

## Checking Adequacy of Model

(21)

- Assuming regression model is adequate  
we can use it to make inferences,  
see previous page. Before doing so,  
we need to check if the assumptions  
hold.

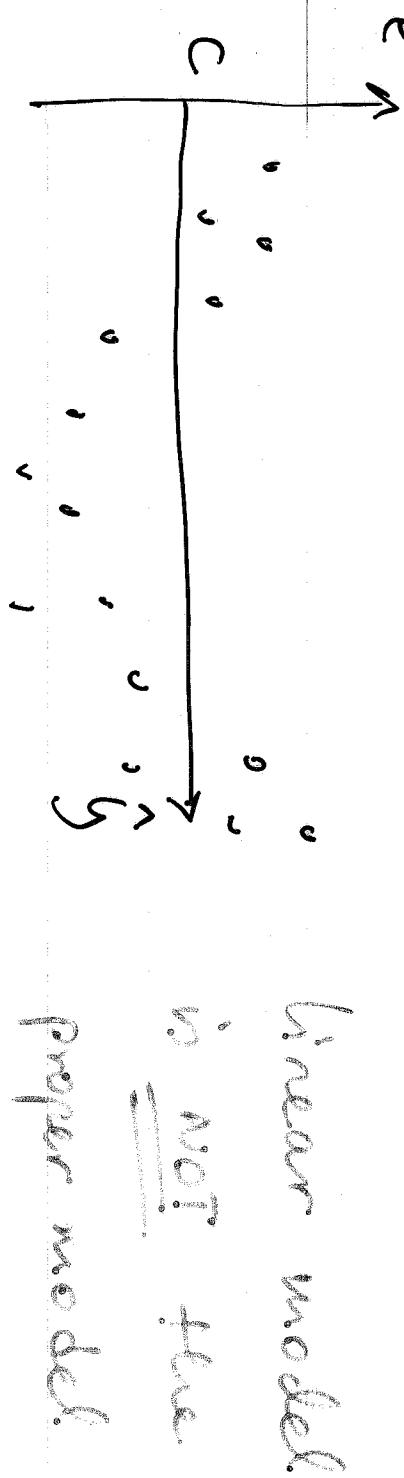
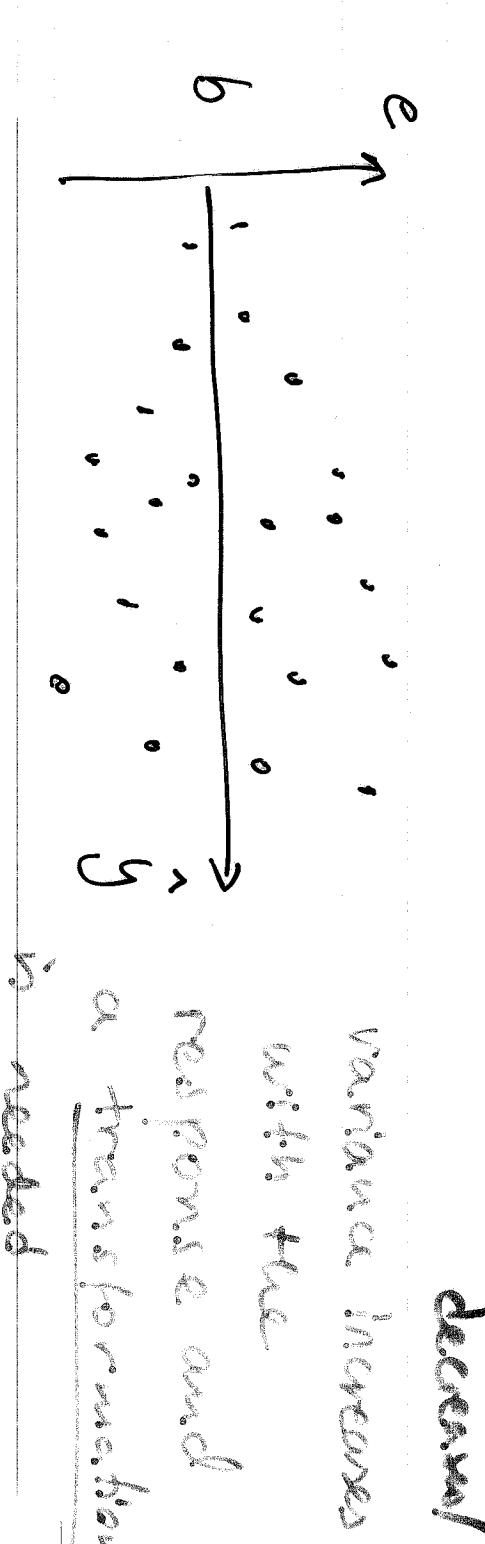
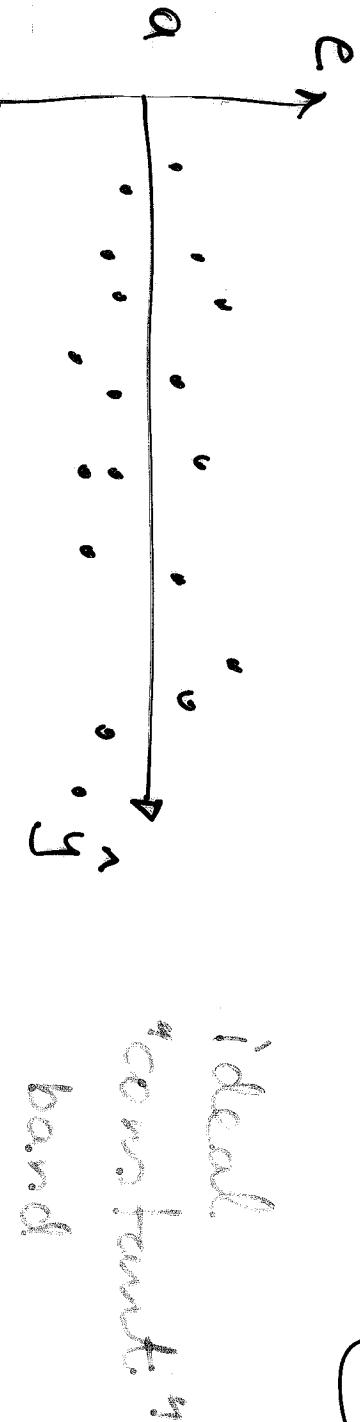
- Are  $\epsilon \sim N(0, \sigma^2)$  ?

First we compute the residuals,  
after  $a, b$  the estimates of  $\alpha$  and  
 $\beta$  have been estimated.

Then

$$e_i = y_i - \hat{y}_i = y_i - (a + b x_i)$$

- Normal plot the residuals  $e_i$  to  
check if the normality assu-  
mption holds
- Plot residuals vs  $\hat{y}_i$  (predicted  
values)



- Plot residuals vs time (or order when this applies) to detect any trends over time.

When variance of residuals is not constant, like in Figure b, we need (2.3)

to do a transformation in order to stabilize the variance.

It is usual practise then to substitute in the regression

$$Y = \alpha + \beta X + \varepsilon \quad \text{or}$$

$$\hat{Y} = \alpha + b \bar{X}$$

the independent variable  $\bar{X}$  by a transformed one. Usually  $\bar{X} = \log X$  or  $\bar{X}' = \sqrt{X}$  or  $\bar{X}' = X^k$  for some  $k$ , are the most common ones.

Be careful,  $\log X$ ,  $\sqrt{X}$  can be used when  $X$  takes only positive values!

## Correlation

(24)

Up to now, we have assumed that the independent variable  $\underline{x}$  was assumed to be known without error. In some cases though both  $x$  and  $y$  have to be thought of as values of random vars  $X, Y$ .

Example Study relationship between input and output of a wastewater treatment plant. / tensile strength and hardness of aluminum ...

To use Correlation.  $(x, y)$   $\rightarrow$   $\underline{x}$  and  $\underline{y}$  are values of a pair of random variables  $X$  and  $Y$  having joint density  $f_{X,Y}(x,y)$ .

Scatterplot provides visual impression of the relation between  $x$  and  $y$ .

Characteristics of the scatter plot to a straight line can be expressed in terms of correlation coefficient (numerically).

Sample correlation coefficient,  $r$ , is defined

$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \cdot \left( \frac{y_i - \bar{y}}{s_y} \right) = \frac{s_{xy}}{\sqrt{s_{xx} \cdot s_{yy}}}$$

$$\text{Where } s_x^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1} = s_{xx}$$

$$s_y^2 = \sum_{i=1}^n \frac{(y_i - \bar{y})^2}{n-1} = s_{yy}$$

- When most of the pairs of obs are s.t. either both  $(x_i, y_i)$  are above their sample means or both below their sample means the products tend to be large and positive so  $r > 0$ .

- When one value of the pair tends to be large when the other is small and vice versa,  $r < 0$ .

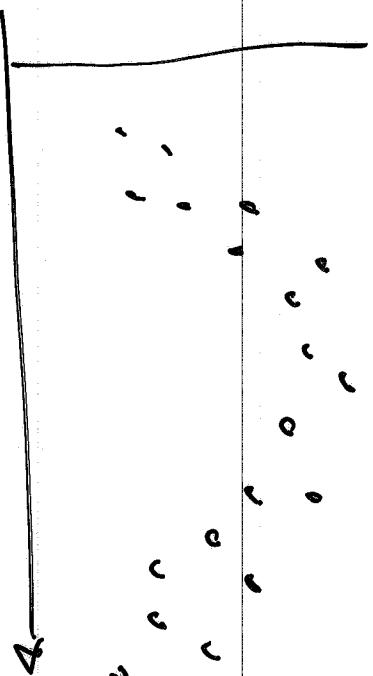
- $-1 \leq r \leq 1$ .

1.) Magnitude of  $r$  describes the strength of a linear relation and sign of  $r$  indicates the direction

- $r = 1$ , all  $(x_i, y_i)$  lie exactly on a straight line having positive slope.
- $r > 0$  - pattern of scatter plot runs from lower left to upper right
- $r < 0$  - pattern of scatter plot runs from upper left to lower right.
- $r = -1$  all pairs  $(x_i, y_i)$  lie on a straight line having a negative slope.

- A value of  $r$  near +1 or -1 describes strong linear relation.

2) A value of  $r$  close to 0  $\Rightarrow$  linear association is weak but there can be still strong association along for example, a curve



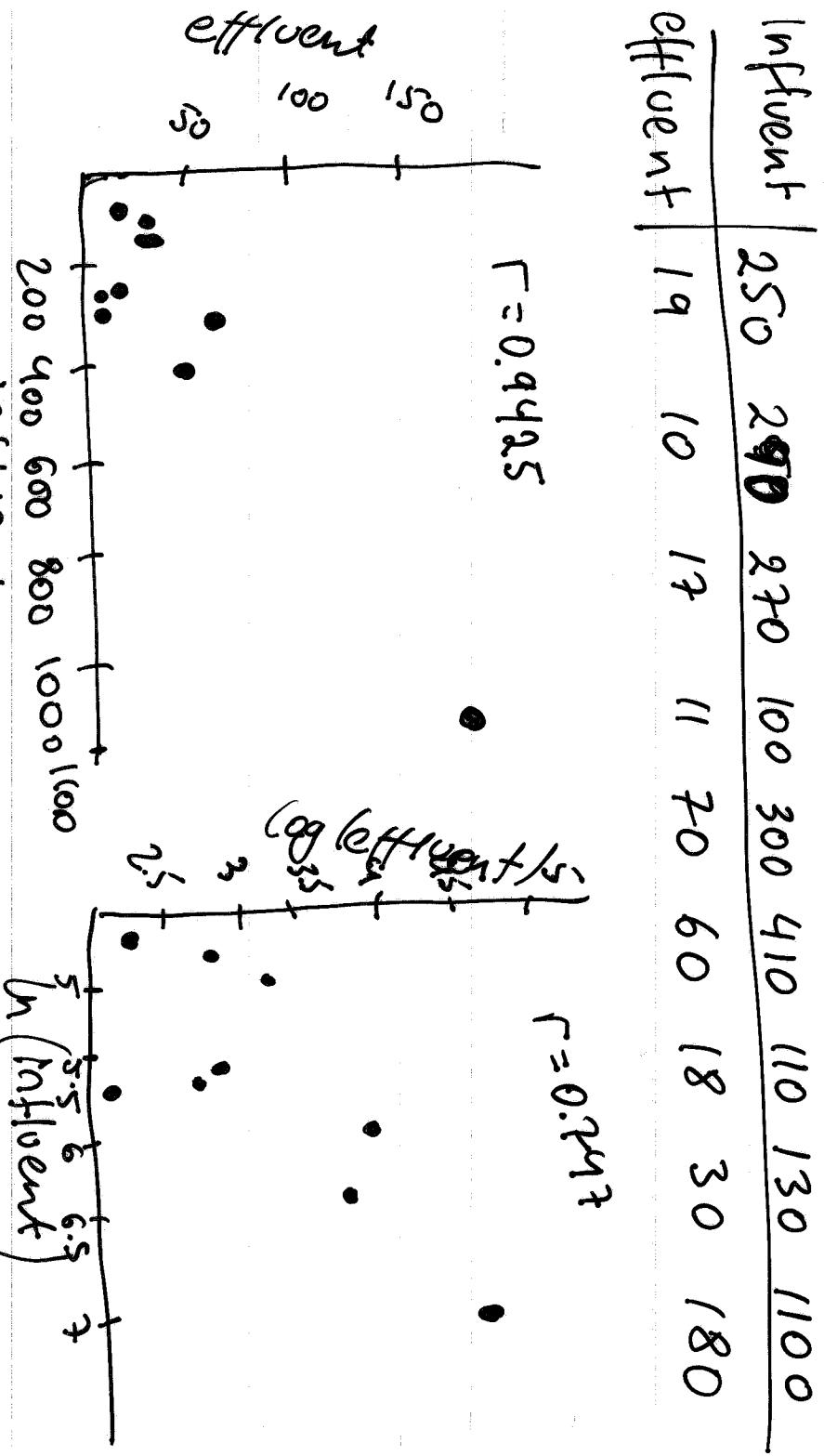
# Exploring and interpreting correlation

(28)

correlation

## Example

Heavy metals can inhibit the biological treatment of waste in municipal treatment plants. Monthly measurements were made at a state-of-the-art treatment plant of the amount of chromium ( $\mu\text{g/l}$ ) in both the influent and effluent.



Notice that  $r$  is not really appropriate for the original data since the one large obs in the upper-right hand corner has too much influence. If the point (160, 180) is dropped,  $r$  drops to 0.578.

## Correlation and Regression

There are two important relationships between  $r$  and the least square fit of a straight line.

$$r = \frac{S_{xy}}{\sqrt{S_{xx} \cdot S_{yy}}} = \frac{\sqrt{S_{xx}}}{\sqrt{S_{yy}}} \frac{S_{xy}}{S_{xx}} = \frac{\sqrt{S_{xx}}}{\sqrt{S_{yy}}} \cdot b \Rightarrow$$

$r$  and  $b$  have the same sign.

total variation in y

$$\sum_{i=1}^n (y_i - \bar{y})^2 = S_{xy}^2 / S_{xx} + \left( S_{yy} - \frac{S_{xy}^2}{S_{xx}} \right) \quad (30)$$



$$S_{yy} = \frac{S_{xy}^2}{S_{xx}} + \left( S_{yy} - \frac{S_{xy}^2}{S_{xx}} \right)$$



total variability

of y

↑  
Variability explained by linear relation

↑  
Residual or unexplained variability.

For the straight line to provide a good fit to the data, the sum of squares due to regression,  $S_{xy}^2 / S_{xx}$ , should be a maximum.

Proportion of the total sum of squares due to regression.

Sum of squares due to regression

$$= \frac{S_{xy}^2}{S_{xx}}$$

Total sum of squares of y

$$= \frac{S_{xy}^2}{S_{xx} S_{yy}} = r^2$$

Correlation coefficient