

# STK4900/9900 - Lecture 2

## Program

1. Comparing two or more groups
  2. One-way analysis of variance (ANOVA)
  3. Covariance and correlation
  4. Simple linear regression
- Section 2.4
  - Sections 3.1.4, 3.2 (not 3.2.2), 3.3
  - Supplementary material on covariance, correlation and one-way ANOVA

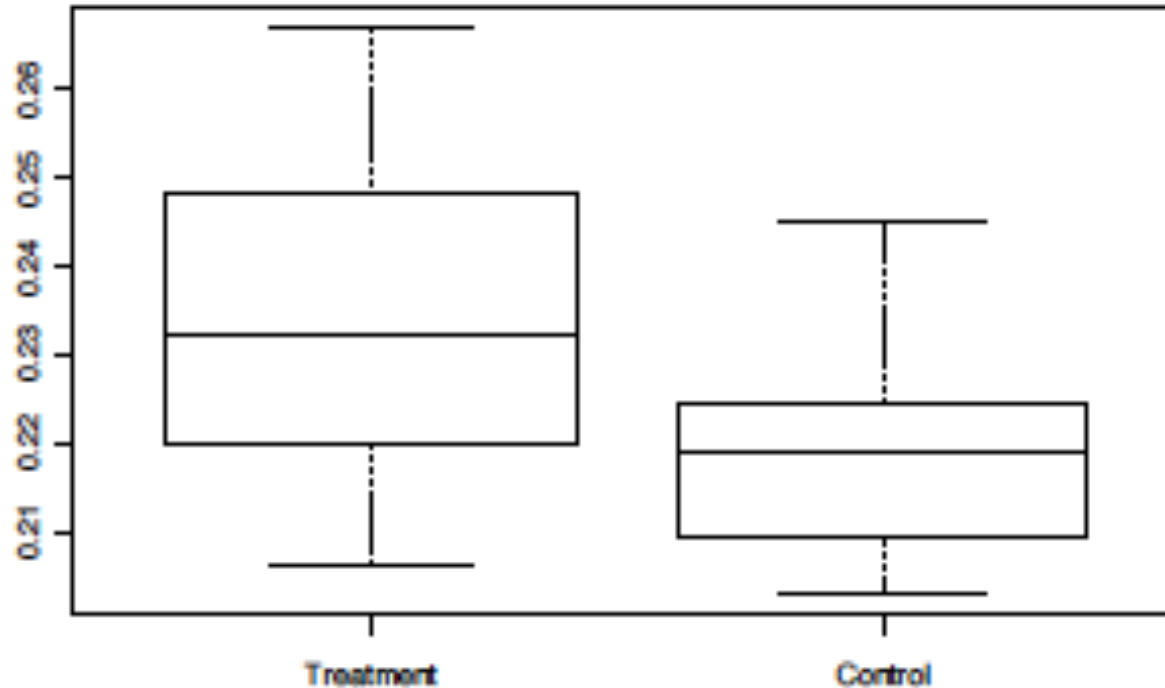
# Comparing two groups

In Lecture 1 we considered an example where we measured bone mineral density (in g/cm<sup>2</sup>) for rats given isoflavone and for rats in a control group:

1) Control ( $n_1 = 15$ )				
0.228	0.207	0.234	0.220	0.217
0.228	0.209	0.221	0.204	0.220
0.203	0.219	0.218	0.245	0.210
2) Isoflavone ( $n_2 = 15$ )				
0.250	0.237	0.217	0.206	0.247
0.228	0.245	0.232	0.267	0.261
0.221	0.219	0.232	0.209	0.255

Question: Does isoflavone have an effect on bone mineral density?

A boxplot gives a graphical comparison of the two groups:



We would like to determine a confidence interval for the treatment effect and test if the difference is statistically significant (cf. next slide)

## R-commands:

```
cont=c(0.228, 0.207, 0.234, 0.220, 0.217, 0.228, 0.209, 0.221, 0.204, 0.220,  
        0.203, 0.219, 0.218, 0.245, 0.210)  
treat=c(0.250, 0.237, 0.217, 0.206, 0.247, 0.228, 0.245, 0.232, 0.267, 0.261,  
        0.221, 0.219, 0.232, 0.209, 0.255)  
boxplot(treat, cont, names=c("Treatment", "Control"))  
t.test(treat, cont, var.equal=T)
```

## R-output (slightly edited)

Two Sample t-test

data: treat and cont

$t = 2.844$ ,  $df = 28$ ,  $p\text{-value} = 0.0082$

alternative hypothesis: true difference in means is not equal to 0

95 percent confidence interval:

0.0045 0.0279

sample estimates:

mean of x	mean of y
-----------	-----------

0.2351	0.2189
--------	--------

Suppose that the data for the two groups are random samples from  $N(\mu_1, \sigma^2)$  and  $N(\mu_2, \sigma^2)$ , respectively

Consider testing the null hypothesis  $H_0 : \mu_1 = \mu_2$  versus the alternative  $H_A : \mu_1 \neq \mu_2$

Test statistic:

$$t = \frac{\bar{x}_2 - \bar{x}_1}{se(\bar{x}_2 - \bar{x}_1)}$$

where

$$se(\bar{x}_2 - \bar{x}_1) = s_p \sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$$

with

$$s_p = \sqrt{\frac{n_1 - 1}{n_1 + n_2 - 2} s_1^2 + \frac{n_2 - 1}{n_1 + n_2 - 2} s_2^2}$$

We reject  $H_0$  for large values of  $|t|$

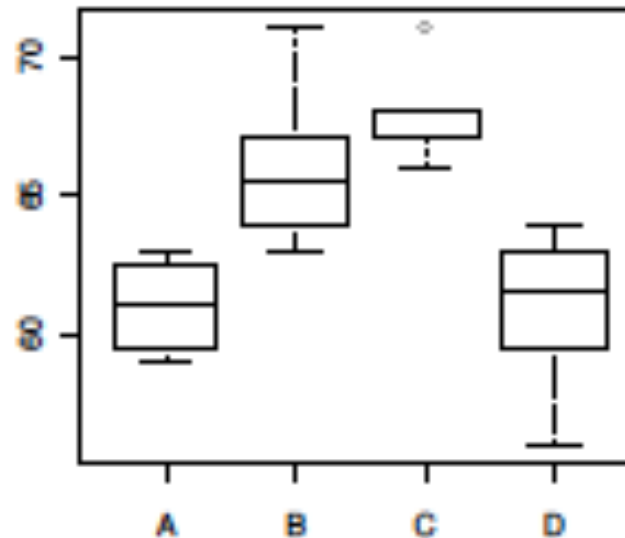
P-value (two-sided) :  $P = 2 P(T > |t|)$ ,

where  $T$  is t-distributed with  $n_1 + n_2 - 2$  df.

# Comparing more than two groups: one-way ANOVA

In an experiment 24 rats were randomly allocated to four different diets, and the blood coagulation time (in seconds) was measured for each animal

Diets (treatment)			
A	B	C	D
62	63	68	56
60	67	66	62
63	71	71	60
59	64	67	61
	65	68	63
	66	68	64
			63
			59



Question: Does diet have an effect on coagulation time?

We may compare two and two diets, using two sample procedures

We would, however, also like to have an overall test

In general we have observations from  $K$  groups:

$x_{ik}$  = observation number  $i$  in group  $k$

$(i = 1, \dots, n_k \quad k = 1, \dots, K)$

We assume that all observations are independent and that the observations from group  $k$  are a random sample from  $N(\mu_k, \sigma^2)$

Notation:

Total number of observations:  $n = \sum_k n_k$

Mean in group  $k$ :  $\bar{x}_k = \frac{1}{n_k} \sum_i x_{ik}$

Overall mean:  $\bar{x} = \frac{1}{n} \sum_{i,k} x_{ik} = \frac{1}{n} \sum_k n_k \bar{x}_k$

We want to test the null hypothesis  $H_0 : \mu_1 = \dots = \mu_K$  versus the alternative that not all the  $\mu_k$  are equal

Introduce the sums of squares:

$$TSS = \sum_{i,k} (x_{ik} - \bar{x})^2 \quad \text{(total sum of squares)}$$
$$MSS = \sum_k n_k (\bar{x}_k - \bar{x})^2 \quad \text{(model sum of squares)}$$
$$RSS = \sum_{i,k} (x_{ik} - \bar{x}_k)^2 \quad \text{(residual sum of squares)}$$

Important decomposition:

$$TSS = MSS + RSS$$



Unbiased estimator of  $\sigma^2$  :

$$s^2 = RSS / (n - K)$$

Under the null hypothesis  $\sigma^2$  may also be estimated by :

$$MSS / (K - 1)$$

However, when the null hypothesis does not hold, the latter estimate tends to be larger than  $\sigma^2$

We reject the null hypothesis for large values of the test statistic

$$F = \frac{MSS / (K - 1)}{RSS / (n - K)}$$

The test statistic is F-distributed with  $K - 1$  and  $n - K$  degrees of freedom under the null hypothesis

This result is used to compute the P-value

The result may be summarized in an ANOVA table:

Source	df	Sum of squares	Mean sum of squares	F statistic	P-value
Model	$K - 1$	$MSS$	$MSS / (K - 1)$	$F = \frac{MSS / (K - 1)}{RSS / (n - K)}$	$P$
Residual	$n - K$	$RSS$	$RSS / (n - K)$		
Total	$n - 1$	$TSS$			

The P-value is found by:

$$P = P(F > \text{observed value of } F)$$

where  $F$  is F-distributed with  $K - 1$  and  $n - K$  degrees of freedom

In Lecture 3 we will see how one-way ANOVA is a special case of multiple linear regression

## R commands for coagulation times:

```
rats=read.table("http://www.uio.no/studier/emner/matnat/math/STK4900/v11/
  rats.txt",header=T)
rats$diet=factor(rats$diet)    # defines diet to be a categorical variable
aov.rats=aov(time~diet,data=rats)
summary(aov.rats)
```

## R output (edited):

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
diet	3	228	76.0	13.6	4.7e-05
Residuals	20	112	5.6		

## Relation to two-sample t-test (two-sided)

Consider the situation with two groups, i.e.  $K = 2$

Will test the null hypothesis  $H_0 : \mu_1 = \mu_2$  versus the alternative hypothesis  $H_A : \mu_1 \neq \mu_2$

t-test statistic:

$$t = \frac{\bar{x}_2 - \bar{x}_1}{se(\bar{x}_2 - \bar{x}_1)}$$

We reject  $H_0$  for large values of  $|t|$

We may show that

$$t^2 = \frac{MSS / (2 - 1)}{RSS / (n - 2)} = F$$

The usual (two-sided) t-test for two samples is a special case of the F-test in one-way ANOVA

## R-commands for bone density example:

```
bonedensity=read.table("http://www.uio.no/studier/emner/matnat/math/  
STK4900/v11/bonedensity.txt",header=T)  
aov.density=aov(density~group,data=bonedensity)  
summary(aov.density)
```

### R-output (edited)

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
group	1	0.00197	0.00197	8.09	0.0082
Residuals	28	0.00681	0.000243		

Note that  $t^2 = 2.844^2 = 8.09 = F$

## Two numerical variables

For one-way ANOVA we study how a numerical variable (e.g. blood coagulation time) depends on a categorical variable (e.g. diet)

Often we want to study the relation between two numerical variables

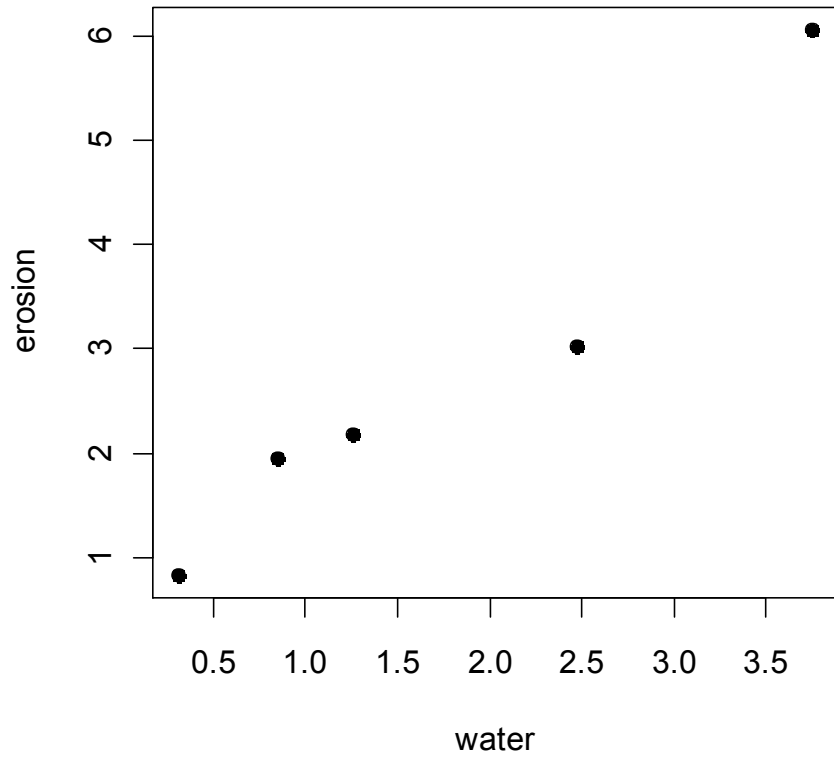
**Example A:** When water flows across a field, some of the soil will be washed away (eroded). An experiment has been performed in order to investigate how the amount of water affects the amount of soil that is eroded.

Amount of water ( <i>l/s</i> )	0.31	0.85	1.26	2.47	3.75
Erosion ( <i>kg</i> )	0.82	1.95	2.18	3.02	6.07

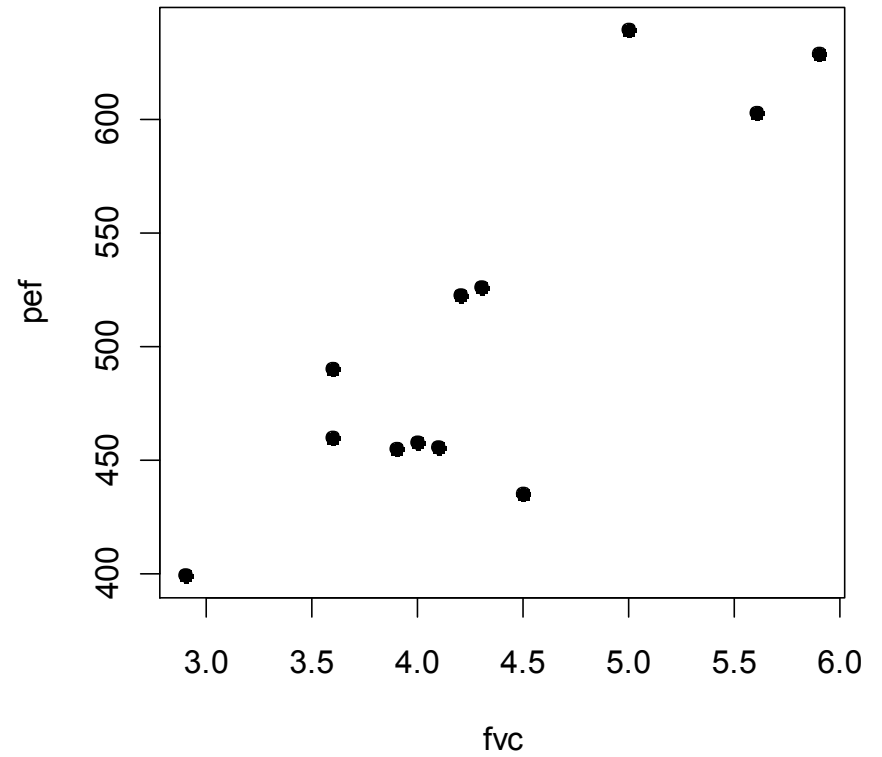
**Example B:** Forced vital capacity (FVC) and peak expiratory flow (PEF) have been measured for 12 adults (in liter and liter per minute, respectively). What is the relation between these two measures of lung function?

Person	1	2	3	4	5	6
FVC	3.9	5.6	4.1	4.2	4.0	3.6
PEF	455	603	456	523	458	460
Person	7	8	9	10	11	12
FVC	5.9	4.5	3.6	5.0	2.9	4.3
PEF	629	435	490	640	399	526

### Example A



### Example B



We will consider two situations:

1. The data  $(x_1, y_1), \dots, (x_n, y_n)$  are considered as independent replications of a pair of random variables  $(X, Y)$
2. The data are described by a linear regression model

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n$$

Here  $y_1, \dots, y_n$  are the outcomes that are considered to be realizations of random variables, while  $x_1, \dots, x_n$  are considered to be fixed (i.e. non-random) and the  $\varepsilon_i$ 's are random errors (noise)

Situation 1 occurs for observational studies (like Example B), while situation 2 occurs for planned experiments, where the values of the  $x_i$ 's are under the control of the experimenter (like Example A)

In situation 1 we will often *condition* on the observed values of the  $x_i$ 's, and analyze the data as if they are from situation 2

We start out by considering situation 1



# Bivariate distributions

We describe the joint distribution of a pair of random variables  $(X, Y)$  through their *bivariate probability density*,  $f(x, y)$

This is defined so that

$$P((X, Y) \in A) = \int_A f(x, y) dx dy$$

The bivariate normal distribution depends on the parameters:

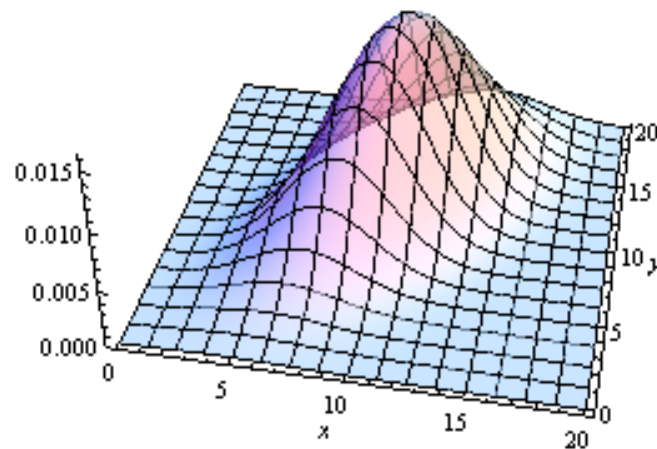
Mean of  $X$  :  $\mu_1$

Mean of  $Y$  :  $\mu_2$

Standard deviation of  $X$  :  $\sigma_1$

Standard deviation of  $Y$  :  $\sigma_2$

Correlation :  $\rho$



# Covariance and correlation

The dependence between  $X$  and  $Y$  may be summarized by the *covariance*:

$$\text{Cov}(X, Y) = E[(X - \mu_1)(Y - \mu_2)]$$

or by the *correlation coefficient*:

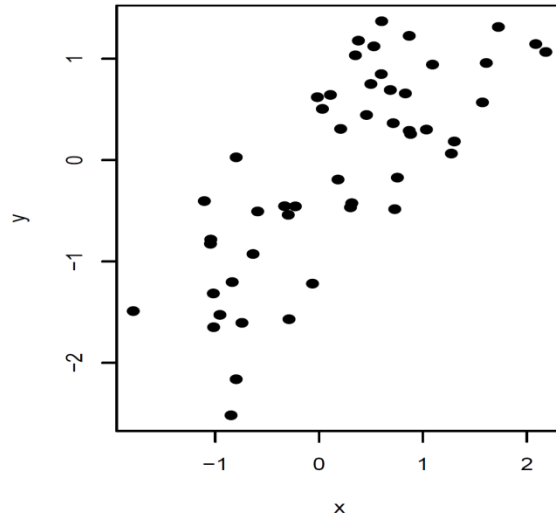
$$\rho = \text{corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\text{sd}(X) \text{sd}(Y)}$$

Important properties of the correlation coefficient:

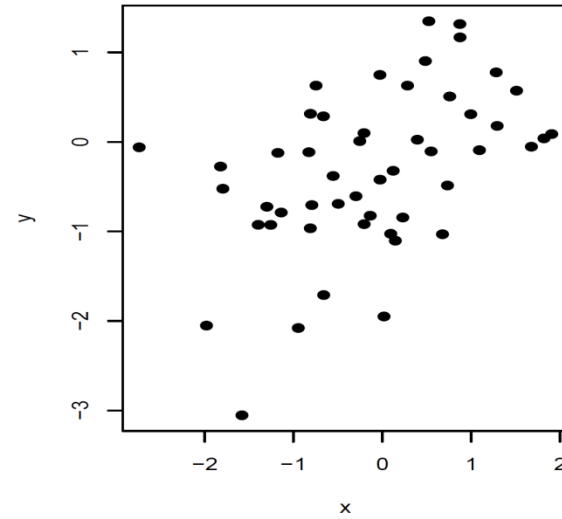
- $\text{corr}(X, Y)$  takes values between -1 and 1
- $\text{corr}(X, Y)$  describes the *linear* relationship between  $Y$  and  $X$
- If  $X$  and  $Y$  are independent, then  $\text{corr}(X, Y) = 0$   
(but not necessarily the other way around)

# Examples of correlated data:

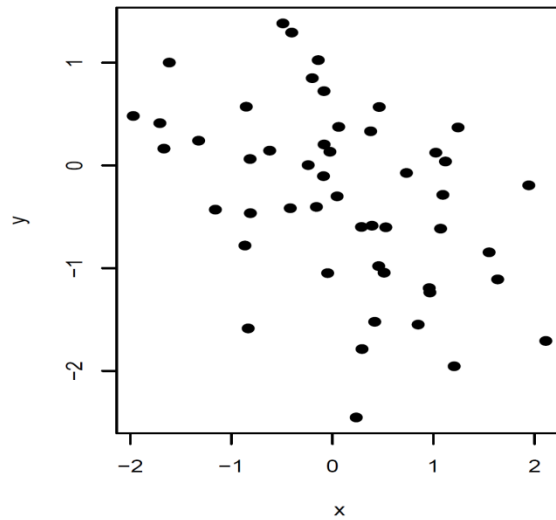
**Correlation 0.9**



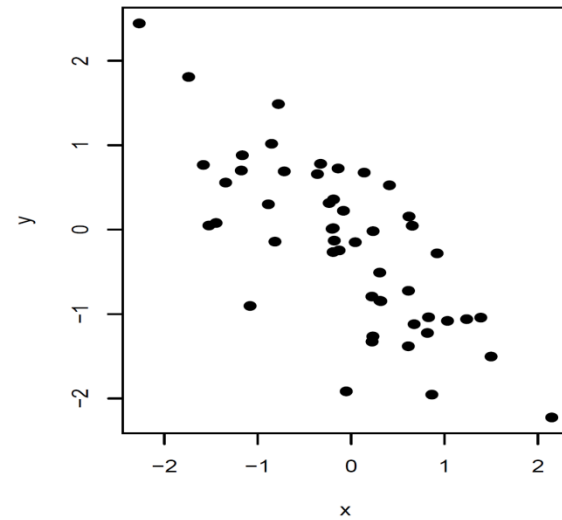
**Correlation 0.5**



**Correlation -0.5**

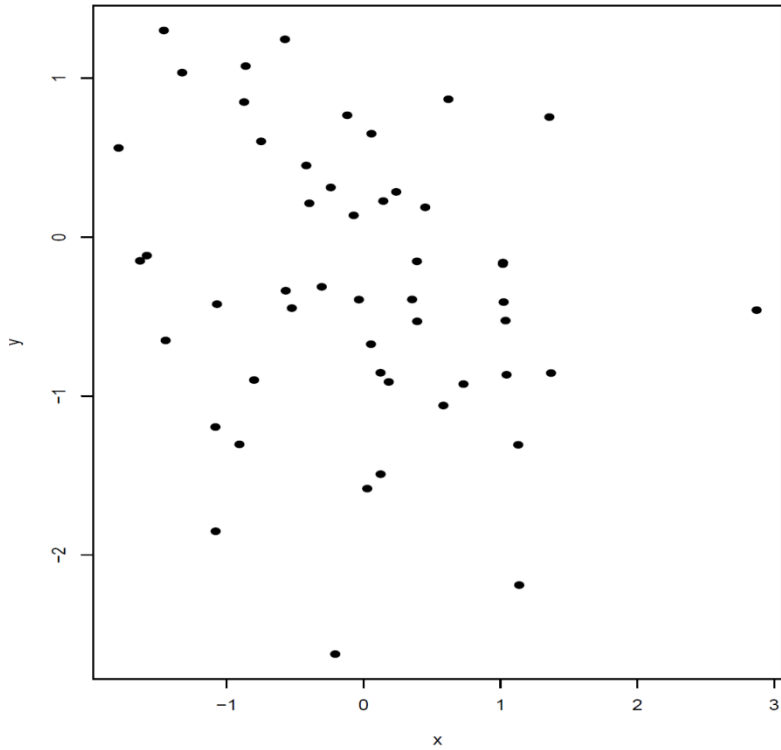


**Correlation -0.9**

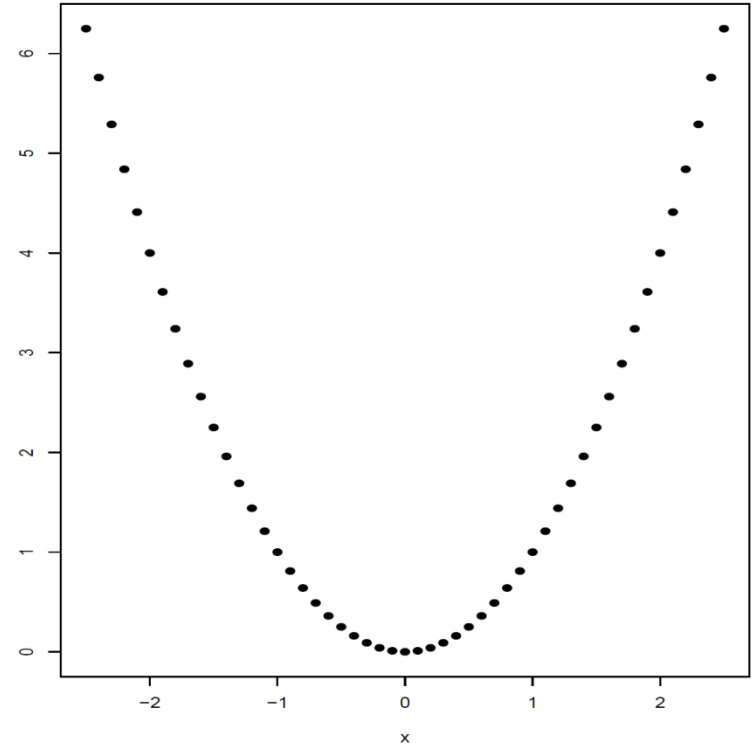


# Examples of uncorrelated data:

Correlation 0.0



Correlation 0.0



## Empirical correlation

The empirical correlation coefficient is an estimator of the theoretical correlation coefficient, and it takes the form

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) / (n - 1)}{s_x \cdot s_y}$$

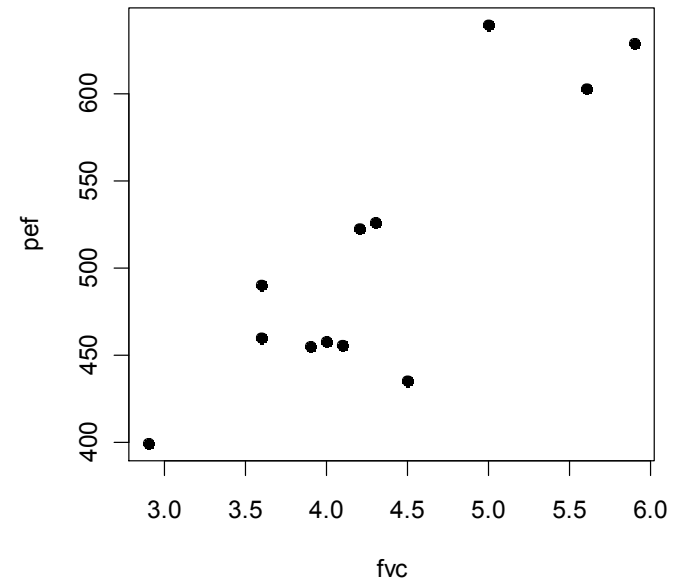
Here  $s_x$  and  $s_y$  are the empirical standard deviations of the  $x_i$ 's and the  $y_i$ 's

$r$  is called the *Pearson correlation coefficient*

The properties of the Pearson correlation coefficient are similar to those of the theoretical correlation coefficient

Consider the example with measures of lung function:

Person	1	2	3	4	5	6
FVC	3.9	5.6	4.1	4.2	4.0	3.6
PEF	455	603	456	523	458	460
Person	7	8	9	10	11	12
FVC	5.9	4.5	3.6	5.0	2.9	4.3
PEF	629	435	490	640	399	526



### R-commands and results:

```
fvc=c(3.9,5.6,4.1,4.2,4.0,3.6,5.9,4.5,3.6,5.0,2.9,4.3)
```

```
pef=c(455,603,456,523,458,460,629,435,490,640,399,526)
```

```
cov(fvc,pef)
```

```
cov(fvc,pef)/(sd(fvc)*sd(pef))
```

```
0.856
```

```
cor(fvc,pef)
```

```
0.856
```

## Test and confidence interval for correlation

We assume that  $(x_1, y_1), \dots, (x_n, y_n)$  are a random sample from a bivariate normal distribution

Consider testing the null hypothesis  $H_0 : \rho = 0$  versus the alternative  $H_0 : \rho \neq 0$

Test statistic:

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$$

We reject  $H_0$  for large values of  $|t|$

Under  $H_0$  the test statistic is t-distributed with  $n - 2$  df

It is more complicated to describe how one may obtain a confidence interval for  $r$  (but one is obtained by the R code on the following slide)

## R-command and results:

```
cor.test(fvc,pef)
```

Pearson's product-moment correlation

data: fvc and pef

t = 5.23, df = 10, p-value = 0.00038

alternative hypothesis: true correlation is not equal to 0

95 percent confidence interval:

0.554 0.959

sample estimates:

cor

0.856

Note that the confidence interval is not symmetric



## Spearman (rank) correlation

The Pearson correlation is sensitive to outliers in the data.

An alternative correlation measure is the Spearman correlation:

The smallest  $x_i$  is replaced by rank  $r_i=1$ ,  
the second smallest  $x_i$  is replaced by rank  $r_i=2$ , and so on to  
the largest  $x_i$  which is replaced by rank  $r_i = n$ .

Similarly, the  $y_i$  are replaced by ranks  $s_i$ .

The Spearman correlation is then simply the Pearson correlation of the ranks  $(r_1, s_1), \dots, (r_n, s_n)$ .

In R:

```
> cor(fvc, pef, method="spearman")  
[1] 0.669
```

# Simple linear regression

We have data  $(x_1, y_1), \dots, (x_n, y_n)$

Here:

$y_i$  = outcome  
(or response)  
(or dependent variable)

$x_i$  = predictor  
(or covariate)  
(or explanatory variable)  
(or independent variable)

Model:

$$y_i = E(y_i | x_i) + \varepsilon_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

where the  $x_i$ 's are considered to be fixed quantities, and the  $\varepsilon_i$ 's are independent error terms ("noise") that are assumed to be  $N(0, \sigma_\varepsilon^2)$ -distributed

Consider the erosion example:

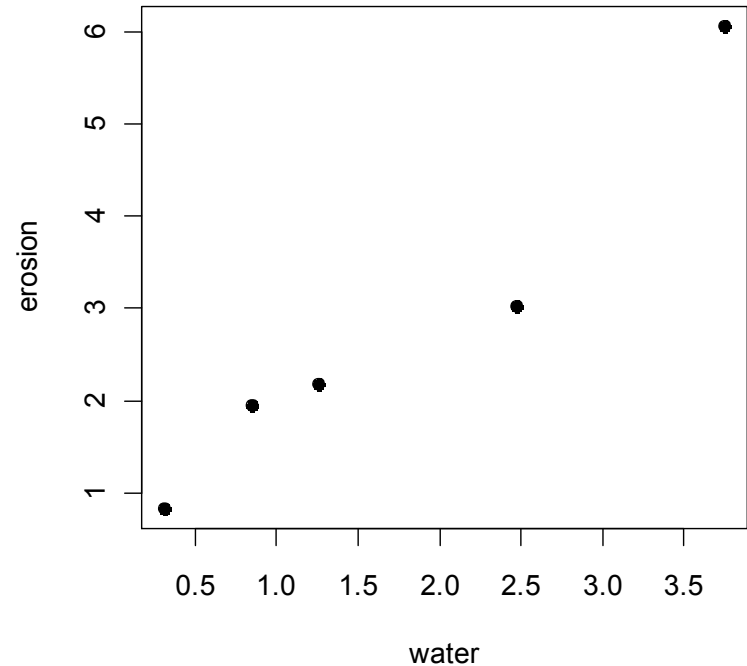
Amount of water ( $l/s$ )	0.31	0.85	1.26	2.47	3.75
Erosion ( $kg$ )	0.82	1.95	2.18	3.02	6.07

Response = erosion

Predictor = amount of water

Model:

$$\text{erosion} = \beta_0 + \beta_1 \text{water} + \varepsilon$$

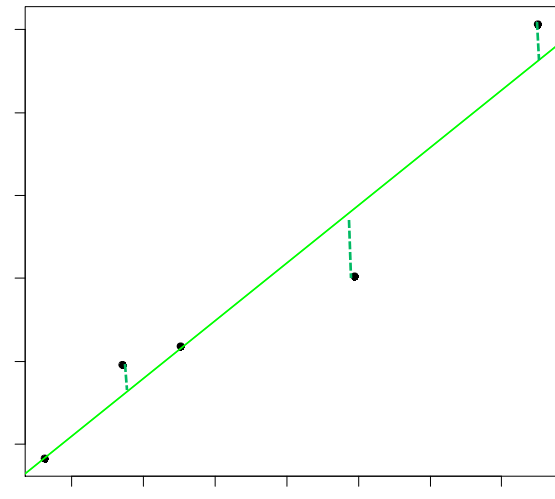
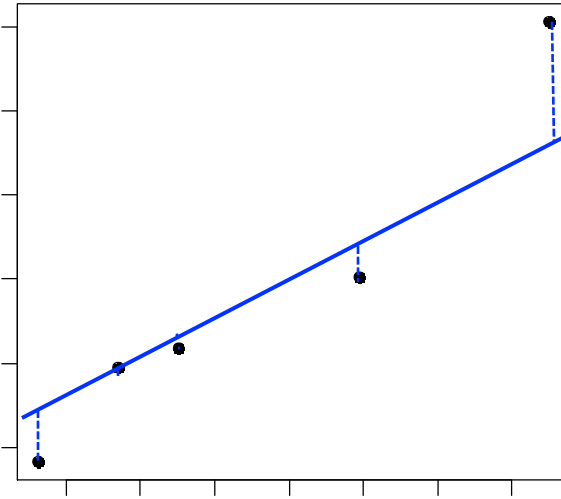
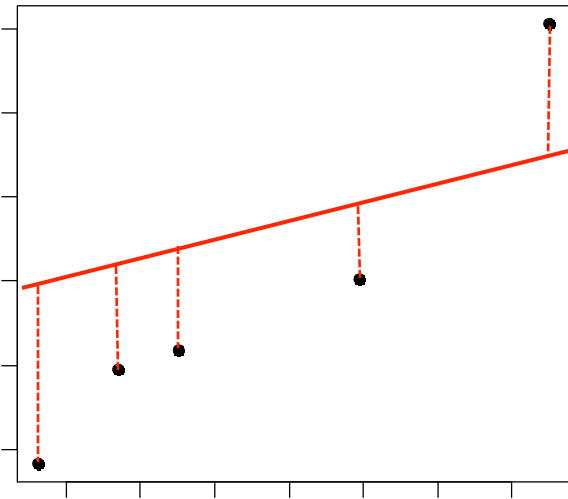


# Least squares

We estimate the regression coefficients using the method of least squares, i.e. the estimates  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are obtained as the values of  $b_0$  and  $b_1$  that minimize the sum of squares

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

Illustration:



## R-commands:

```
water=c(0.31,0.85,1.26,2.47,3.75)
erosion=c(0.82,1.95,2.18,3.02,6.07)
fit=lm(erosion~water)
summary(fit)
plot(water,erosion,pch=19)
abline(fit)
```

## R-output (edited)

Coefficients:

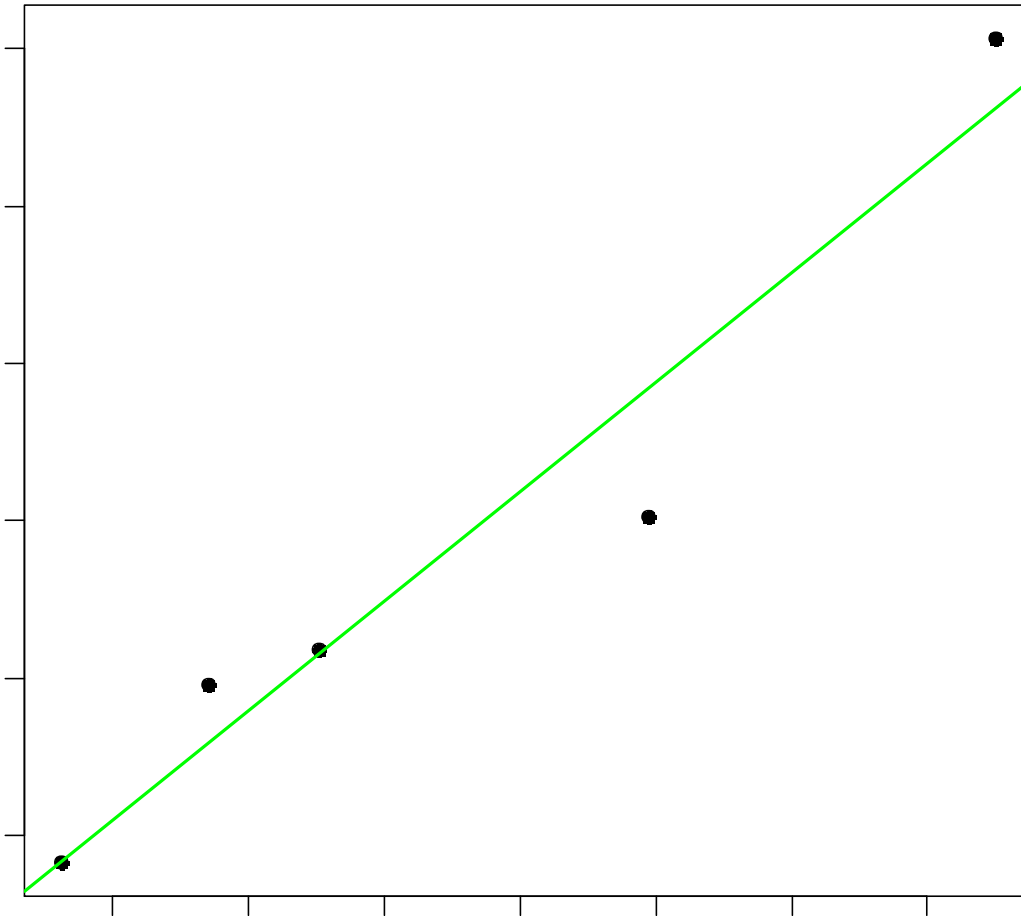
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.406	0.445	0.912	0.429
water	1.390	0.210	6.630	0.007

Residual standard error: 0.580 on 3 degrees of freedom

Multiple R-squared: 0.936, Adjusted R-squared: 0.915

F-statistic: 44.0 on 1 and 3 DF, p-value: 0.007

"Estimate" denotes the least squares estimates (the meaning of the other parts of the output will be made clear in the following)



Fitted regression line:  $\text{erosion} = 0.406 + 1.390 \times \text{water}$

# Fitted values and residuals

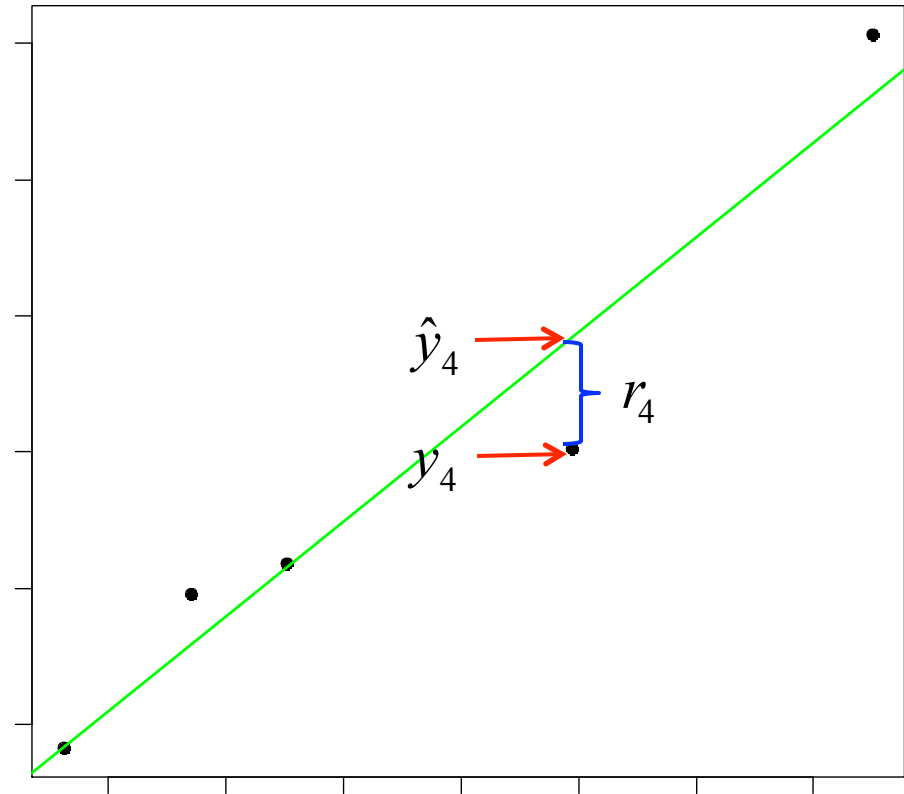
Fitted values:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

Residuals:

$$r_i = y_i - \hat{y}_i$$

The residuals are estimates of the unobserved  $e_i$ 's



## Sums of squares

In a similar manner as for one-way ANOVA, we have the sums of squares:

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2 \quad \text{(total sum of squares)}$$
$$MSS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad \text{(model sum of squares)}$$
$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{(residual sum of squares)}$$

Decomposition:

$$TSS = MSS + RSS$$



## Standard errors

Unbiased estimator of  $\sigma_\varepsilon^2$  :

$$\widehat{\text{Var}}(\varepsilon) = s_{y|x}^2 = \text{RSS} / (n - 2)$$

$s_{y|x}$  is the "residual standard error" in the R output

The variance of  $\hat{\beta}_1$  is estimated by :

$$\widehat{\text{Var}}(\hat{\beta}_1) = \frac{s_{y|x}^2}{(n - 1)s_x^2}$$

where  $s_x^2 = \sum_{i=1}^n (x_i - \bar{x})^2 / (n - 1)$  is the sample variance of the  $x_i$ 's

Standard error:  $se(\hat{\beta}_1) = \sqrt{\widehat{\text{Var}}(\hat{\beta}_1)}$

Similar formulas hold for the variance and standard error of  $\hat{\beta}_0$

The standard errors are denoted "Std. Error" in the R output

## Hypothesis tests

Consider testing the null hypothesis  $H_0 : \beta_1 = 0$  versus the alternative  $H_A : \beta_1 \neq 0$

Test statistic:

$$t = \frac{\hat{\beta}_1}{se(\hat{\beta}_1)}$$

We reject  $H_0$  for large values of  $|t|$

Under  $H_0$  the test statistic is t-distributed with  $n - 2$  df

P-value (two-sided) :  $P = 2 P(T > |t|)$ ,

where  $T$  is t-distributed with  $n - 2$  df.

Testing the null hypothesis  $H_0 : \beta_0 = 0$  is performed similarly (but is usually not of much interest)

t-statistics and P-values are given in the R output as "t value" and "Pr(>|t|)"

## Confidence intervals

95% confidence interval for  $\beta_1$  :

$$\hat{\beta}_1 \pm c \cdot se(\hat{\beta}_1)$$

where  $c$  is the upper 97.5% percentile in the t-distribution with  $n - 2$  df

95% confidence interval in the erosion example:

$$1.39 \pm 3.18 \cdot 0.210$$

i.e. from 0.72 to 2.06

Note that the confidence interval does not contain 0 if and only if the P-value for the test is less than 5%

# Correlation and regression

The least squares estimate for the slope is given by:

$$\hat{\beta}_1 = r \frac{s_y}{s_x}$$

where

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) / (n - 1)}{s_x \cdot s_y}$$

is the Pearson correlation coefficient (and  $s_x$  and  $s_y$  are the empirical standard deviations of the  $x_i$ 's and the  $y_i$ 's)

Further the test for  $H_0 : \beta_1 = 0$  in a linear regression model (slide 33) is numerically equivalent to the test for  $H_0 : \rho = 0$  for bivariate data (slide 23)

## Coefficient of determination

The coefficient of determination is given by

$$R^2 = \frac{MSS}{TSS} = 1 - \frac{RSS}{TSS}$$

This may be interpreted as the proportion of the total variability in the outcomes (TSS) that is accounted for by the model (MSS)

$R^2$  is given as "Multiple R-squared" in the R output

For the simple linear regression model  $R^2$  is just the square of the Pearson correlation coefficient:

$$R^2 = r^2$$