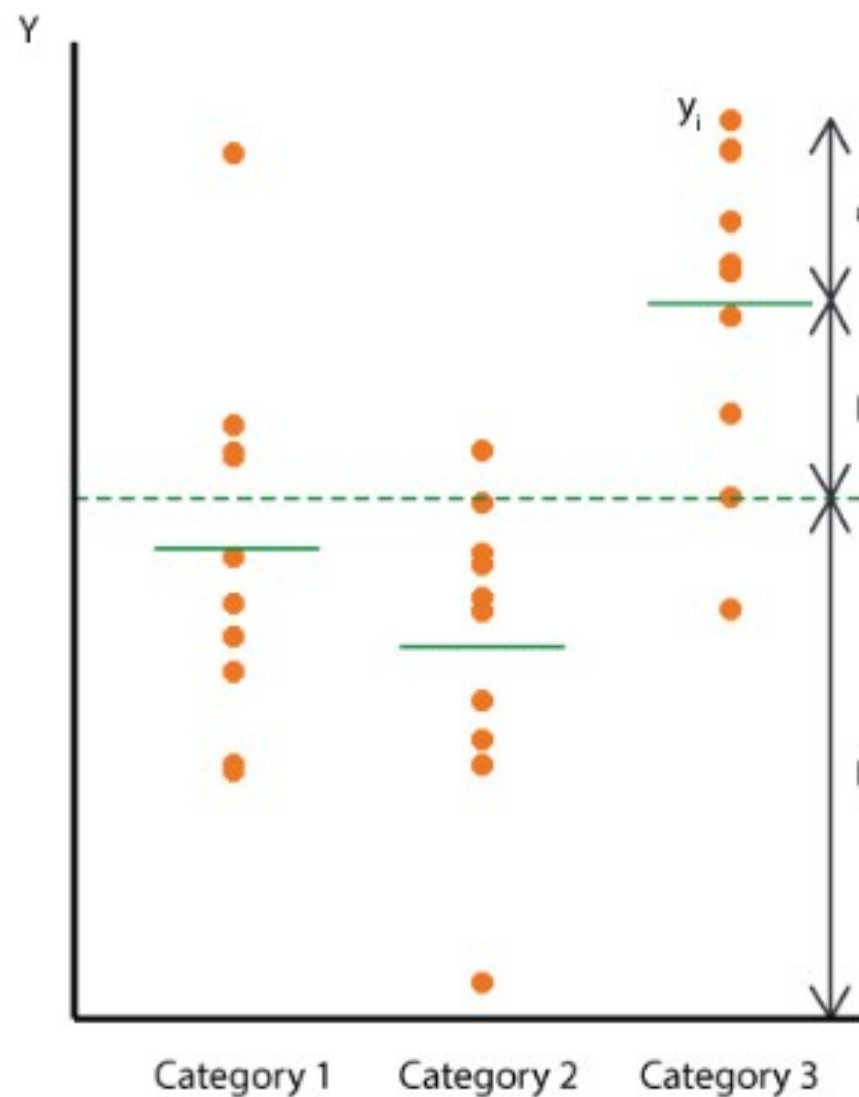
$$H_0 : \mu_1 = \mu_2 = \mu_3$$

$$H_1 : \mu_1 \neq \mu_2 \neq \mu_3$$



概念

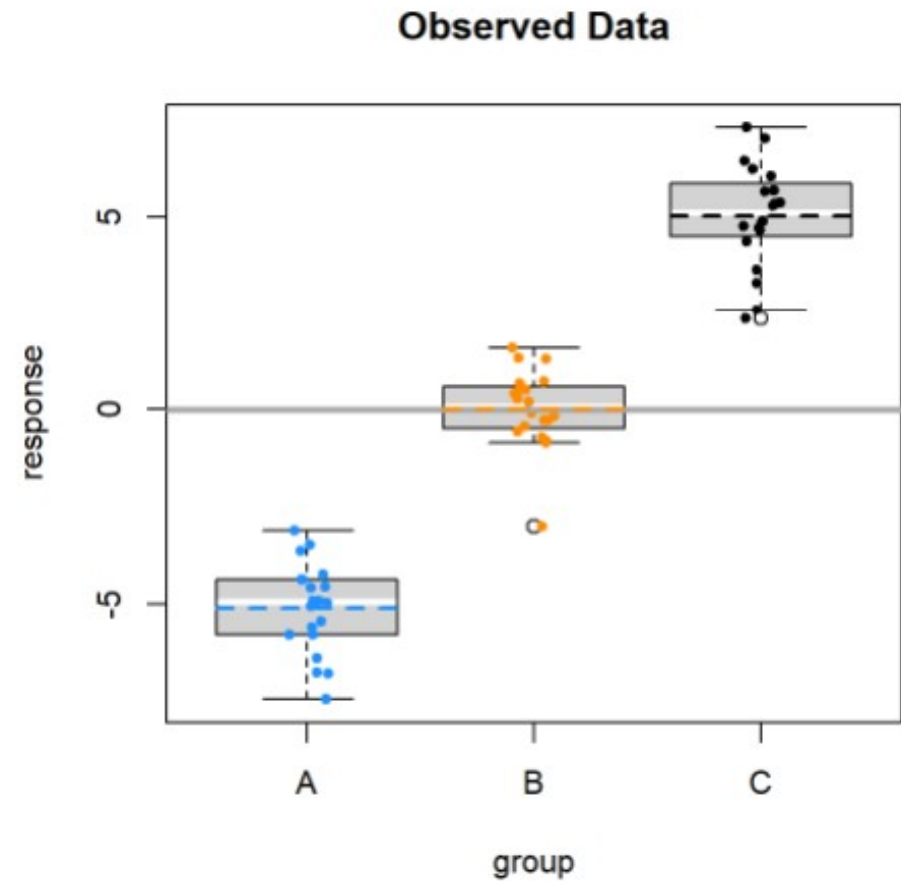
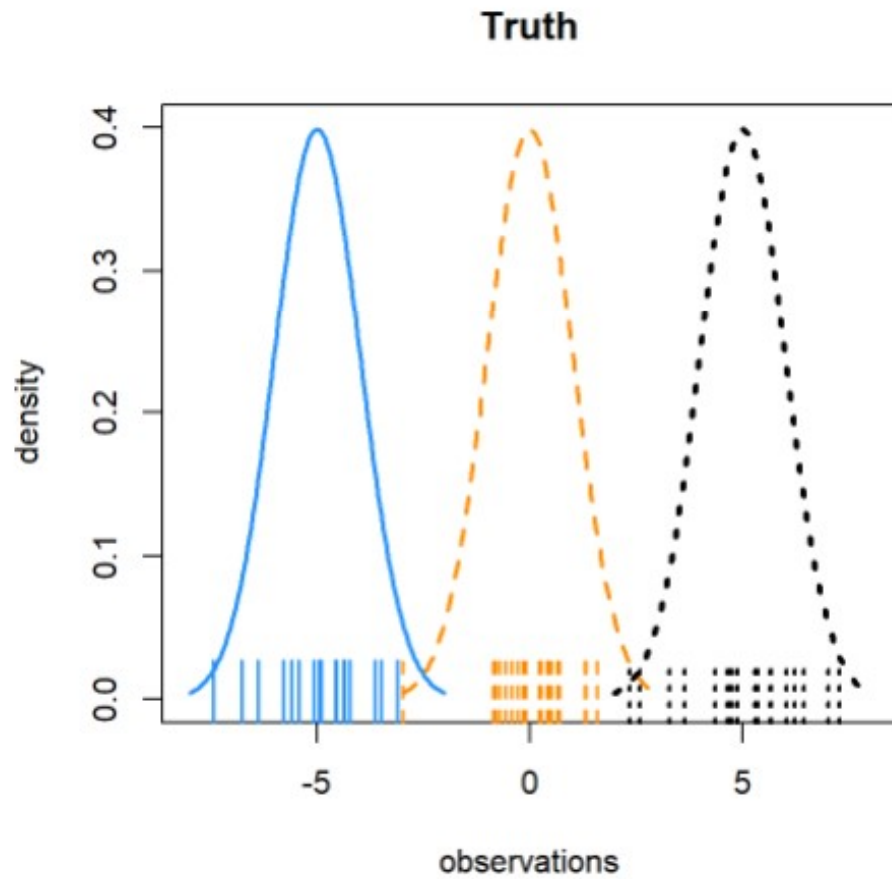
- 變異數分析
 - 組間變異 (Between-group variance)
 - 組內變異 (Within-group variance)
- 組間變異是否大於組內變異
- 補充：變異 = 離平均的距離



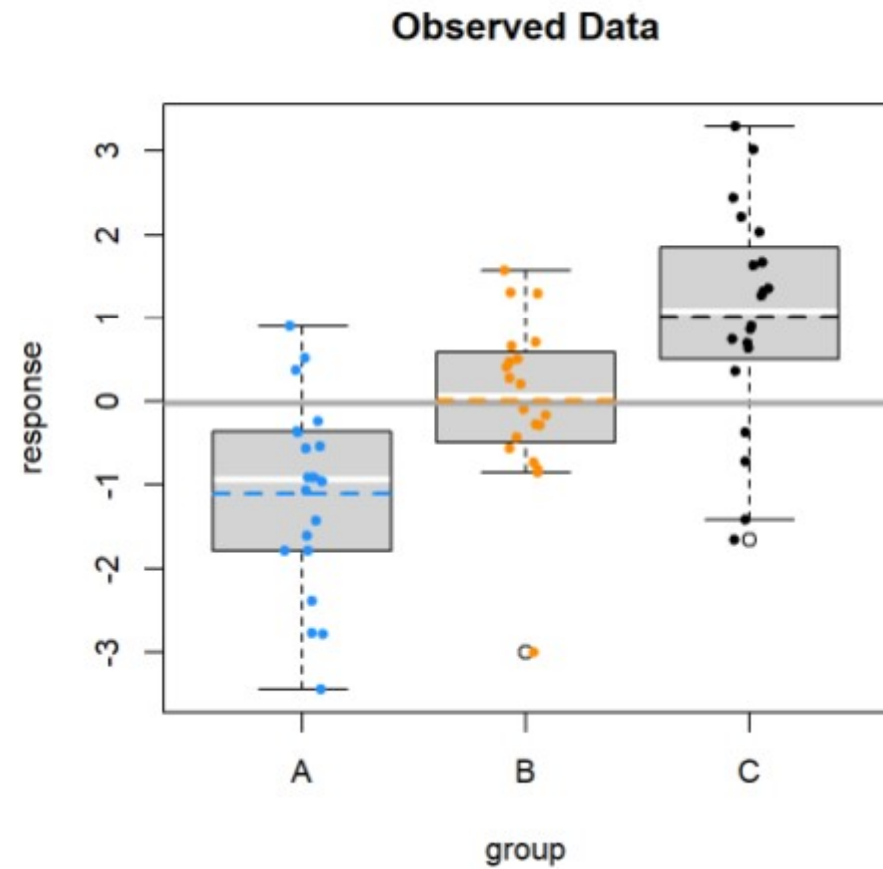
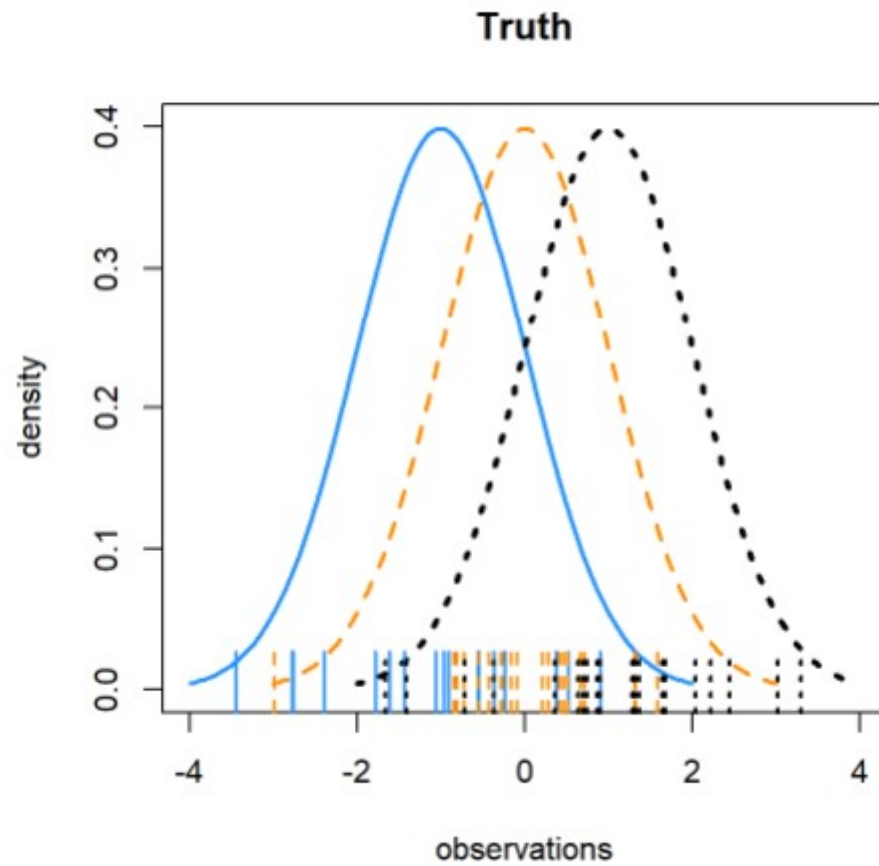
假想的例子

	Data 1			Data 2			Data 3		
	Group 1	Group 2	Group3	Group 1	Group 2	Group3	Group 1	Group 2	Grou
	2	300	20	1	1	1	2.1	5.1	6.1
	3	350	30	2	2	2	2.2	5.2	6.2
	4	400	40	2	5	6	2.2	5.5	6.6
	5	450	50	2	8	10	2.2	5.8	7.0
	6	500	60	3	9	11	2.3	5.9	7.1
均	4	400	40	2	5	6	2.2	5.5	6.6
異數	2	5000	200	0.4	10	16.4	0.004	0.1	0.16

概念 - 組間 vs 組內變異



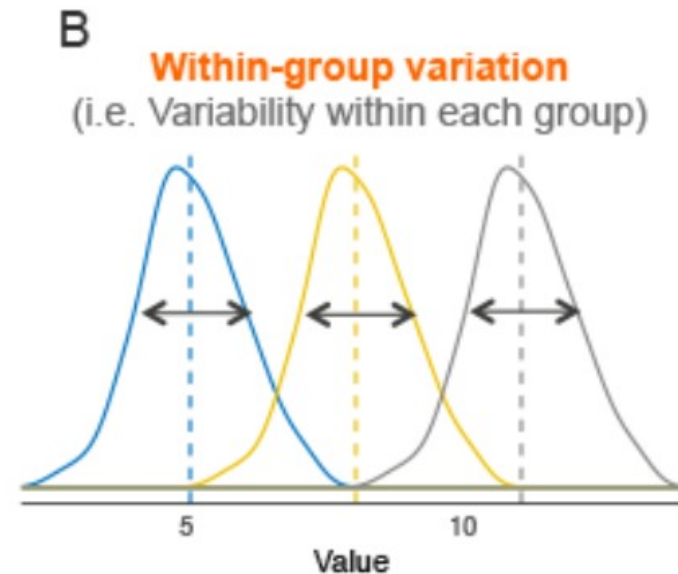
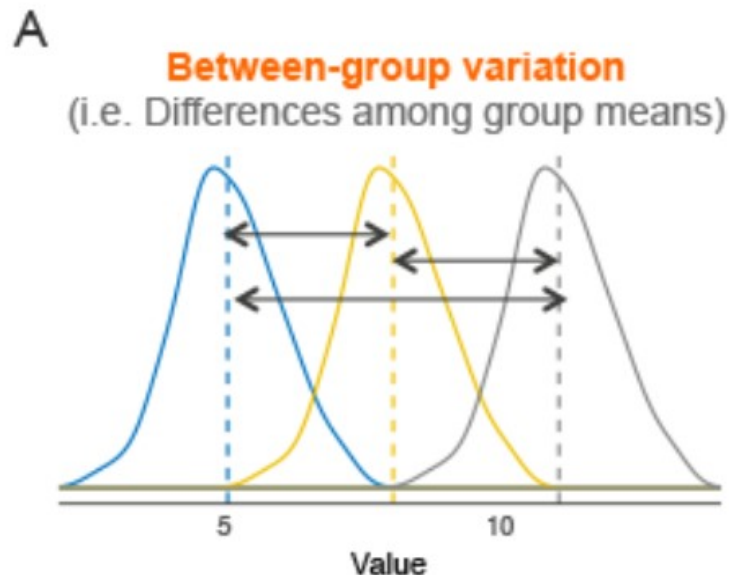
概念 - 組間 vs 組內變異



拆解變異數

$$\begin{array}{ccccc} \text{Variation in } Y & = & \text{Explained Variation} & + & \text{Unexplained Variation} \\ \text{(around its mean)} & & \text{(due to factors)} & & \text{(random error)} \end{array}$$

$$\text{變異數 (Variance)} = SD^2 = \frac{\text{平方和 } \textit{Sum of squares (SS)}}{\text{自由度 } \textit{Degree of freedom (df)}}$$



拆解變異數

Partitioned Sum of Squares

Sum of Squares Total (<i>SST</i>)	=	Sum of Squares between Treatments (<i>SSB</i>)	+	Sum of Squares within Treatments (<i>SSE</i>)
		↑		↑
		Explained by Treatments		Unexplained Random Error

組間變異

組內變異

如何拆解

- ANOVA table

<i>Source of Variation</i>	<i>Sum of Squares</i>	<i>Degrees of Freedom</i>	<i>Mean Square</i>	<i>F Statistic</i>
Treatment (between groups)	$SSB = \sum_{j=1}^c n_j (\bar{y}_j - \bar{y})^2$	$c - 1$	$MSB = \frac{SSB}{c - 1}$	$F = \frac{MSB}{MSE}$
Error (within groups)	$SSE = \sum_{j=1}^c \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_j)^2$	$n - c$	$MSE = \frac{SSE}{n - c}$	
Total	$SST = \sum_{j=1}^c \sum_{i=1}^{n_j} (y_{ij} - \bar{y})^2$	$n - 1$		

F檢定統計量 (F分布)

Between groups
(explained)

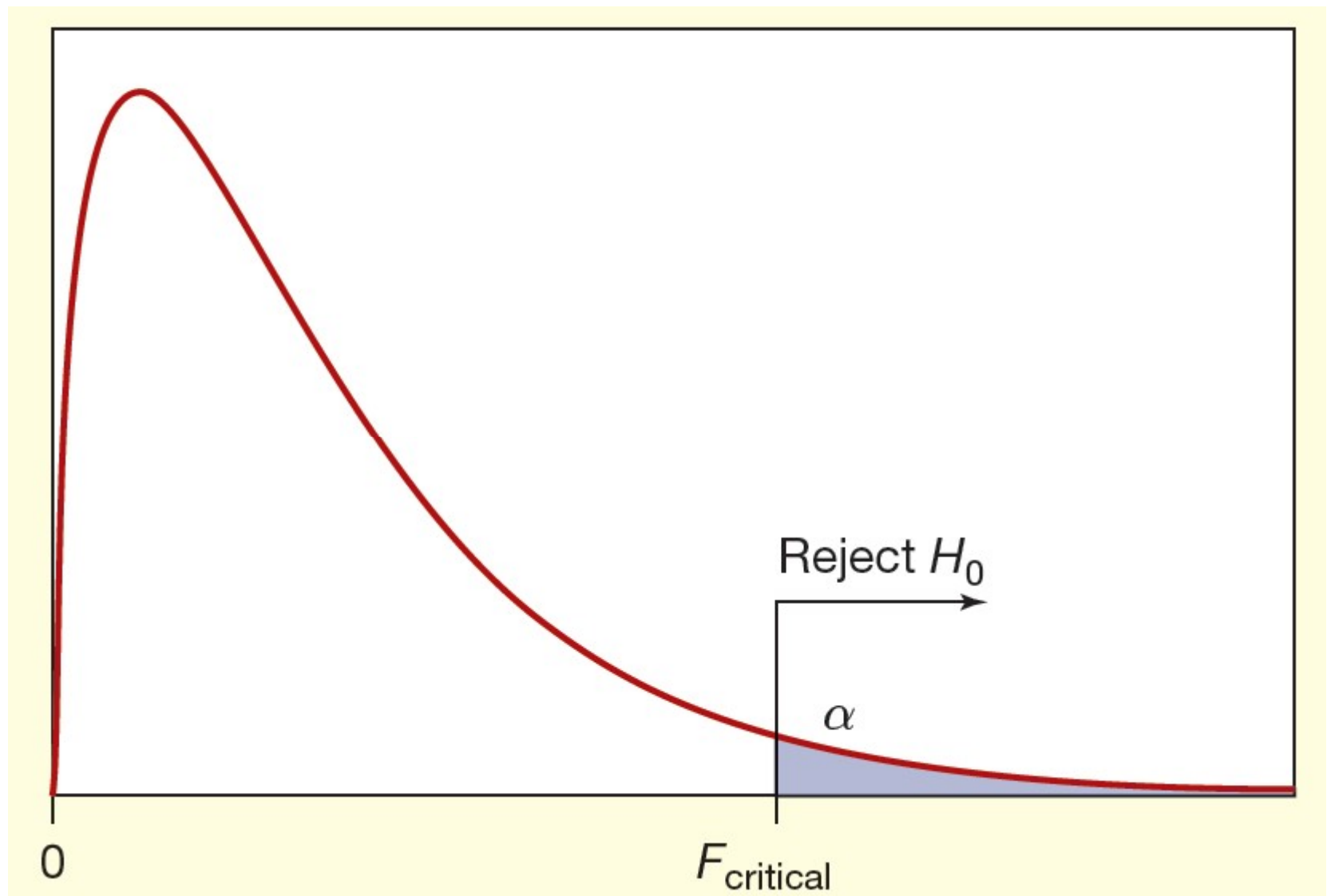
Within groups
(unexplained)

$$F = \frac{MSB}{MSE} = \frac{\left(\frac{SSB}{c - 1} \right)}{\left(\frac{SSE}{n - c} \right)}$$

$df_1 = c - 1$ (numerator)

$df_2 = n - c$ (denominator)

F分布



範例

- 某公司有四個工作站，負責將商品裝箱，運送給零售商，每個工作站每天可以完成200箱以上。
- 請參考data_w11.xlsx檔案
- 請比較四個工作站平均裝箱數是否有差異？

範例解答

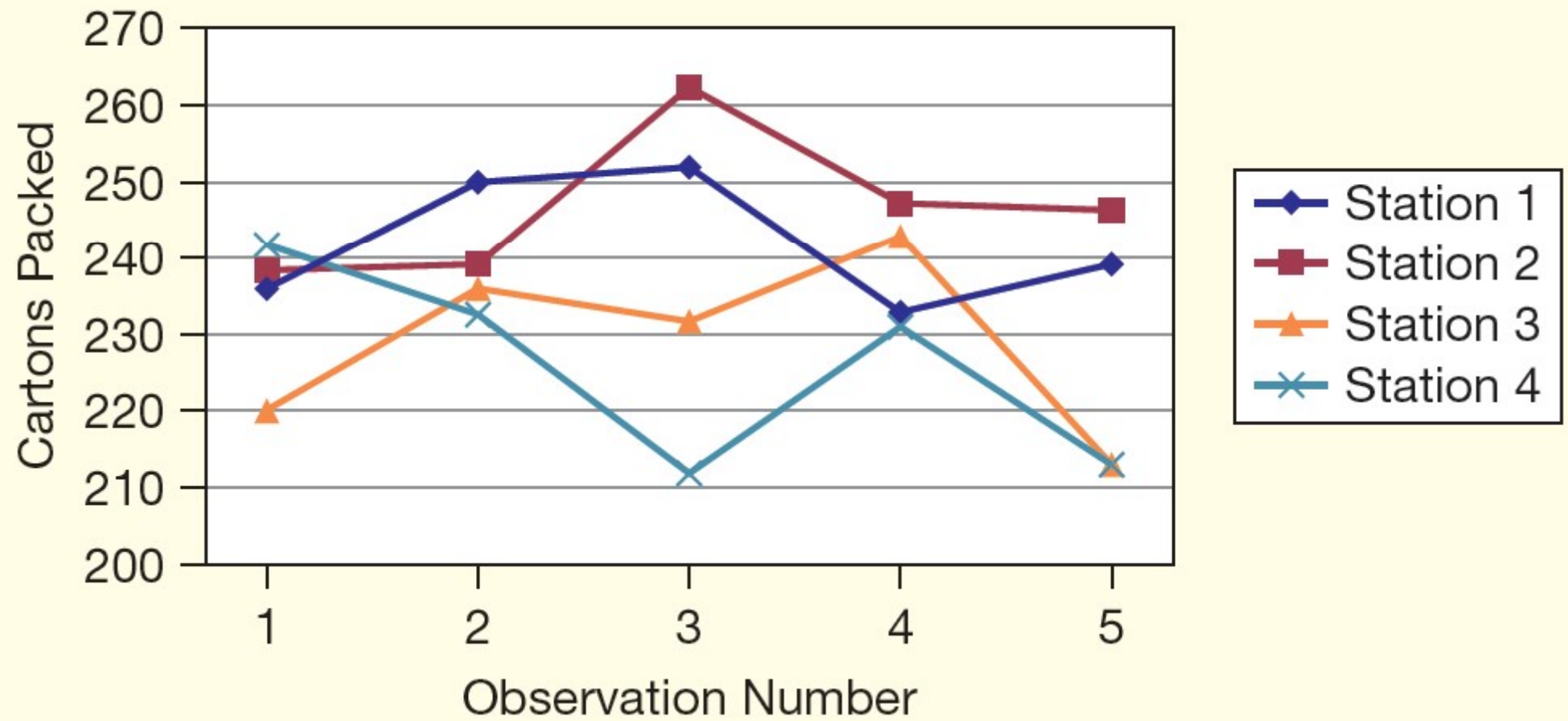
- 假說建立

$H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$ (the means are the same)

H_1 : Not all the means are equal (at least one mean is different)

- 設定決策規則
 - P值法
 - 型I誤差 = 0.05

資料描述



anova: Single Factor

SUMMARY

<i>Groups</i>	<i>Count</i>	<i>Sum</i>	<i>Average</i>	<i>Variance</i>
station 1	5	1210	242.0	72.5
station 2	5	1232	246.4	92.3
station 3	5	1144	228.8	147.7
station 4	5	1130	226.0	166.0

ANOVA

<i>Source of Variation</i>	<i>SS</i>	<i>df</i>	<i>MS</i>	<i>F</i>	<i>P-value</i>	<i>F cr</i>
Between Groups	1479.2	3	493.0667	4.121769	0.024124	3.23887
Within Groups	1914.0	16	119.6250			
Total	3393.2	19				

結論

- P值小於0.05，達到統計顯著，拒絕虛無假說。
- 至少有兩個工作站的平均裝箱數量有差異。

哪幾組有差異？

多重比較 (事後檢定)

$$H_0: \mu_j = \mu_k$$

$$H_1: \mu_j \neq \mu_k$$

- Bonferroni correction

- $\alpha^* = \frac{0.05}{6}$

- Tukey事後檢定

$$T_{\text{calc}} = \frac{|\bar{y}_j - \bar{y}_k|}{\sqrt{MSE \left[\frac{1}{n_j} + \frac{1}{n_k} \right]}}$$

變異數分析的統計假設

- 獨立樣本 (Independence)
- 常態分布 (Normal distribution)
- 變異數相等 (Equal variance; homoscedasticity)