# TMA4267 Linear Statistical Models V2014 (12)
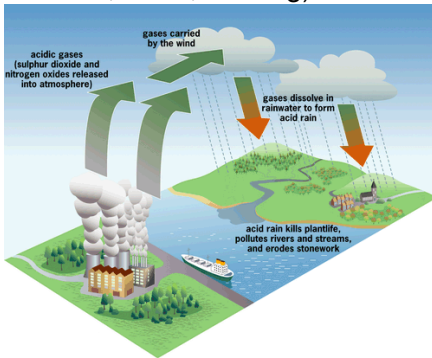## Multiple linear regression, normal equations [3.1-3.2]

Mette Langaas

To be lectured: February 11, 2014
wiki.math.ntnu.no/emner/tma4267/2014v/start/

# Acid rain

occurs when emissions of sulfur dioxide ($SO_2$) and oxides of nitrogen ($NO_x$) react in the atmosphere with water, oxygen, and oxidants to form various acidic compounds. These compounds then fall to the earth in either dry form (such as gas and particles) or wet form (such as rain, snow, and fog).
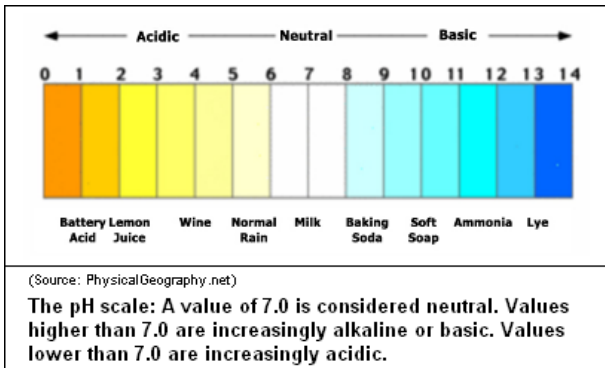
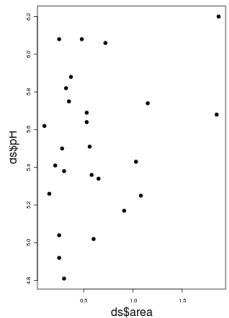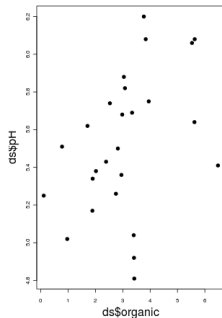Source: http://myecoproject.org/get-involved/pollution/acid-rain/

Mette.Langaas@math.ntnu.no, TMA4267V2014

http://www.eoearth.org/view/article/149814/

# **Acid rain in Norwegian lakes**

Measured pH in Norwegian lakes explained by content of

— $x_1$: $SO_4$: sulfate (the salt of sulfuric acid),

— $x_2$: $NO_3$: nitrate (the conjugate base of nitric acid),

— $x_3$: $Ca$: calsium,

— $x_4$: latent $Al$: aluminium,

— $x_5$: organic substance,

— $x_6$: area of lake,

— $x_7$: position of lake (Telemark or Trøndelag),

pH is a measure of the acidity of alkalinity of water, expressed in terms of its concentration of hydrogen ions. The pH scale ranges from 0 to 14. A pH of 7 is considered to be neutral. Substances with pH of less that 7 are acidic; substances with pH greater than 7 are basic.

Acidic ——— Neutral ——— Basic

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14

Battery Lemon    Wine    Normal    Milk    Baking    Soft    Ammonia    Lye
Acid Juice          Rain        Soda    Soap

(Source: PhysicalGeography.net)

The pH scale: A value of 7.0 is considered neutral. Values higher than 7.0 are increasingly alkaline or basic. Values lower than 7.0 are increasingly acidic.

http://www.eoearth.org/view/article/149814/

0=Telemark, 1=Trondelag

# Acid rain data

# Part 4: Multiple linear regression

## [Bingham & Fry chapter 3]

In ch1 we studied simple linear regression model

$$Y_i = \alpha + \beta x_i + \varepsilon_i \quad \text{where } i = 1,..,n$$

$$\varepsilon_i \text{ s.i.d } N(0, \sigma^2)$$

We found
$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{x}$$

$$\hat{\beta} = \frac{S_{xY}}{S_{xx}} = \frac{\sum_{i=1}^{n}(Y_i - \bar{Y})(x_i - \bar{x})}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$$
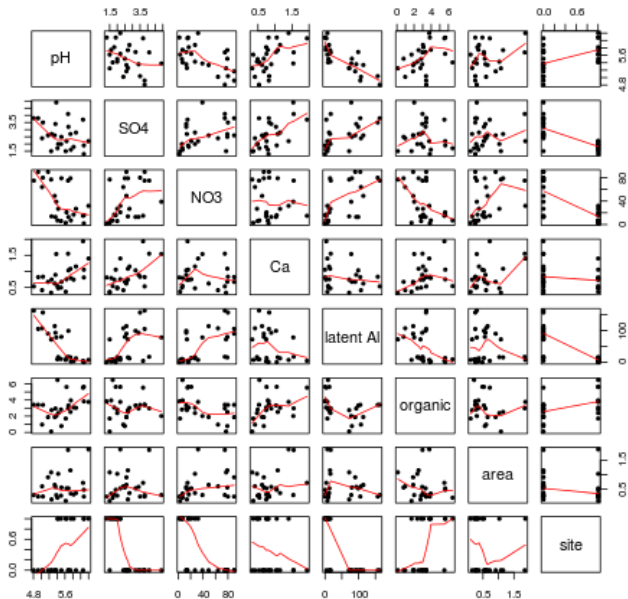
(LS)

are the least squares and maximum likelihood (ML) estimators
of $\alpha$ and $\beta$, and

$$\hat{S}^2 = \frac{1}{n}\sum_{i=1}^{n}(Y_i - \hat{Y}_i)^2$$

$$\uparrow \hat{\alpha} + \hat{\beta}x_i$$

is the ML estimator for $\sigma^2$.

We also saw that
$$\begin{bmatrix} X \\ Y \end{bmatrix} \sim N_2 \left( \begin{bmatrix} \mu_X \\ \mu_Y \end{bmatrix}, \begin{bmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{bmatrix} \right)$$ then

$$E(Y \mid X = x) = \alpha + \beta x = \left( \mu_Y - \rho\frac{\sigma_y}{\sigma_x} \right) + \left( \rho\frac{\sigma_y}{\sigma_x} \right) x$$

is a linear function in $x$.

Now: $x_{1i} = 1$ $\forall i = 1, \ldots, n$ ← number of observation

to treat the intercept in the same way as the slope, and

$$Y_i = \overbrace{\beta_1 \cdot x_{i1}}^{\text{old } \alpha} + \overbrace{\beta_2 \cdot x_{i2}}^{\text{old } \beta x_i} + \ldots + \beta_p x_{ip} + \varepsilon_i \qquad i = 1, \ldots, n$$

$$\varepsilon_i \text{ i.i.d } N(0, \sigma^2)$$

$$Cov(\varepsilon_i, \varepsilon_j) = 0 \quad \forall i \neq j$$

We want vectors and matrices!

Let

$$\underset{n \times 1}{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$ be a vector of RVs

$n = \#$ observations

$$\underset{n \times p}{\underline{X}} = \begin{bmatrix} x_{11} & x_{12} & \ldots & x_{1p} \\ x_{21} & & & \\ \vdots & & & \vdots \\ x_{n1} & x_{n2} & & x_{np} \end{bmatrix}$$ is a design matrix X

↑

reflecting the design of the experiment.

The x's may be set by design - as we will see in Design of Experiments (DOE) or be observed together with Y in an observational study (acid rain).

We look at X as a matrix of **constants** — not RV's

$$\beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} \quad \text{(p×1)}$$

vector of regression parameters

Remark: more common to use $\beta_0, \ldots, \beta_p$, but we stay with the notation of the book.

$$\varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix} \quad \text{(n×1)}$$

is a random vector of errors, where

$E(\varepsilon) = 0$ and $\text{Var}(\varepsilon_i) = \sigma^2$

and $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \quad \forall i \neq j \Rightarrow \text{Cov}(\varepsilon) = \sigma^2 I$ (n×n)

The regression model:

$$Y = X \beta + \varepsilon \quad , \quad \varepsilon \sim N_n(0, \sigma^2 I)$$

$Y$: (n×1) — random vector — response vector

$X$: (n×p) — design matrix

$\beta$: (p×1) — parameter vector

$\varepsilon$: (n×1) — random vector — error vector

# Questions about *X*

1. Why do we want to assume that the design matrix *X* has full rank?
2. Can we find $X^{-1}$?

Q: what are typical values of $n$ and $p$?

$$10 - 10\,000 \qquad 1 - 10$$

We assume that $n \gg p$, i.e. the number of observations $n$ is much larger than the number of covariates $p$.

And we will assume that $\underset{n \times p}{X}$ has full rank $p$.

Q: Why do we want to assume that $X$ has full rank $p$?

We don't want to include covariates that are linear combinations of eachother — that adds no information.

Q: Can we find $\underset{n \times p}{X}^{-1}$?  $X$ is not quadratic when $n \gg p$

NO.

$\Rightarrow$ proceed to estimate $\underline{\beta \ \& \ \sigma^2}$

# Multiple linear regression model

$$Y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip} + \varepsilon_i, \text{ for } i = 1, ..., n$$

Matrix formulation:

$$\underset{(n \times 1)}{\boldsymbol{Y}} = \underset{(n \times p)(p \times 1)}{\boldsymbol{X}\,\boldsymbol{\beta}} + \underset{(n \times 1)}{\boldsymbol{\varepsilon}}$$

$$E(\varepsilon) = \underset{(n \times 1)}{\boldsymbol{0}} \quad \text{and} \quad Cov(\varepsilon) = \underset{(n \times 1)}{\sigma^2 \boldsymbol{I}}$$

where

— $\beta$ and $\sigma^2$ are unknown parameters and

— the design matrix $\boldsymbol{X}$ has $i$th row $[x_{i1}, x_{i2}, ..., x_{ip}]$.

# Estimation of $\beta$

We have $Y = X\beta + \varepsilon$, with elements

$$Y_i = \underset{\underset{\text{row vector}: 1\times p}{\uparrow}}{X_i} \underset{p\times 1}{\beta} + \varepsilon_i \qquad \varepsilon_i \text{ i.i.a } N(0, \sigma^2)$$

and thus $Y_i$'s are independent and $N(X_i\beta, \sigma^2)$.

## ML estimation

$$L(\beta, \sigma) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} \exp\left\{ -\frac{1}{2\sigma^2} (y_i - X_i\beta)^2 \right\}$$

$$= \left(\frac{1}{2\pi}\right)^{n/2} \frac{1}{\sigma^n} \exp\left\{ -\frac{1}{2\sigma^2} \underbrace{\sum_{i=1}^{n} (y_i - X_i\beta)^2}_{\underbrace{(y - X\beta)^T (y - X\beta)}_{Q(\beta)}} \right\}$$

Maximizing $L$ wrt $\beta$ will be the same as minimizing $Q(\beta)$ [as before... ].

Remark:

$$e = y - Xb$$

residuals $n \times 1$ — response $n \times 1$ — estimates $n \times 1$

$$Q(b) = e^T e = (y - Xb)^T (y - Xb)$$

Task: $\dfrac{\partial Q(b)}{\partial b} = 0$    solve to find b.

↖ set of p equations

$\dfrac{\partial Q(b)}{\partial b_1}$

$\dfrac{\partial Q(b)}{\partial b_2}$

$\vdots$

$\dfrac{\partial Q(b)}{\partial b_p}$

Need two simple rules for derivatives wrt vectors.

$$\frac{\partial}{\partial b} (d^T b) = d$$

$\sum_{i=1}^{p} d_i b_i$    $p \times 1$

$$\frac{\partial}{\partial b} (b^T D b) = (D + D^T) b$$

$\sum_{j=1}^{p} \sum_{k=1}^{p} b_j d_{jk} \cdot b_k$

if $D = D^T$

$\Rightarrow 2Db$

$$Q(b) = (y - Xb)^\top (y - Xb)$$

$$= y^\top y - y^\top Xb - b^\top X^\top y + b^\top X^\top Xb$$

$$= y^\top y - 2\underbrace{y^\top X}_{d^\top} b + b^\top \underbrace{X^\top X}_{D} b \qquad (X^\top X)^\top = X^\top X$$

$$\frac{\partial Q(b)}{\partial b} = 0 - 2(y^\top X)^\top + (X^\top X + (X^\top X)^\top)b = 0$$

$$-2 X^\top y + 2(X^\top X)b = 0$$

"$Ax = b$"

$$\underline{(X^\top X)b = X^\top y} \qquad \text{Normal equations (NE)}$$

$$\underbrace{\overbrace{p \times n}^{} \ \overbrace{n \times p}^{}}_{p \times p \ \ (p \times)} \qquad \underbrace{p \times n \ \ n \times 1}_{(p \times 1)}$$

How can we solve this? Since $X$ has full rank, then $X^\top X$ will also have full rank and $(X^\top X)^{-1}$ exists.

$$b = (X^\top X)^{-1} X^\top y$$

And the LS estimator

$$\hat{\beta} = (X^\top X)^{-1} X^\top Y \qquad \swarrow RV$$

Bonk: $(X^\top X)$ called information matrix

# Least squares estimation

— $\boldsymbol{Y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ with $E(\boldsymbol{\varepsilon}) = \boldsymbol{0}$ and $\mathrm{Cov}(\boldsymbol{\varepsilon}) = \sigma^2 \boldsymbol{I}$.

— Let $\boldsymbol{X}$ has full rank $p \leq n$.

— The Least Squares Estimate (LSE) of $\boldsymbol{\beta}$ is given by

$$\hat{\boldsymbol{\beta}} = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{Y}$$

Since X has full rank $X^T X$ is PD i.e.

$$z^t (X^T X) z > 0 \qquad \forall \ z \neq 0$$

Proof: book p 64-65

Near linear dependence, called multicollinearity,
will make multiple linear regression numerically unstable
and the interpretation of $\hat{\beta}$ will be difficult.