

STA 302 / 1001

Note Title

10/4/2011

Reminder Anatomy in RW 107/109 today 11:00 - 13:00  
in RW 109 Friday 13:00 - 15:00

Test date

Thurs Oct 20 10:10 - 11:40 26

~~Thurs Oct 25 10:10 - 11:40~~

$X, Y$  both random variables

BIVARIATE NORMAL Distribution

$X, Y$  are jointly normally distributed if their joint density function is

$$f(x, y) = \frac{1}{\pi \sigma_x \sigma_y \sqrt{1-\rho^2}} \exp \left\{ -\frac{1}{2(1-\rho^2)} \left[ \left( \frac{x-\mu_x}{\sigma_x} \right)^2 - 2\rho \left( \frac{x-\mu_x}{\sigma_x} \right) \left( \frac{y-\mu_y}{\sigma_y} \right) + \left( \frac{y-\mu_y}{\sigma_y} \right)^2 \right] \right\}$$

$$-\infty < x < \infty, \quad -\infty < y < \infty$$

Can show (STA 257 exercises) marginal distributions:  
 $X \sim N(\mu_X, \sigma_X^2), \quad Y \sim N(\mu_Y, \sigma_Y^2)$

$\rho$  "rho"  
is the CORRELATION COEFFICIENT between  $X$  and  $Y$

i.e. 
$$\rho = \frac{E[(X - E(X))(Y - E(Y))]}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

It can be shown (STA 257 exercise) that the conditional distribution of  $Y$  given  $X = x$

$$Y|X=x \sim N\left(\mu_Y + \rho \frac{\sigma_Y}{\sigma_X}(x - \mu_X), (1 - \rho^2) \sigma_Y^2\right)$$

- Conditional distribution of  $Y|X=x$
- $E(Y|X=x)$  is of the form  $\beta_0 + \beta_1 x$

- Variance is constant over  $x$

Estimate mean of  $Y|X=x$  by  $\hat{\mu}_{Y|X=x}$  estimating

$$\hat{\mu}_{Y|X=x} = \frac{\mu_Y + \rho \frac{\sigma_Y}{\sigma_X}(x - \mu_X)}{\sqrt{\frac{1}{n-1} \sum (y_i - \bar{y})^2}} = \text{s.d.}(y)$$

$$\hat{\mu}_{Y|X=x} = \frac{\mu_Y + \rho \frac{\sigma_Y}{\sigma_X}(x - \mu_X)}{\sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2}} = \text{s.d.}(x)$$

$$\text{and } \rho_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}{\text{s.d.}(x) \text{s.d.}(y)}$$

you get (check this)

$$\text{Est. of } E(Y | X=x) = b_0 + b_1 x$$

$$\text{where } b_0 = \bar{y} - b_1 \bar{x}, \quad b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

i.e. the estimator gives the least squares estimate of the slope and intercept

# CORRELATION

Suppose you have 2 r.v.'s  $X$  and  $Y$   
Interested in relationship between  $X$  &  $Y$ , believe it's  
linear

- If there is a clear choice for  $Y$  being the response and  $X$  being the explanatory variable, can carry out a regression
- An option when  $X$  and  $Y$  don't have a clear choice for dependent/independent variables, can summarise the strength of the linear relation by correlation

• Correlation is a symmetric measure

$n$  observations  $(x_i, y_i)$  of r.v.'s  $X$  and  $Y$

Correlation (Pearson's Product Moment Correlation)

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

When  $X, Y$  have a bivariate normal distribution,  
 $r$  is the maximum likelihood estimate of  $\rho$

Facts about  $r$ :

- measure of degree of linear association between  $X$  &  $Y$

- it is dimension free

- always between  $+1$  and  $-1$ , inclusive

•  $r = +1$  means  $(x_i, y_i)$  fall exactly on a straight line with positive relationship

•  $r = -1$  means negative relationship

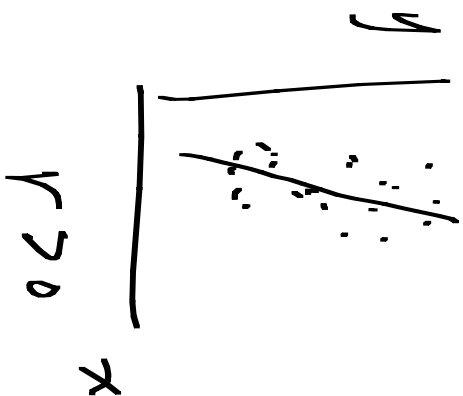
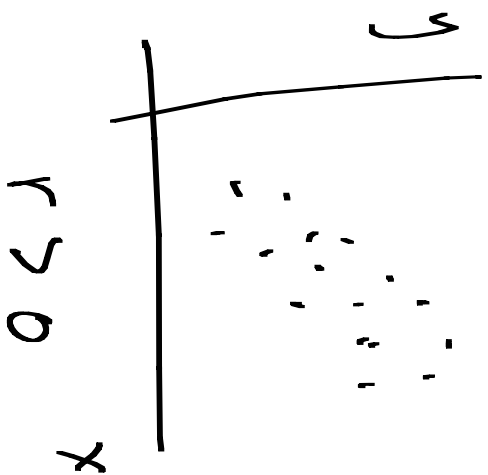
•  $r = 0$  means no linear relationship

SAs proc corr ;

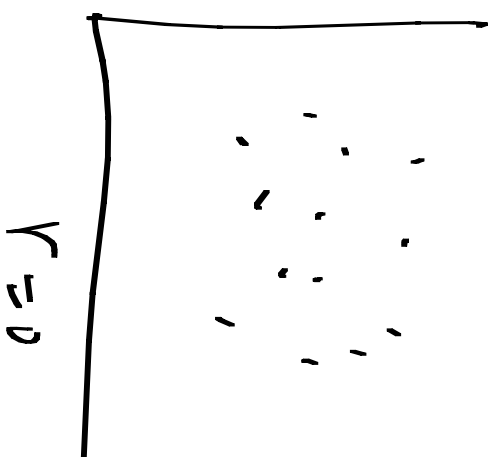
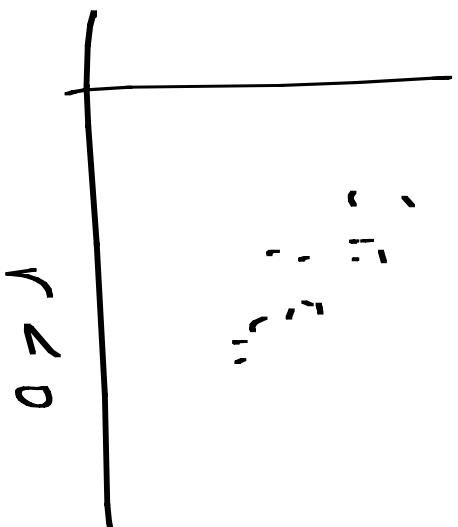
gives correlations and the  $p$ -value  
for test:  $H_0: \rho = 0$  versus  $H_a: \rho \neq 0$

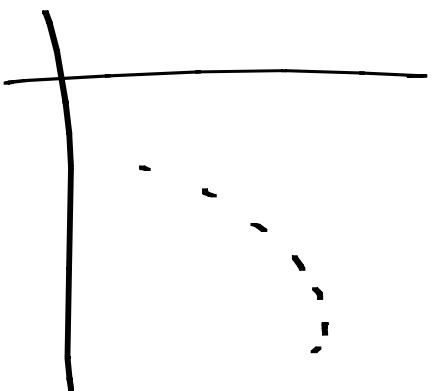
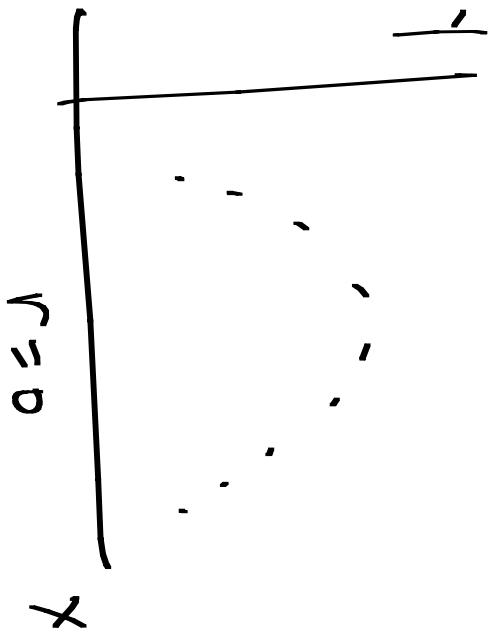


Assuming  $X, Y$  bivariate normal



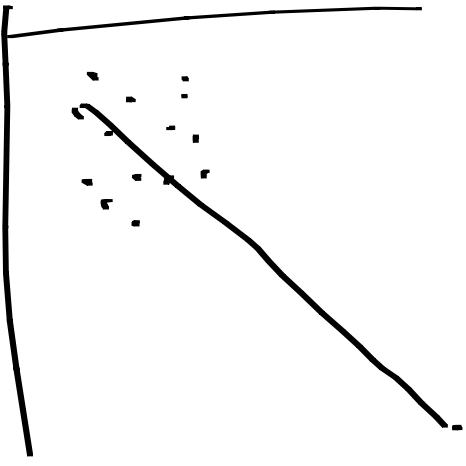
These 2 have similar values of  $r$ .





$0.2 < r < 1$

Not a linear relationship  
so  $r$  not appropriate



$r > 0$   
useless

Relationship between  $R^2$  and  $r$

---

$$R^2 = r^2$$

$$\begin{aligned} R^2 &= \frac{SS_{\text{Reg}}}{SS_T} = \frac{b_1^2 \sum (x_i - \bar{x})^2}{\sum (y_i - \bar{y})^2} \\ &= \left[ \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \right]^2 \frac{\sum (x_i - \bar{x})^2}{\sum (y_i - \bar{y})^2} \\ &= r^2 \end{aligned}$$

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

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2}} \cdot \frac{\sqrt{\sum (y_i - \bar{y})^2}}{\sqrt{\sum (x_i - \bar{x})^2}} / \sqrt{n-1}$$

$$= r \frac{s_y}{s_x}$$

If  $r > 0 \Rightarrow b_1 > 0$ ,  $r < 0 \Rightarrow b_1 < 0$ ,  
 $r > 0 \Rightarrow b_1 > 0$ ,  
 $r < 0 \Rightarrow b_1 < 0$

Units of  $b_1$  are " $^n y$  per unit  $x$ "

Example

CFCs pre-MP  $b_1 = 9.71152$

indicates for every increase in time ( $x$ ) by 1 year, the estimate of the mean of CFC concentration increases by 9.71152 parts per trillion

Interpretation  
of  $b_1$

Exercise:

Known:  $\sum \hat{e}_i x_i = 0$  and  $\sum \hat{e}_i y_i = 0$

Why does this imply  $\hat{e}_i, x_i$  are uncorrelated?  
and  $\hat{e}_i, y_i$  are uncorrelated?