

Assumptions

for ANOVA

- ① Pop. variances are equal
- ② " normal
- ③ samples are independent

<http://profs.degroote.mcmaster.ca/ads/parlar/courses/q600/ChapterComments/documents/ANOVA-Calcs-Scan-Colour.pdf>

← print it

<http://profs.degroote.mcmaster.ca/ads/parlar/courses/q600/ChapterComments/documents/Fertilizer.xls>

One factor ANOVA					
	Mean	n	Std. Dev		
	3.0	6	1.90	Group 1	
	7.0	6	1.26	Group 2	
	5.0	6	1.41	Group 3	
	5.0	18	2.22	Total	
ANOVA table					
Source	SS	df	MS	F	p-value
Treatment	48.00	2	24.000	10.00	.0017
Error	36.00	15	2.400		
Total	84.00	17			
Post hoc analysis					
p-values for pairwise t-tests					
		Group 1	Group 3	Group 2	
		3.0	5.0	7.0	
Group 1	3.0				
Group 3	5.0	.0410			
Group 2	7.0	.0004	.0410		

Pasted from <file:///C:/DOCUME~1/parlar/LOCALS~1/Temp/Fertilizer.xls>

c) Conf. Int. for  $\mu_i - \mu_j$

Result:  $100(1-\alpha)\%$  CI

$$\left[ (\bar{x}_i - \bar{x}_j) \pm t_{\alpha/2, n-p} \sqrt{MSE \left( \frac{1}{n_i} + \frac{1}{n_j} \right)} \right], \quad df = n - p$$

L-M (n-p)

$$L-H (\mu_1 - \mu_2)$$

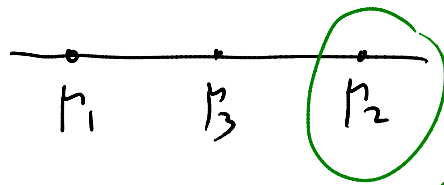
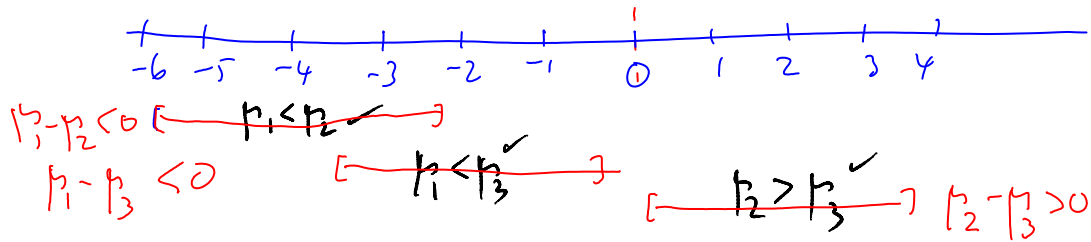
$$[(3-7) \mp 2.13 \sqrt{2.4(\frac{1}{6} + \frac{1}{6})}] = [-5.90, -2.09]$$

$$L-H (\mu_1 - \mu_3)$$

$$[(3-5) \mp \dots] = [-3.90, -0.09]$$

$$M-H (\mu_2 - \mu_3)$$

$$[(7-5) \mp \dots] = [0.09, 3.90]$$



M is best!

MegaStat

$$H_0: \mu_i = \mu_j$$

$$H_a: \mu_i \neq \mu_j$$

Post hoc analysis				
p-values for pairwise t-tests				
	L			
	Group 1	Group 3	Group 2	
		3.0	5.0	7.0
Group 1	3.0			
Group 3	5.0	.0410		
M Group 2	7.0	.0004	.0410	

$\mu_1 \neq \mu_2$ , etc.

Pasted from <file:///C:/DOCUME~1/parla/LOCALS~1/Temp/Fertilizer.xls>

Ch. 11 Correlation Coefficient & Simple linear regression

Two variables & how they relate

## a) Covariance & correlation coefficient

Ex, Height vs. handspan

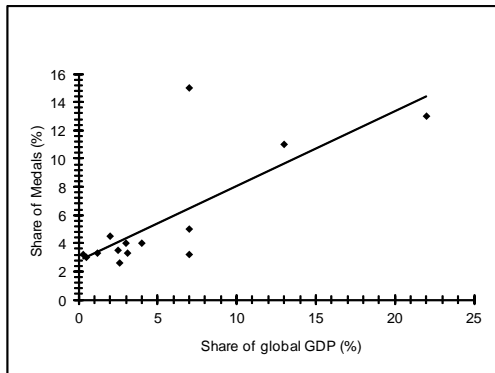
<http://profs.degroote.mcmaster.ca/ads/parlar/courses/q600/ChapterComments/documents/Q600-2013-Scanned-Height-Gender-Handspan.pdf>

<http://profs.degroote.mcmaster.ca/ads/parlar/courses/q600/ChapterComments/documents/Q600-2013-Height-Gender-Handspan-Regression.xlsx>

Ex. Olympic medals vs GDP

2008-08-08 8:08 pm 14

[http://profs.degroote.mcmaster.ca/ads/parlar/courses/q600/ChapterComments/documents/MedalsAndEconomy\\_000.xls](http://profs.degroote.mcmaster.ca/ads/parlar/courses/q600/ChapterComments/documents/MedalsAndEconomy_000.xls)



Country	Share of global GDP (%)	Share of Medals (%)
USA	22	13
China	13	11
Russia	7	15
Great Britain	7	5
Australia	2	4.5
Germany	4	4
France	3	4
Korea	2.5	3.5
Italy	3.1	3.3
Ukraine	1.2	3.3
Japan	7	3.2
Cuba	0.3	3.2
Belarus	0.5	3
Canada	2.6	2.6

Pasted from <file:///C:/DOCUME~1/parlar/LOCALS~1/Temp/MedalsAndEconomy\_000.xls>

Covariance

Variance (single var)

$$S_x^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

Co-variance (two vars)

$$S_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

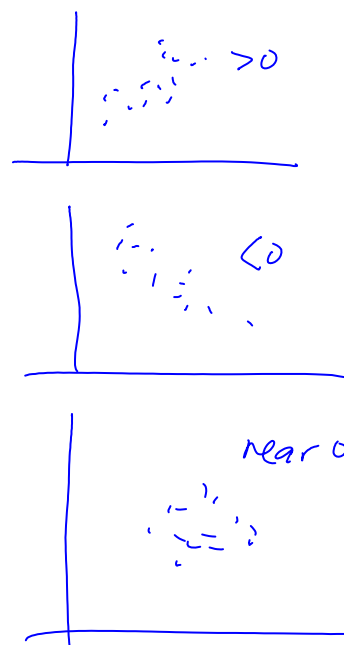
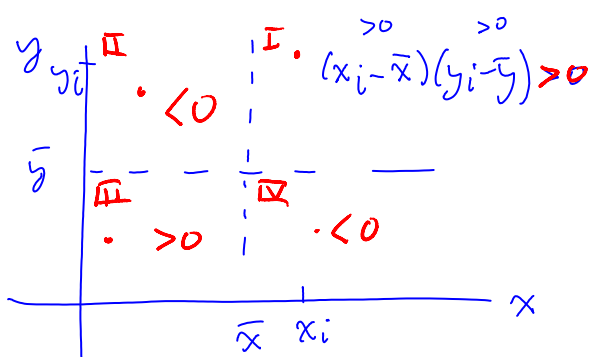
	Share of global GDP (%)	Share of Medals (%)
count	14	14
mean	$\bar{x}$ 5.371	$\bar{y}$ 5.614
sample variance	34.575	17.018
sample standard deviation	5.880	4.125
minimum	0.3	2.6
maximum	22	15
range	21.7	12.4

Pasted from <file:///C:/DOCUME~1/paran/LOCALS~1/Temp/MedalsAndEconomy\_000.xls>

$n=14$

$x_i$	$y_i$	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$
22	13	16.62	7.38	122.81
'	'	'	'	'
'	'	'	'	'
2.6	2.6	-2.77	-3.01	8.35
				<hr/> 238.36

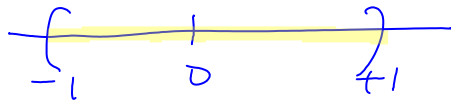
$$S_{xy} = \frac{238.36}{13} = 18.33 > 0$$



Better measure: Correlation Coefficient ( $r$ )

$$r = \frac{S_{xy}}{S_x S_y}$$

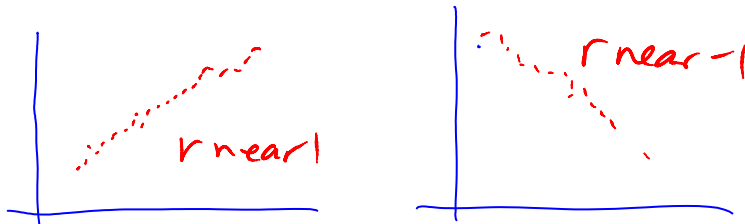
$$S_x = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2}, \quad S_y = \sqrt{\frac{1}{n-1} \sum (y_i - \bar{y})^2}$$



Ex. Medals

$$S_{xy} = 18.33, \quad S_x = 5.88, \quad S_y = 4.12$$

$$r = \frac{18.33}{(5.88)(4.12)} = 0.75 \quad \text{somewhat strong}$$



Correlation doesn't imply causation

## b) Simple Linear Regression

x ind't var  
y dep't var



Ex. Naive model (no uncertainty)

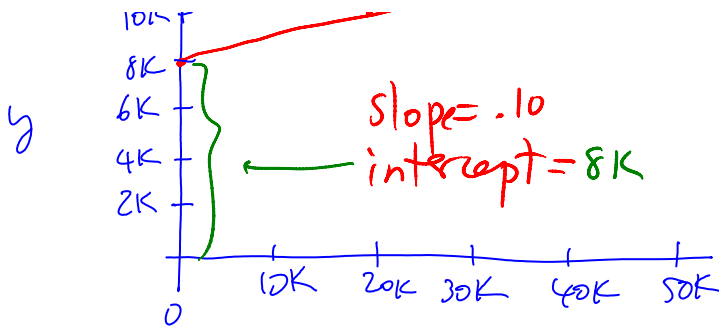
x: <sup>gross</sup> sales at local Esso (monthly)

y: payments to parent comp. (")

Agreement: \$8,000/month

+ 10% of gross

$$y = 8000 + 0.10x \quad \text{Linear model}$$



$$x=0, y=8$$

$$x=20, y=8 + .1(20)$$

$$= 8 + 2 = 10$$

$$y = \beta_0 + \beta_1 x$$

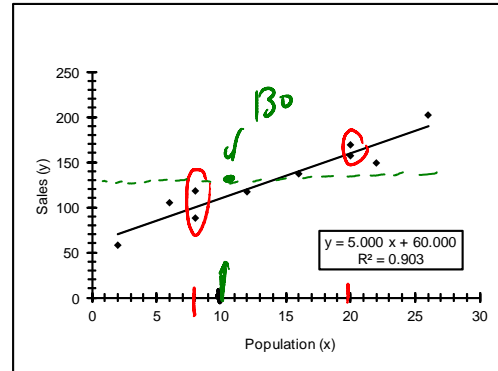
$\downarrow$                        $\uparrow$   
 Intercept              slope

Ex. Statistical model - Harvey's restaurant

Student Population (x)	Monthly Sales (y)
2	58
6	105
8	88
8	118
12	117
16	137
20	157
20	169
22	149
26	202

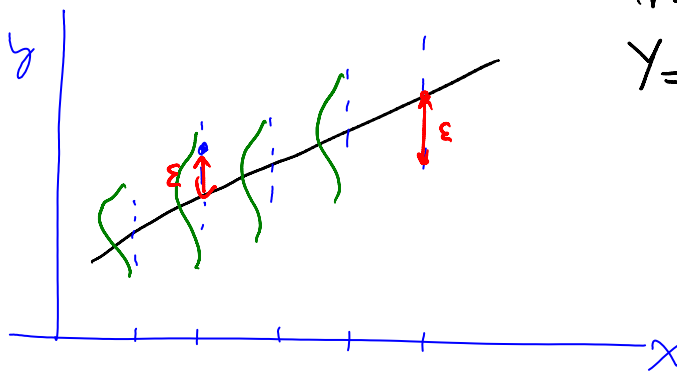
Pasted from <file:///C:/DOCUME~1/nparal/LOCALS~1/Temp/Harveys-1.xls>

$$\bar{y} = 130$$



Statistical model needed

Pop'n of 300+ restaurants

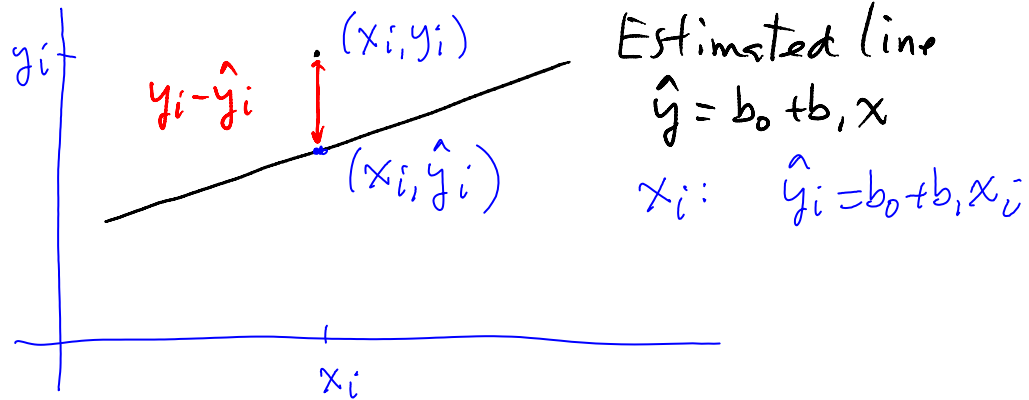
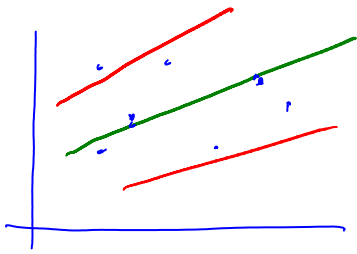


True model

$$Y = \beta_0 + \beta_1 X + \epsilon$$

$\downarrow$                $\downarrow$   
 $b_0$                $b_1$

c) Best method to find regression line

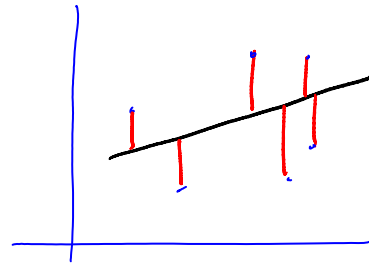


Consider  $(x_i, y_i)$  : Actual  $y_i$   
 Estimate  $\hat{y}_i$

Residual (error).  $e_i = y_i - \hat{y}_i$   
 $= y_i - (b_0 + b_1 x_i)$

Minimize

$$SSE = \sum_{i=1}^n [y_i - (b_0 + b_1 x_i)]^2$$



Solution for finding  $b_0$  &  $b_1$  to min. SSE

$$\textcircled{1} \quad b_1 = \frac{SS_{xy}}{SS_{xx}}, \quad SS_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$$

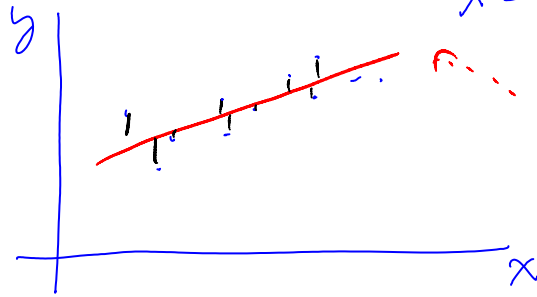
$$= \sum x_i y_i - \frac{1}{n} (\sum x_i)(\sum y_i)$$

$$SS_{xx} = \sum (x_i - \bar{x})^2$$

$$= \sum x_i^2 - \frac{1}{n} (\sum x_i)^2$$

$$\textcircled{2} \quad b_0 = \bar{y} - b_1 \bar{x}, \quad \bar{y} = \frac{1}{n} \sum y_i$$

$$\bar{x} = \frac{1}{n} \sum x_i$$



$$\hat{y} = b_0 + b_1 x$$

