# TMA4267 Linear Statistical Models V2017 [L7]

Part 2: Linear regression [F p73-86] Model definition [F3.1], Parameters and residuals [F3.1.1], Model check [F3.1.2]

# Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 7, 2017

1/20

# Age-predicted maximal heart rate in healthy subjects: The HUNT Fitness Study

B. M. Nes, I. Janszky, U. Wisløff, A. Støylen, T. Karlsen (2012) in Scandinavian Journal of Medicine and Science in Sports.

- ▶ HRmax describes the highest heart rate achieved by a subject exercising to exhaustion and is verified by a plateau of heart rate despite increasing workload. In the literature, HRmax commonly refers to the peak heart rate at termination of a graded maximal exercise test.
- However, in clinical settings, a maximal exercise test is not always feasible and there is a need to predict HRmax from age prior to testing to be able to adequately assess heart rate response and relative intensity of effort at submaximal levels.

# Part 2: Linear regression

#### Part 2: Linear regression

► Fahrmeir et al (2013): Regression. Chapter 3.1, 3.2, 3.4 and required parts of 3.5 and Appendix B.

#### Part 3: Hypothesis testing and analysis of variance

- ► Fahrmeir et al (2013): Regression. Chapter 3.3 and required parts of 3.5 and Appendix B.
- ► Härdle et al (2015): Applied Multivariate Statistical Analysis. Chapter 8.1.1. (ANOVA).
- A short note on multiple testing (to be written).

File TMA4267Part2and3.pdf available from course www-page.

1/20

# Age-predicted maximal heart rate in healthy subjects: The HUNT Fitness Study

- HRmax at a given age is frequently estimated by the "220 age" formula.
- ▶ The aim of the present study was to develop a new prediction formula for HRmax through analysis of HRmax measured at VO2peak in a diverse population of 4635 healthy subjects and compare this formula with three commonly used prediction formulas. Furthermore, we wanted to investigate the relationship between HRmax and gender, physical activity status, BMI, and objectively measured aerobic fitness.

# Age-predicted maximal heart rate in healthy subjects: The HUNT Fitness Study - Statistical procedures

- ▶ Only subjects that fulfilled the criteria of a maximal test, with registered maximal heart rate (HRmax), were included in the analysis (n = 3320).
- ▶ General linear modeling was used to determine the effect of age on HRmax. HRmax was entered as the dependent variable and age as the independent variable. Nonlinearity of the relationship between age and HRmax was investigated by including polynomial terms to the regression model.
- ▶ In a subsequent analysis, the effects of gender, BMI, physical activity status, and maximal oxygen uptake were examined by entering these factors as independent variables in addition to age. In further subsequent models, interaction terms were included as well to assess effect modification.
- The continuous variables were checked for normality, homogeneity of variances, and heteroscedasticity of the residuals.

4/20

#### Munich Rent Index data set

3rd Qu.: 559.36 3rd Qu.: 8.8408

Max. :1843.38 Max. :17.7216

1:1794 0:2891 0:2951 0:321 Min. : 113 2:1210 1: 191 1: 131 1:2761 1st Qu.: 561

location bath

described in Fahrmeir et al (2013) on pages 19-20.

```
> library("gamlss.data")
> ds=rent99
> dim(ds)
Γ17 3082
> colnames(ds)
[1] "rent" "rentsqm" "area"
                                   "yearc"
                                                "location" "bath"
[7] "kitchen" "cheating" "district"
    rent
 Min. : 40.51 Min. : 0.4158 Min. : 20.00
 1st Qu.: 322.03 1st Qu.: 5.2610
Median: 426.97 Median: 6.9802
                         Median : 65.00
Mean : 459.44 Mean : 7.1113
                         Mean : 67.37
```

3rd Qu.: 81.00

Max. :160.00

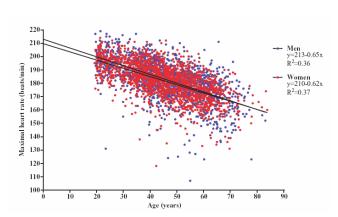
Median :1025 Mean :1170 3rd Qu.:1714

kitchen cheating district

3rd Qu.:1972

Max. :1997

6 / 20



Nes et al (2012): Age-predicted maximal heart rate in healthy subjects: The HUNT Fitness Study. n = 3320 individuals.

5 / 20

### The classical linear model

The model

$$Y = X\beta + \varepsilon$$

is called a classical linear model if the following is true:

- 1.  $E(\varepsilon) = 0$ .
- 2.  $Cov(\varepsilon) = E(\varepsilon \varepsilon^T) = \sigma^2 I$ .
- 3. The design matrix has full rank, rank(X) = k + 1 = p.

The classical *normal* linear regression model is obtained if additionally

4. 
$$\varepsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$$

holds. For random covariates these assumptions are to be understood conditionally on  $\boldsymbol{X}$ .

### Conditional mean and covariance

If we believe that the vector with elements Y and X are multivariate normal  $N_{k+1}(\mu, \Sigma)$  we may look at the partition

$$\begin{pmatrix} Y \\ X \end{pmatrix} \sim N_{k+1} \begin{pmatrix} \mu_Y \\ \mu_X \end{pmatrix}, \begin{pmatrix} \Sigma_{YY} & \Sigma_{YX} \\ \Sigma_{XY} & \Sigma_{XX} \end{pmatrix} \end{pmatrix}$$

The conditional distributions of the components are (multivariate) normal, with conditional mean and variance of  $Y \mid \mathbf{X} = \mathbf{x}$  are

$$E(Y \mid \mathbf{X} = \mathbf{x}) = \mu_Y + \mathbf{\Sigma}_{YX} \mathbf{\Sigma}_{XX}^{-1} (\mathbf{x} - \mu_X)$$

$$Var(Y \mid \mathbf{X} = \mathbf{x}) = \mathbf{\Sigma}_Y - \mathbf{\Sigma}_{YX} \mathbf{\Sigma}_{XX}^{-1} \mathbf{\Sigma}_{XY}$$

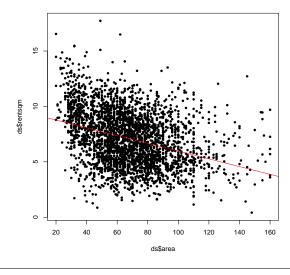
Observe: mean is linear in x and variance independent of x.

8 / 20

10/20

# Linearity of covariates: Covariate vs. response

Munich Rent Index: area vs rentsqm



Model assumptions for the classical linear model [F:3.1.2]

What are our model assumptions, how can we spot violations and what can we do to amend the violations.

- 1. Linearity of covariates:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
- 2. Homoscedastic error variance:  $Cov(\varepsilon) = \sigma^2 I$ .
- 3. Uncorrelated errors:  $Cov(\varepsilon_i, \varepsilon_i) = 0$ .
- 4. Additivity of errors:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

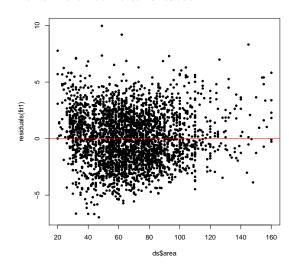
We mainly use plots to assess this (more on model fit in F:3.4 Model choice and variable seletion)

- ► Covariate vs response (for each covariate)
- Covariate vs error (when we have simulated data and know the truth)
- Covariate vs residual (estimated error),
- ▶ Predicted response vs residual (to be popular later).

9 / 20

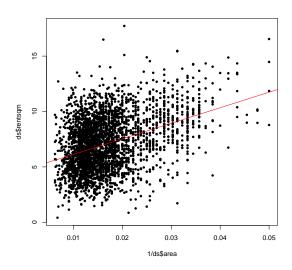
# Linearity of covariates: Covariate vs. residual (residual plot)

Munich Rent Index: area vs residual



# Linearity of covariates: Transformed covariate vs. response

Munich Rent Index: 1/area vs rentsqm



12/20

#### 3.2 Modeling Nonlinear Covariate Effects Through Variable Transformation

If the continuous covariate z has an approximately nonlinear effect  $\beta_1 f(z)$ with known transformation f, then the model

$$v_i = \beta_0 + \beta_1 f(z_i) + \ldots + \varepsilon_i$$

can be transformed into the linear regression model

$$y_i = \beta_0 + \beta_1 x_i + \ldots + \varepsilon_i,$$

where  $x_i = f(z_i) - \bar{f}$ . By subtracting

$$\bar{f} = \frac{1}{n} \sum_{i=1}^{n} f(z_i),$$

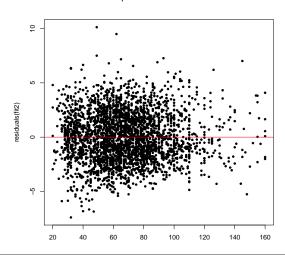
the estimated effect  $\hat{\beta}_1 x$  is automatically centered around zero. The estimated curve is best interpreted by plotting  $\hat{\beta}_1 x$  against z (instead of x).

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.94)

14/20

# Linearity of covariates: Transformed covariate vs. residual (residual plot)

Munich Rent Index: 1/area vs residual



13 / 20

#### 3.3 Modeling Nonlinear Covariate Effects Through Polynomials

If the continuous covariate z has an approximately polynomial effect  $\beta_1 z$  +  $\beta_2 z^2 + \ldots + \beta_l z^l$  of degree l, then the model

$$y_i = \beta_0 + \beta_1 z_i + \beta_2 z_i^2 + \ldots + \beta_l z_i^l + \ldots + \varepsilon_i$$

can be transformed into the linear regression model

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_1 x_{i2} + \ldots + \beta_l x_{il} + \ldots + \varepsilon_i$$

where  $x_{i1} = z_i$ ,  $x_{i2} = z_i^2$ , ...,  $x_{il} = z_i^l$ . The centering (and possibly orthogonalization) of the vectors  $\mathbf{x}^j =$  $(x_{1i},\ldots,x_{ni})', j=1,\ldots,l,$  to  $x^1-\bar{x}_1,\ldots,x^l-\bar{x}_l$  with the mean vector  $\bar{x}_i = (\bar{x}_i, \dots, \bar{x}_i)'$  facilitates interpretation of the estimated effects. A graphical illustration of the estimated polynomial is a useful way to interpret the estimated effect of z.

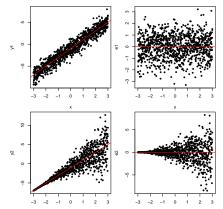
Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.95)

# Homoscedastic errors

```
n=1000
x=seq(-3,3,length=n)
beta0=-1
beta1=2
xbeta=beta0+beta1*x
sigma=1
e1=rnorm(n,mean=0,sd=sigma)
y1=xbeta+e1
ehat1=residuals(lm(y1~x))
plot(x,y1,pch=20)
abline(beta0,beta1,col=1)
plot(x,e1,pch=20)
abline(h=0,col=2)
```

16/20

# Homo- and heteroscedastic errors



Top: homoscedastic errors. Bottom: heteroscedastic errors. Right: x vs y. Left: x vs error. Example from Fahrmeir et al (2013): Regression. Springer. (p.79). R code from TMA4267 lectures tab.

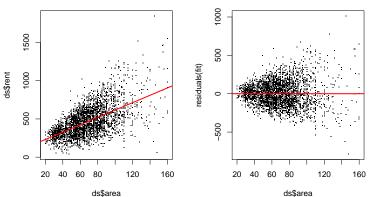
18 / 20

# Heteroscedastic errors

```
sigma=(0.1+0.3*(x+3))^2
e2=rnorm(n,0,sd=sigma)
y2=xbeta+e2
ehat2=residuals(lm(y2~x))
plot(x,y2,pch=20)
abline(beta0,beta1,col=2)
plot(x,e2,pch=20)
abline(h=0,col=2)
```

17 / 20

# Homoscedastic errors?



Left: area vs rent, right: area vs residuals. Fahrmeir et al (2013): Regression. Springer. (p.80). R code from TMA4267 lectures tab.

# Today

- ▶ Normal linear model: implication for Y.
- ▶ Model parameters  $\beta$ ,  $\sigma^2$ , parameter estimators  $\hat{\beta}$ ,  $\hat{\sigma}^2$ , residuals  $\hat{\boldsymbol{\varepsilon}} = \boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}}.$
- ► Model assumptions.
- ▶ Next: covariates- how to include in linear regression, and then parameter estimation.

20 / 20

Deta and design matrix

We collect independent data 
$$(Y_i, X_i)$$
 for  $i=1,...,n$ 

response

 $Y_i = \begin{bmatrix} Y_i \\ Y_i \end{bmatrix}, E : \begin{bmatrix} E_1 \\ E_2 \\ E_n \end{bmatrix}$ 
 $X = \begin{bmatrix} 1 & X_{11} & X_{12} & ... & X_{12} \\ 1 & X_{13} & X_{12} & ... & X_{12} \end{bmatrix}$ 

The properties of the value assume that  $n >> p$ .

The number of observations  $p = n$  when  $p = n$  and  $p = n$  when  $p = n$  when  $p = n$  when  $p = n$  and  $p = n$  and  $p = n$  when  $p = n$  and  $p = n$  and

THAYRET L7 PART 2: Of.02.2017 LINEAR REBRESSION

#### Model definition [F31.0]

Y = variable of primary interest (response, dependent

X1, X2,... , X1 = regressors, explenetory vanables independent variables, coveration

Assumptions:

1) Systematic component is a linear combination of the coveriors. f(x1, x2,.., x4) = \$0+\$1x1+\$2x2+ ...+ \$1. Xe

2) Additive erros Y= xTB+E Restrictive? Haybe > Fronformations?

I is chosen so that we have a good noted

#### The classical linear model

1) E(E)=0

Var (Ei)= of for all i ← homoscedastic error Cov (εί, εί) = 0 for itje-uncorrelated errors

- S) The design metrix has rank rank (X)= h+1= p
- 4) If we in addition essume that En Nn (0, 0. I) then we have a normal linear regression

What does the imply for the distribution of Y?

3

Yn Nn eina Y=
$$X_p$$
 +  $E_1$  and combon  $W_n$ 
 $E(Y) = E(X_p + e) = X_p + E(e) = X_p$ 
 $Cov(Y) = Cov(X_p + e) = 0 + Cov(e) = 0 + Cov(e)$ 

The converses X may be regarded as rendom vanishing, and then the assimplians (1)+(2) are made conditional on  $X=\times$ , so  $E(E|X\circ x)=0$  and  $Cov(E|X\circ x)=0$  T

If we leterd assure that

4

Be aware: don't mix errors e (unabserved) with residuals & ("observed")

The residuals will be used to assess model assurptions as proxices for the errors.

#### Model essumption [F3.12]

1) Linearity of coverists Y= Xp+E

If the relationship between Y and X1 IS
ASA here or on to use polymonial (or similar)
IN X1. More advanced; nonperental tragentian

#### 2) Homoscedestic error variance: (ov(E)= o'T

Need to check that Var(c) does not very suptemblic across observations

Look at covener is reciduale -> trend? for out, for in. Solution of problem: If we know Cov(cs)= I

Hodel paremeters, estimates endresiduals

Y= X8+8 , Eco)=0, Cor(0)=0-I

The model parameter as p, or the enhann px1

We will develop abmaton:

\$\beta = (\times T)^{-1} \times T \times \text{by least square and maximum limitions.} \\

\hat{G} = \frac{1}{h^{2}} \left( \gamma - T \beta \right)^{-1} \left( \gamma - T \beta \right) \text{ by resolved } \\

\hat{G} = \frac{1}{h^{2}} \left( \gamma - T \beta \right)^{-1} \left( \gamma - T \beta \right) \text{ by resolved } \\

\hat{G} = \frac{1}{h^{2}} \left( \gamma - T \beta \right)^{-1} \left( \gamma - T \beta \right) \text{ by resolved } \\

\hat{G} = \frac{1}{h^{2}} \left( \gamma - T \beta \right)^{-1} \left( \gamma - T \beta \r

Further: Y is a rendom vector with mean XB, and estimator for E(Y)=XB is Y=XB.

the error  $\epsilon$  is a rendom vector with  $E(\epsilon)$  = 0 and  $Ch(\epsilon)$  =  $6^{\circ}I$ , but  $\epsilon$  is not absenced.

Y= Xx+E randbserred Observed

Our best guess for the error is the residual vector  $(\hat{\epsilon}, \epsilon)$ 

E= Y-9= Y-ZB

So, the residuals can be calculated, and we may think of the residuals as predictions of the errors

2

we may use a socalled general linear model, with veighted [eest squares (lecEx3.P4), but Z is in general unknown.

Remaining: 3) uncorrelated error 4) additive error.

Ŧ

# TMA4267 Linear Statistical Models V2017 (L8)

Part 2: Linear regression:

Modelling the effects of covariates [F:3.1.3] Parameter estimation: Estimator for  $\beta$  [F:3.2.1]

# Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 10, 2017

1/18

# Model assumptions for the classical linear model [F:3.1.2]

What are our model assumptions, how can we spot violations and what can we do to amend the violations

- 1. Linearity of covariates:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
- 2. Homoscedastic error variance:  $Var(\varepsilon_i) = \sigma^2$ .
- 3. Uncorrelated errors:  $Cov(\varepsilon_i, \varepsilon_i) = 0$ .
- 4. Additivity of errors:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

We mainly use plots to assess this (more on model fit in F:3.4 Model choice and variable seletion)

- ► Covariate vs response (for each covariate)
- Covariate vs error (when we have simulated data and know the truth)
- Covariate vs residual (estimated error),
- ► Predicted response vs residual.

2/18

### The classical linear model

The model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

is called a classical linear model if the following is true:

- 1.  $E(\varepsilon) = 0$ .
- 2.  $Cov(\varepsilon) = E(\varepsilon \varepsilon^T) = \sigma^2 I$ .
- 3. The design matrix has full rank rank(X) = k + 1 = p.

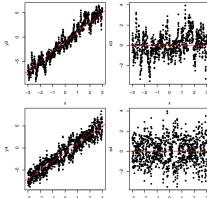
The classical *normal* linear regression model is obtained if additionally

4. 
$$\varepsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$$

holds. For random covariates these assumptions are to be understood conditionally on  $\boldsymbol{X}$ .

1/18

### Uncorrelated errors?



Top: positively autocorrelated errors. Bottom: negatively correlated errors. Right: x vs y. Left: x vs error. Example from Fahrmeir et al (2013): Regression. Springer. (p.81). R code from TMA4267 lectures tab.

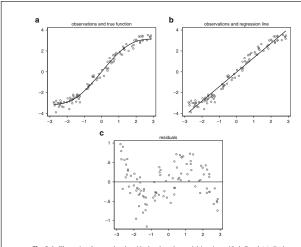
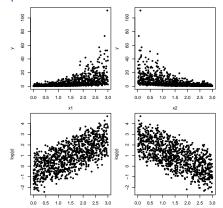


Fig. 3.4 Illustration for correlated residuals when the model is misspecified: Panel (a) displays (simulated) data based on the function  $E(y_1 \mid x_i) = \sin(x_i) + x_i$  and  $\varepsilon_i \sim N(0, 0.3^2)$ . Panel (b) shows the estimated regression line, i.e., the nonlinear relationship is ignored. The corresponding residuals can be found in panel (c)

Fahrmeir et al (2013): Regression. Springer. (p.82)

4/18

# Multiplicative errors



Top: x1 and  $x^{\prime}2$  vs y. Bottom: x1 and x2 vs log(y). Example from Fahrmeir et al (2013): Regression. Springer. (p.85). R code from TMA4267 lectures tab.

6/18

# Multiplicative errors

```
x1=runif(n,0,3)
x2=runif(n,0,3)
e=rnorm(n,0,0.4)
y=exp(1+x1-x2+e)
plot(x1,y,pch=20)
plot(x2,y,pch=20)
plot(x1,log(y),pch=20)
plot(x2,log(y),pch=20)
```

5/18

# Covariates - how to include in the linear regression?

- 1. Continuous covariates: as is, transformed or using polynomials.
- 2. Categorical covariates: dummy variable or effect coding.
- 3. Interactions between covariates.

#### Munich rent index data

```
> colnames(ds)
[1] "rent" "rentsqm" "area" "yearc" "location" "bath"
[7] "kitchen" "cheating" "district"
> apply(ds[,1:4],2,summary)
          rent rentsqm area yearc
Min.
         40.51 0.4158 20.00 1918
1st Qu. 322.00 5.2610 51.00 1939
Median 427.00 6.9800 65.00 1959
        459.40 7.1110 67.37 1956
3rd Qu. 559.40 8.8410 81.00 1972
      1843.00 17.7200 160.00 1997
> unlist(apply(ds[,5:8],2,table))
location.1 location.2 location.3 bath.0 bath.1 kitchen.0
      1794
                1210
                            78
                                     2891 191 2951
 kitchen.1 cheating.0 cheating.1
      131
                 321
                          2761
```

8/18

# Linear coding

```
> fit1=lm(rentsqm~as.numeric(location),data=ds)
> summarv(fit1)
Call:
lm(formula = rentsqm ~ as.numeric(location), data = ds)
Coefficients:
                    Estimate Std. Error t value Pr(>|t|)
(Intercept)
                     6.54390
                             0.12368 52.911 < 2e-16 ***
as.numeric(location) 0.39312 0.08016 4.904 9.88e-07 ***
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 ' 1
Residual standard error: 2.427 on 3080 degrees of freedom
Multiple R-squared: 0.007748, Adjusted R-squared: 0.007425
F-statistic: 24.05 on 1 and 3080 DF, p-value: 9.878e-07
```

10/18

# How to code categorical covariates: rentsqm vs location with linear coding

► Location average=1, good=2 and top=3, and regression model

rentsqm<sub>i</sub> = 
$$\beta_0 + \beta_1 location_i + \varepsilon_i$$

- ▶ Parameter estimate:  $\hat{\beta}_1 = 0.39$ . What does that mean?
  - ▶ Flat of average location:  $\widehat{\text{rentsqm}} = \hat{\beta}_0 + \hat{\beta}_1 \cdot 1$
  - Flat of good location:  $\widehat{\text{rentsqm}} = \hat{\beta}_0 + \hat{\beta}_1 \cdot 2$ Flat of top location:  $\widehat{\text{rentsqm}} = \hat{\beta}_0 + \hat{\beta}_1 \cdot 3$
- ▶ What is the difference in predicted rentsgm between top and good location, and between good and average location?
- ▶ So, the difference between a top and a good location is the same as the difference between good and average. Is this what we want?

9/18

# rentsqm vs location with dummy variable coding

$$\begin{aligned} \mathsf{aloc}_i &= \left\{ \begin{array}{l} 0 & \mathsf{location}_i \ \mathsf{is} \ \mathsf{not} \ \mathsf{average} \\ 1 & \mathsf{location}_i \ \mathsf{is} \ \mathsf{average} \end{array} \right. \\ \mathsf{gloc}_i &= \left\{ \begin{array}{l} 0 & \mathsf{location}_i \ \mathsf{is} \ \mathsf{not} \ \mathsf{good} \\ 1 & \mathsf{location}_i \ \mathsf{is} \ \mathsf{good} \end{array} \right. \\ \mathsf{tloc}_i &= \left\{ \begin{array}{l} 0 & \mathsf{location}_i \ \mathsf{is} \ \mathsf{not} \ \mathsf{top} \\ 1 & \mathsf{location}_i \ \mathsf{is} \ \mathsf{top} \end{array} \right. \end{aligned}$$

rentsqm<sub>i</sub> = 
$$\beta_0 + \beta_1 \operatorname{aloc}_i + \beta_2 \operatorname{gloc}_i + \beta_3 \operatorname{tloc}_i + \varepsilon_i$$

- ▶ Write down the design matrix for this regression model, when we have 1794 flats with average location, 1210 with good and 78 with top location.
- ▶ What is the rank of this design matrix?
- ▶ Is there a problem, and a solution?

#### 3.4 Dummy Coding for Categorical Covariates

For modeling the effect of a covariate  $x \in \{1, ..., c\}$  with c categories using dummy coding, we define the c-1 dummy variables

$$x_{i1} = \begin{cases} 1 & x_i = 1, \\ 0 & \text{otherwise,} \end{cases} \dots x_{i,c-1} = \begin{cases} 1 & x_i = c - 1, \\ 0 & \text{otherwise,} \end{cases}$$

for i = 1, ..., n, and include them as explanatory variables in the regression model

$$y_i = \beta_0 + \beta_1 x_{i1} + \ldots + \beta_{i,c-1} x_{i,c-1} + \ldots + \varepsilon_i.$$

For reasons of identifiability, we omit one of the dummy variables, in this case the dummy variable for category c. This category is called reference category. The estimated effects can be interpreted by direct comparison with the (omitted) reference category.

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.97)

12 / 18

# Effect coding via contr.sum

F-statistic: 13.77 on 2 and 3079 DF. p-value: 1.109e-06

14/18

# Dummy coding via contr.treatment

```
> contrasts(ds$location)=contr.treatment(3)
> fit2=lm(rentsqm~location,data=ds)
> summary(fit2)
Call:
lm(formula = rentsqm ~ location, data = ds)
Coefficients:
           Estimate Std. Error t value Pr(>|t|)
(Intercept) 6.95654
                      0.05728 121.456 < 2e-16 ***
                      0.09025 3.498 0.000475 ***
location2
            0.31570
location3 1.21579
                      0.28060 4.333 1.52e-05 ***
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 ' 1
Residual standard error: 2.426 on 3079 degrees of freedom
```

Multiple R-squared: 0.008867, Adjusted R-squared: 0.008223

F-statistic: 13.77 on 2 and 3079 DF, p-value: 1.109e-06

13 / 18

# Response: birth weight

Covariates: glucose level of mother and BMI of mother.

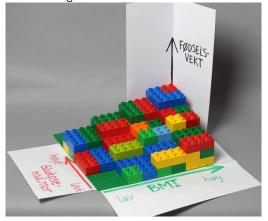


Figure from Kathrine Frey Frøslie.

# Response: birth weight

Covariates: glucose level of mother and BMI of mother - with

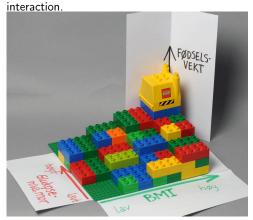


Figure from Kathrine Frey Frøslie.

16/18

# Today

- ► Model assessment: residual plots.
- ► Covariates: how to include in linear regression?
- $\blacktriangleright$  Least squares and maximum likelihood estimator for  $\beta$ .

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

18 / 18

# The classical linear model

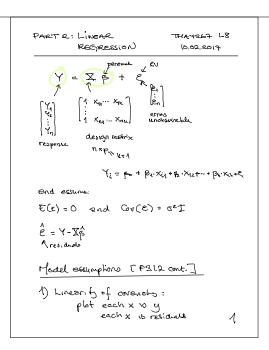
$$\begin{array}{ccc} \mathbf{Y} & = & \mathbf{X} & \boldsymbol{\beta} + \boldsymbol{\varepsilon} \\ (n \times n) & (n \times \rho)(\rho \times 1) & (n \times 1) \end{array}$$

$$E(\varepsilon) = \begin{matrix} \mathbf{0} \\ (n \times 1) \end{matrix} \quad \text{and} \quad Cov(\varepsilon) = \begin{matrix} \sigma^2 \mathbf{I} \\ (n \times n) \end{matrix}$$

#### where

- ightharpoonup and  $\sigma^2$  are unknown parameters and
- ▶ the design matrix X has ith row  $[x_{i1}x_{i2}\cdots x_{ip}]$ .

Next: find the estimator  $\hat{\beta}$ .



- 2) Homoscodestic error ver. enco
- 3) Uncorrelated errors: (or (E;, E)=0

If data comes from newsurements in fine or spece arms may be autocorrelated.

$$\mathcal{E}_{i}^{*} = \mathcal{G} \cdot \mathcal{E}_{i-1} + \mathcal{U}_{i}$$
popular line serve model

But, autocorrelation may also be due to nicipecification of the model, E.g. a musins (unobserved) corenate, or by modelling a linear instead of a nonlinear relationship.

4) Additionly of errors: Y= XB+E

= exp(po).exp(q, x).... exp(p, x). exp(e) has multiplibative or or. .
Transformy the model (Y) using in given additive or or.

2

Problem: renk (X) = 3, not unique solution. How to solve:

- a) not include intrapt
- b) conit one of the during varieties, the anited collegary is called the reference certagory

Ex: let average be omitted: contribustment rentsqui; = po+ps yloci + pr. tloc: + z.

c) add a restriction: sum-to-zero Z pj=0 R. Only. sum

Effect cooling (important in \$24 Set) We have  $X_1 = \begin{cases} 1 & \text{average} \\ 2 & \text{good} \end{cases}$ 

\_

Coveriates: how to include in the linear regression [FS.13]

- 1) Continuous xy
  - \* Y VS X1 linear
  - \* transform xy (lnx, xy, txy)
  - \* use polynomial 11 x1.

dishict

2) Categorical (green, red, blue) ordinal (green, red, blue) (cocher)

b) Dummy reneble coding

Ex: location -> aloc & < 1794

glic & < 1210

See stide there is the first than the control of the co

2 0

$$Z_{1} = \begin{cases} 1 & \text{if } x_{1} = 1 \\ -1 & \text{if } x_{1} = 3 \\ 0 & \text{else } (x_{1} = 2) \end{cases} \qquad Z_{2} = \begin{cases} 1 & \text{if } x_{1} = 2 \\ -1 & \text{if } x_{1} = 3 \\ 0 & \text{else } (x_{1} = 4) \end{cases}$$

$$U_{1} = d_{0} + d_{1} Z_{1} + d_{2} \cdot Z_{2} + C$$

#### 3) Interections

d3= - d1- 02.

Is the effect (on Y) of a change in Xy dependent on the value of another overrele Xz?

Lego: Y= birth waght child

X1: Glucose level of nother

X2: BMI of mother (2)

Figure 1: Y= po+ pyx++ pexe+ &

Figure 2: tigh glucon will have a different effect on brith weight when 611 is low composation when Gri to high.  $\Rightarrow$  we have an interction between  $x_1$  and  $x_2$ .

a) Continuous x end x. te simplest solution: Yi=pon fixint paxit p

meny complex solutions possible

b) categorical: may do the same as for continuous. Easiest solution, define new veneble with all combinetors of x1 and x2

#### Peremeter estimation [F3.2]

#### Eshmetar for B [F3.2.1]

1) Maximum likelihood

If e-Na(0,0°I) then Y~Na(Xp,0°I)

Alt1: Y1, Y2, ", Yn independent

E(Yi)= xi p, Ve (Yi) = 02

5

2) Least squares: minimize LS(3) with LS(p)= (y-XB)T(y-XB)

ii) To minimize LS(p) wit p we may solve

"Need' two rules for donvehires

maximizing L wit & is the seme a minimizing LS(B) wit ps with respect to

Att2: 
$$Y \sim N_{\eta}(\mu, \mathbb{Z})$$
  
 $f(y), \mu, \mathbb{Z}) = (\frac{2\pi}{2\pi})^{\frac{1}{2}} \left[ \operatorname{det}(\mathbb{Z})^{\frac{1}{2}} \right]$   
 $\exp \left\{ -\frac{1}{2} \left( y - \mu\right)^{\frac{1}{2}} \sum_{i=1}^{n} (y \mu_{i})^{\frac{1}{2}} \right\}$   
Homawork:  $\mu = \mathbb{X}_{p}, \ Z = \sigma \in \mathbb{Z} \Rightarrow \operatorname{glek} \text{ the}$   
 $\operatorname{same} \ L(\rho, \sigma^{\epsilon}) \text{ as } \otimes$ 

Honework: 3 LS(p) with these two rules!

8

# TMA4267 Linear Statistical Models V2017 (L9)

Part 2: Linear regression: Parameter estimation [F:3.2]

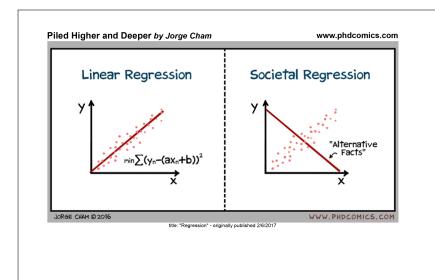
# Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 14, 2017

1/22

2 / 22



The classical linear model

$$F(\varepsilon) = \begin{cases} \mathbf{X} & \mathbf{\beta} + \varepsilon \\ (n \times p)(p \times 1) \end{cases}$$

$$E(\varepsilon) = \begin{cases} \mathbf{0} & \text{and} \quad Cov(\varepsilon) = \sigma^2 \mathbf{I} \\ (n \times n) \end{cases}$$

where

- $\triangleright$   $\beta$  and  $\sigma^2$  are unknown parameters and
- the design matrix **X** has full rank, with *i*th row  $[x_{i1}x_{i2}\cdots x_{ip}]$ .

#### Today

- 1. find estimator for  $\beta$ ,
- 2. find estimator for  $\sigma^2$ , and
- 3. look at two idempotent matrices  $\boldsymbol{H}$  and  $\boldsymbol{I} \boldsymbol{H}$  to arrive at
- 4. geometric interpretation.

1/22

# Rules for derivatives with respect to a vector

- $\blacktriangleright$  Let  $\beta$  be a p-dimensional column vector of interest,
- ▶ and let  $\frac{\partial}{\partial \beta}$  denote the *p*-dimensional vector with partial derivatives wrt the *p* elements of  $\beta$ .
- Let **d** be a p-dimensional column vector of constants and
- ▶ **D** be a  $p \times p$  symmetric matrix of constants.

Rule 1:

$$rac{\partial}{\partialoldsymbol{eta}}(oldsymbol{d}^{ op}oldsymbol{eta}) = rac{\partial}{\partialoldsymbol{eta}}(\sum_{i=1}^{oldsymbol{
ho}}d_{j}eta_{j}) = oldsymbol{d}$$

Rule 2:

$$\frac{\partial}{\partial \boldsymbol{\beta}} (\boldsymbol{\beta}^\mathsf{T} \boldsymbol{D} \boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}} (\sum_{i=1}^p \sum_{k=1}^p \beta_j d_{jk} \beta_k) = 2 \boldsymbol{D} \boldsymbol{\beta}$$

See Härdle and Simes (2015), page 65, Equation (2.23) and (2.24).

# Two questions

Have found least squares and maximum likelihood estimator for  $\beta$ :

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

and we have assumed that the rank( $\boldsymbol{X}$ ) = p for  $n \times p$  design matrix (where n > p).

- ightharpoonup Q1: What can we say about  $X^T X$ ?
- ▶ Q2: Why is the following wrong?

Using  $(AB)^{-1} = B^{-1}A^{-1}$ ,

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

4/22

### Acid rain

occurs when emissions of sulfur dioxide (SO2) and oxides of nitrogen (NOx) 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/

The classical linear model

The model

$$oldsymbol{Y} = oldsymbol{X}oldsymbol{eta} + oldsymbol{arepsilon}$$

is called a classical linear model if the following is true:

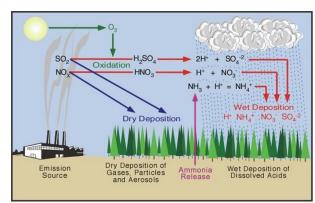
- 1.  $E(\varepsilon) = 0$ .
- 2.  $Cov(\varepsilon) = E(\varepsilon \varepsilon^T) = \sigma^2 I$ .
- 3. The design matrix has full rank rank(X) = k + 1 = p.

The classical *normal* linear regression model is obtained if additionally

4. 
$$\varepsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$$

holds. For random covariates these assumptions are to be understood conditionally on  $\boldsymbol{X}$ .

5 / 22



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

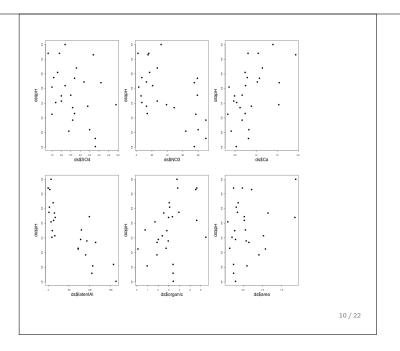
# Acid rain in Norwegian lakes

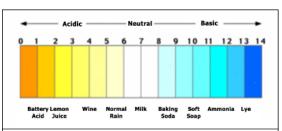
Measured pH in Norwegian lakes explained by content of

- ► x1: SO<sub>4</sub>: sulfate (the salt of sulfuric acid),
- ► x2: N0<sub>3</sub>: nitrate (the conjugate base of nitric acid),
- x3: Ca: calsium,
- ► x4: latent AI: aluminium,
- ▶ x5: organic substance,
- x6: area of lake,
- ► x7: 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.

8 / 22

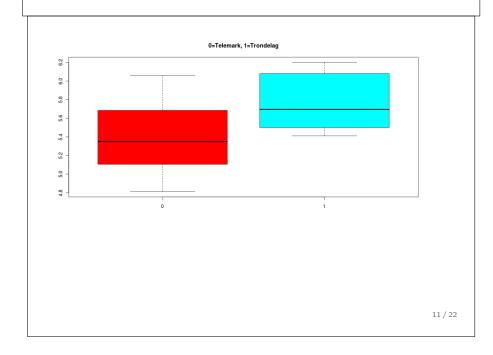


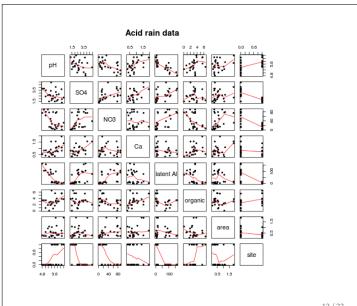


(Source: Physical Geography.net)

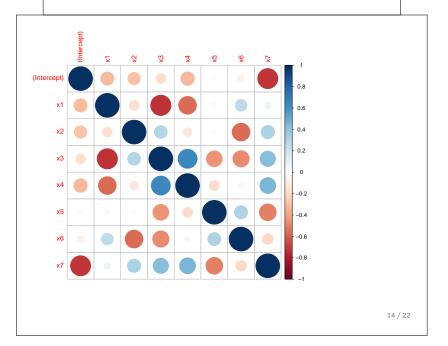
The pH scale: A value of 7.0 is considered neutral. Values higher than 7.0 are increasingly alkaline or basic. Values lower than 7.0 are increasingly acidic.

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





12 / 22



# Output from fitting the full model in R

```
> fit=lm(y~.,data=ds)
> summary(fit)
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 5.6764334 0.1389162 40.862 < 2e-16 ***
           -0.3150444 0.0587512 -5.362 4.27e-05 ***
x2
            -0.0018533 0.0012587 -1.472
                                           0.158
x3
            0.9751745 0.1449075
                                  6.730 2.62e-06 ***
x4
           -0.0002268 0.0010038 -0.226
                                           0.824
           -0.0334242 0.0225009 -1.485
                                           0.155
           -0.0039399 0.0724339
                                 -0.054
                                           0.957
            0.0888722 0.1025724
x7
                                  0.866
                                           0.398
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 ' 1
Residual standard error: 0.1165 on 18 degrees of freedom
Multiple R-squared: 0.93, Adjusted R-squared: 0.9027
```

Question: explain how to interpret  $\hat{\beta}_0$  and  $\hat{\beta}_3$ .

13 / 22

# 3.10 Asymptotic Properties of the Least Squares Estimator

F-statistic: 34.15 on 7 and 18 DF, p-value: 3.904e-09

- 1. The least squares estimator  $\hat{\beta}_n$  for  $\beta$  and the ML or REML estimator  $\hat{\sigma}_n^2$ for the variance  $\sigma^2$  are consistent.
- 2. The least squares estimator asymptotically follows a normal distribution, specifically

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \stackrel{d}{\to} N(\boldsymbol{0}, \sigma^2 V^{-1}).$$

That is the difference  $\hat{\beta}_n - \beta$  normalized with  $\sqrt{n}$  converges in distribution to the normal distribution on the right-hand side.

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.120)

# Projection matrix: definition and properties

▶ A matrix **A** is a *projection matrix* if it is idempotent,  $\mathbf{A}^2 = \mathbf{A}$ .

16/22

# Projection matrix: definition and properties

- A matrix **A** is a *projection matrix* if it is idempotent,  $\mathbf{A}^2 = \mathbf{A}$ .
- An idempotent matrix is an *orthogonal* projection matrix if, in the decomposition of a vector,  $\mathbf{v} = \mathbf{A}\mathbf{v} + (\mathbf{v} \mathbf{A}\mathbf{v})$ ,  $\mathbf{A}\mathbf{v}$  and  $\mathbf{v} \mathbf{A}\mathbf{v} = (\mathbf{I} \mathbf{A})\mathbf{v}$  are always orthogonal, that is,  $(\mathbf{A}\mathbf{v})^T(\mathbf{v} \mathbf{A}\mathbf{v}) = 0$ .
- ► A symmetric projection matrix is orthogonal.

Projection matrix: definition and properties

- A matrix **A** is a projection matrix if it is idempotent,  $\mathbf{A}^2 = \mathbf{A}$ .
- An idempotent matrix is an *orthogonal* projection matrix if, in the decomposition of a vector,  $\mathbf{v} = \mathbf{A}\mathbf{v} + (\mathbf{v} \mathbf{A}\mathbf{v})$ ,  $\mathbf{A}\mathbf{v}$  and  $\mathbf{v} \mathbf{A}\mathbf{v} = (\mathbf{I} \mathbf{A})\mathbf{v}$  are always orthogonal, that is,  $(\mathbf{A}\mathbf{v})^T(\mathbf{v} \mathbf{A}\mathbf{v}) = 0$ .

16 / 22

# Projection matrix: definition and properties

- ▶ A matrix **A** is a *projection matrix* if it is idempotent,  $\mathbf{A}^2 = \mathbf{A}$ .
- An idempotent matrix is an *orthogonal* projection matrix if, in the decomposition of a vector,  $\mathbf{v} = \mathbf{A}\mathbf{v} + (\mathbf{v} \mathbf{A}\mathbf{v})$ ,  $\mathbf{A}\mathbf{v}$  and  $\mathbf{v} \mathbf{A}\mathbf{v} = (\mathbf{I} \mathbf{A})\mathbf{v}$  are always orthogonal, that is,  $(\mathbf{A}\mathbf{v})^T(\mathbf{v} \mathbf{A}\mathbf{v}) = 0$ .
- ► A symmetric projection matrix is orthogonal.
- ▶ The eigenvalues of a projection matrix are 0 and 1.

16 / 22

# Projection matrix: definition and properties

- A matrix **A** is a projection matrix if it is idempotent,  $\mathbf{A}^2 = \mathbf{A}$ .
- An idempotent matrix is an *orthogonal* projection matrix if, in the decomposition of a vector,  $\mathbf{v} = A\mathbf{v} + (\mathbf{v} A\mathbf{v})$ ,  $A\mathbf{v}$  and  $\mathbf{v} A\mathbf{v} = (\mathbf{I} \mathbf{A})\mathbf{v}$  are always orthogonal, that is,  $(A\mathbf{v})^T(\mathbf{v} A\mathbf{v}) = 0$ .
- ▶ A symmetric projection matrix is orthogonal.
- ▶ The eigenvalues of a projection matrix are 0 and 1.
- ▶ If a  $(n \times n)$  symmetric projection matrix **A** has rank r then r eigenvalues are 1 and n r are 0.

16/22

#### Results so far

Least squares and maximum likelihood estimator for β:

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

▶ Restricted maximum likelihood estimator for  $\sigma^2$ :

$$\hat{\sigma^2} = \frac{1}{n-p} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \frac{\mathsf{SSE}}{n-p}$$

▶ Projection matrices: idempotent, symmetric/orthogonal:

$$H = X(X^TX)^{-1}X^T$$
$$I - H = I - X(X^TX)^{-1}X^T$$

with important connection:

$$\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$$
 $\hat{\mathbf{\varepsilon}} = \mathbf{I} - \mathbf{H}\mathbf{Y}$ 

17 / 22

# Projection matrix: definition and properties

- A matrix **A** is a projection matrix if it is idempotent,  $\mathbf{A}^2 = \mathbf{A}$ .
- An idempotent matrix is an *orthogonal* projection matrix if, in the decomposition of a vector,  $\mathbf{v} = \mathbf{A}\mathbf{v} + (\mathbf{v} \mathbf{A}\mathbf{v})$ ,  $\mathbf{A}\mathbf{v}$  and  $\mathbf{v} \mathbf{A}\mathbf{v} = (\mathbf{I} \mathbf{A})\mathbf{v}$  are always orthogonal, that is,  $(\mathbf{A}\mathbf{v})^T(\mathbf{v} \mathbf{A}\mathbf{v}) = 0$ .
- A symmetric projection matrix is orthogonal.
- ▶ The eigenvalues of a projection matrix are 0 and 1.
- ▶ If a  $(n \times n)$  symmetric projection matrix **A** has rank r then r eigenvalues are 1 and n r are 0.
- The trace and rank of a symmetric projection matrix are equal: tr(A) = rank(A).

16 / 22

### Results from Mathematics 3

#### Best approximation theorem

The vector  $\hat{\mathbf{Y}}$  in the column space of  $\mathbf{X}$  that makes  $||\mathbf{Y} - \hat{\mathbf{Y}}||$  as small as possible, is the orthogonal projection of  $\mathbf{Y}$  on the column space of  $\mathbf{X}$ .

#### Results from Mathematics 3

#### Best approximation theorem

The vector  $\hat{\mathbf{Y}}$  in the column space of  $\mathbf{X}$  that makes  $||\mathbf{Y} - \hat{\mathbf{Y}}||$  as small as possible, is the orthogonal projection of  $\mathbf{Y}$  on the column space of  $\mathbf{X}$ .

#### Orthogonal decomposition

We want  $\hat{\boldsymbol{\beta}}$  to minimize  $||\boldsymbol{Y} - \hat{\boldsymbol{Y}}|| = (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})$  (least squares principle).

18 / 22

### Results from Mathematics 3

#### Best approximation theorem

The vector  $\hat{\mathbf{Y}}$  in the column space of  $\mathbf{X}$  that makes  $||\mathbf{Y} - \hat{\mathbf{Y}}||$  as small as possible, is the orthogonal projection of  $\mathbf{Y}$  on the column space of  $\mathbf{X}$ .

#### Orthogonal decomposition

We want  $\hat{\boldsymbol{\beta}}$  to minimize  $||\boldsymbol{Y} - \hat{\boldsymbol{Y}}|| = (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})$  (least squares principle).

The column space of X consists of vectors of the form  $X\hat{\beta}$ , so  $X\hat{\beta}$  is the orthogonal projection of Y onto the column space of X.

This is equivalent to observing that  $\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  is in the orthogonal complement of the column space of  $\mathbf{X}$ .

18 / 22

#### Results from Mathematics 3

#### Best approximation theorem

The vector  $\hat{\mathbf{Y}}$  in the column space of  $\mathbf{X}$  that makes  $||\mathbf{Y} - \hat{\mathbf{Y}}||$  as small as possible, is the orthogonal projection of  $\mathbf{Y}$  on the column space of  $\mathbf{X}$ .

#### Orthogonal decomposition

We want  $\hat{\boldsymbol{\beta}}$  to minimize  $||\boldsymbol{Y} - \hat{\boldsymbol{Y}}|| = (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})$  (least squares principle).

The column space of X consists of vectors of the form  $X\hat{\beta}$ , so  $X\hat{\beta}$  is the orthogonal projection of Y onto the column space of X.

18 / 22

### Results from Mathematics 3

#### Best approximation theorem

The vector  $\hat{\mathbf{Y}}$  in the column space of  $\mathbf{X}$  that makes  $||\mathbf{Y} - \hat{\mathbf{Y}}||$  as small as possible, is the orthogonal projection of  $\mathbf{Y}$  on the column space of  $\mathbf{X}$ .

#### Orthogonal decomposition

We want  $\hat{\boldsymbol{\beta}}$  to minimize  $||\boldsymbol{Y} - \hat{\boldsymbol{Y}}|| = (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})$  (least squares principle).

The column space of  $\boldsymbol{X}$  consists of vectors of the form  $\boldsymbol{X}\hat{\boldsymbol{\beta}}$ , so  $\boldsymbol{X}\hat{\boldsymbol{\beta}}$  is the orthogonal projection of  $\boldsymbol{Y}$  onto the column space of  $\boldsymbol{X}$ .

This is equivalent to observing that  $\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  is in the orthogonal complement of the column space of  $\mathbf{X}$ .

That is,  $\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  is orthogonal to all columns of  $\mathbf{X}$ , so  $\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{0}$  and  $\mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{Y}$ .

#### Results from Mathematics 3

#### Best approximation theorem

The vector  $\hat{\mathbf{Y}}$  in the column space of  $\mathbf{X}$  that makes  $||\mathbf{Y} - \hat{\mathbf{Y}}||$  as small as possible, is the orthogonal projection of  $\mathbf{Y}$  on the column space of  $\mathbf{X}$ .

#### Orthogonal decomposition

We want  $\hat{\boldsymbol{\beta}}$  to minimize  $||\boldsymbol{Y} - \hat{\boldsymbol{Y}}|| = (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})$  (least squares principle).

The column space of X consists of vectors of the form  $X\hat{\beta}$ , so  $X\hat{\beta}$  is the orthogonal projection of Y onto the column space of X.

This is equivalent to observing that  $Y - X\hat{\beta}$  is in the orthogonal complement of the column space of X.

That is,  $\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  is orthogonal to all columns of  $\mathbf{X}$ , so  $\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{0}$  and  $\mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{Y}$ .

18 / 22

### Results from Mathematics 3

#### Best approximation theorem

The vector  $\hat{\mathbf{Y}}$  in the column space of  $\mathbf{X}$  that makes  $||\mathbf{Y} - \hat{\mathbf{Y}}||$  as small as possible, is the orthogonal projection of  $\mathbf{Y}$  on the column space of  $\mathbf{X}$ .

# Orthogonal decomposition

We want  $\hat{\boldsymbol{\beta}}$  to minimize  $||\boldsymbol{Y} - \hat{\boldsymbol{Y}}|| = (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})$  (least squares principle).

The column space of  $\boldsymbol{X}$  consists of vectors of the form  $\boldsymbol{X}\boldsymbol{\hat{\beta}}$ , so  $\boldsymbol{X}\boldsymbol{\hat{\beta}}$  is the orthogonal projection of  $\boldsymbol{Y}$  onto the column space of  $\boldsymbol{X}$ .

This is equivalent to observing that  $Y - X\hat{\beta}$  is in the orthogonal complement of the column space of X.

 $\hat{\varepsilon} = Y - HY = (I - H)Y$ , and I - H projects onto the space orthogonal to the column space of X. Observe: (I-H)X=0

That is,  $\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  is orthogonal to all columns of  $\mathbf{X}$ , so  $\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{0}$  and  $\mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{Y}$ .

#### Results from Mathematics 3

#### Best approximation theorem

The vector  $\hat{\mathbf{Y}}$  in the column space of  $\mathbf{X}$  that makes  $||\mathbf{Y} - \hat{\mathbf{Y}}||$  as small as possible, is the orthogonal projection of  $\mathbf{Y}$  on the column space of  $\mathbf{X}$ .

#### Orthogonal decomposition

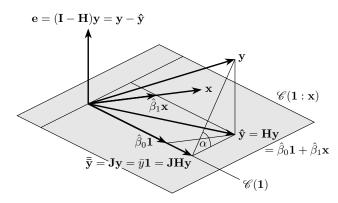
We want  $\hat{\boldsymbol{\beta}}$  to minimize  $||\boldsymbol{Y} - \hat{\boldsymbol{Y}}|| = (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})$  (least squares principle).

The column space of X consists of vectors of the form  $X\hat{\beta}$ , so  $X\hat{\beta}$  is the orthogonal projection of Y onto the column space of  $X.\hat{Y} = HY$ , and  $H = X(X^TX)^{-1}X^T$  projects onto the column space of X. Observe: HX = X.

This is equivalent to observing that  $\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  is in the orthogonal complement of the column space of  $\mathbf{X}$ .

That is,  $\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  is orthogonal to all columns of  $\mathbf{X}$ , so  $\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{0}$  and  $\mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{Y}$ .

18 / 22



Putanen, Styan and Isotalo: Matrix Tricks for Linear Statistical Models: Our Personal Top Twenty, Figure 8.3.

#### 3.7 Geometric Properties of the Least Squares Estimator

The method of least squares has the following geometric properties:

- 1. The predicted values  $\hat{y}$  are orthogonal to the residuals  $\hat{\epsilon}$ , i.e.,  $\hat{y}'\hat{\epsilon} = 0$ .
- 2. The columns  $x^j$  of X are orthogonal to the residuals  $\hat{\epsilon}$ , i.e.,  $(x^j)'\hat{\epsilon} = 0$  or  $X'\hat{\epsilon} = 0$
- 3. The average of the residuals is zero, i.e.,

$$\sum_{i=1}^{n} \hat{\varepsilon}_i = 0 \quad \text{or} \quad \frac{1}{n} \sum_{i=1}^{n} \hat{\varepsilon}_i = 0.$$

4. The average of the predicted values  $\hat{y}_i$  is equal to the average of the observed response  $y_i$ , i.e.,

$$\frac{1}{n}\sum_{i=1}^n \hat{y}_i = \bar{y}.$$

5. The regression hyperplane runs through the average of the data, i.e.,

$$\bar{v} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}_1 + \dots + \hat{\beta}_k \bar{x}_k.$$

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.112)

20 / 22

# Today

▶ Least squares and maximum likelihood estimator for  $\beta$ :

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

has mean  $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$  and  $Cov(\hat{\boldsymbol{\beta}}) = \sigma^2(\boldsymbol{X}^T\boldsymbol{X})^{-1}$ .

- ▶ For the normal model:  $\hat{\boldsymbol{\beta}} \sim N_n(\boldsymbol{\beta}, \sigma^2(\boldsymbol{X}^T\boldsymbol{X})^{-1})$ .
- Asymptotic properties of the least squares estimator: normality.
- ► Orthogonal projection matrices **H** and **I** − **H** with geometric interpretation.

Next time: properties of residuals and  $\hat{\sigma}^2$ , confidence intervals and hypothesis testing for regression coefficients.

22 / 22

# Alternative summery of Geometry of Least Squares

- ▶ Mean response vector:  $E(Y) = X\beta$
- As β varies, Xβ spans the model plane of all linear combinations. I.e. the space spanned by the columns of X: the column-space of X.
- ▶ Due to random error (and unobserved covariates), Y is not exactly a linear combination of the columns of X.
- LS-estimation chooses  $\hat{\beta}$  such that  $X\hat{\beta}$  is the point in the column-space of X that is closes to Y.
- ► The residual vector  $\hat{\boldsymbol{\varepsilon}} = \boldsymbol{Y} \hat{\boldsymbol{Y}} = (\boldsymbol{I} \boldsymbol{H})\boldsymbol{Y}$  is perpendicular to the column-space of  $\boldsymbol{X}$ .
- ▶ Multiplication by  $\boldsymbol{H} = \boldsymbol{X}(\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T$  projects a vector onto the column-space of  $\boldsymbol{X}$ .
- ▶ Multiplication by  $I H = I X(X^TX)^{-1}X^T$  projects a vector onto the space perpendicular to the column-space of X.

21 / 22

iii) solving the normal equations (substitutes

io) mun or max

= 2 XTX

If this matrix has only positive ergenvalue this will be the minimum.

#### Questions on slide:

Stenup: XTX pxp notinx

symmetric

positive definite

whose exists

If UT XXXV > 0 for all U40 than XXX possibive definite.

UXXv= (Xv) (Xv) ≥ 0

2

#### Ex: Acid rain: fit in R wing Im

βο= Estimate Talerapt = 5.67
= estimate of the pH in a lake when

x=0, x=0, , x=0.

A= 0.975

I if a necess by one wit, and all the other covernatione hapt content, then we predict that plf will increase by Bs = 9743.

St. 50(8)

4

Assume that UTETEV = 0 than XU=0. If X has full rent than XU=0 only has U=0 or solution.

Than XTE must be positive definite.

B=(XTX) BN is the less squares estrates of p. If we assume a normal linear needs than \$ 15 also the maximum likelihood estrater

3

5

### Properties of &

$$\beta = \frac{(X \cdot X)^{-1} X^{-1} Y}{C} \quad \text{and} \quad \delta(Y) = X p$$

$$C \quad (CV) = \sigma^{2} D$$

In a normal model: for Np (p, or (XTX)-1)

Last infunction on 
$$\hat{\beta}$$
: From part 1:  

$$(\hat{\beta} - \Xi(\hat{\beta}))^{\dagger} \text{ Gav}(\hat{\beta})^{\dagger}(\hat{\beta} - \Xi(\hat{\beta})) \sim \mathcal{K}_{p}$$

$$\vec{\sigma}^{\dagger} (\hat{\beta} - \beta)^{\dagger} (X^{\dagger}X) (\hat{\beta} - \beta) \sim \mathcal{K}_{p}$$

$$Co(\frac{1}{2})_{-1} = \frac{0}{7}(X_1X)_{-1}$$

$$Co(\frac{1}{2})_{-1} = 0, (X_1X)_{-1}$$

#### Estimator for CF322]

Æ

#### Predicted values and residuals [F3.2.1]

$$E(Y) = X_{\beta}$$
, so  $E(Y) = X_{\beta}^{2} = \hat{Y} \leftarrow \text{prediction}$   
 $\hat{\beta} = (x_{1}x_{2}^{2})^{2}Y$ 

$$\lambda = \chi v = \overline{\chi(\chi_L \chi_{J-1} \chi_L)} = H\lambda$$

H=  $X(X^TX)^TX^T$  is collect the "hat matrix" for pulting the linet on Y.

Residuals:  $\hat{\varepsilon} = Y - \hat{Y} = Y - HY = (I - H)Y$ IY

Obsence (see also Rector 3.9%) that

H is symmetric

1s idempolent H<sup>2</sup>=H

has Both P => Show this

(I-H) is also symmetric and idempotenty and ond rank (I-H) = n-p. => show this

⇒ Work with Rec. Ex 3. P3 - and be reacy for L10! supervision Thurs. 16:15 at Smla.

8

# TMA4267 Linear Statistical Models V2017 (L10)

Part 2: Linear regression: Parameter estimation [F:3.2], Properties of residuals and distribution of estimator for error variance Confidence interval and hypothesis for one regression coefficient

# Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 17, 2017

### Today

- 1. Properties for residuals (from the hat matrix), leading to properties for  $\hat{\sigma}^2$ ,
- 2. Then, confidence interval and hypothesis test for regression coefficient

1/17

### Results so far

Least squares and maximum likelihood estimator for β:

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

with mean  $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$  and  $Cov(\hat{\boldsymbol{\beta}}) = \sigma^2(\boldsymbol{X}^T\boldsymbol{X})^{-1}$ .

▶ Restricted maximum likelihood estimator for  $\sigma^2$ :

$$\hat{\sigma^2} = \frac{1}{n-p} (\mathbf{Y} - \mathbf{X}\hat{\beta})^T (\mathbf{Y} - \mathbf{X}\hat{\beta}) = \frac{\mathsf{SSE}}{n-p}$$

▶ Projection matrices: idempotent, symmetric/orthogonal:

$$\boldsymbol{H} = \boldsymbol{X} (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T$$

projects onto column space of X

$$I - H = I - X(X^TX)^{-1}X^T$$

projects onto space orthogonal to column space of  $\boldsymbol{X}$ 

with important connection: predictions  $\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$  and residuals  $\hat{\boldsymbol{\varepsilon}} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$ 

3 / 17

### The classical linear model

The model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

is called a classical linear model if the following is true:

1.  $E(\varepsilon) = 0$ .

2. 
$$Cov(\varepsilon) = E(\varepsilon \varepsilon^T) = \sigma^2 I$$
.

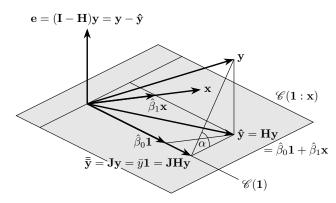
3. The design matrix has full rank rank( $\boldsymbol{X}$ ) = k+1=p.

The classical *normal* linear regression model is obtained if additionally

1. 
$$\varepsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$$

holds. For random covariates these assumptions are to be understood conditionally on  $\boldsymbol{X}$ .

2/17



Putanen, Styan and Isotalo: Matrix Tricks for Linear Statistical Models: Our Personal Top Twenty, Figure 8.3.

# Quadratic forms [F:B3.3, Theorem B.2]

Random vector  $\boldsymbol{X}$  with mean  $\mu$  and covariance matrix  $\boldsymbol{\Sigma}$ , symmetric constant matrix  $\boldsymbol{A}$ .

- Quadratic form: X<sup>T</sup>AX.
- ► The "trace-formula":  $E(X^TAX) = tr(AΣ) + μ^TAμ$ .

Then, let  $X \sim N_p(0, I)$ , and R is a symmetric and idempotent matrix with rank r.

$$\boldsymbol{X}^T \boldsymbol{R} \boldsymbol{X} \sim \chi_r^2$$

Now, also  $\boldsymbol{S}$  is a symmetric and idempotent matrix with rank  $\boldsymbol{s}$ , and  $\boldsymbol{RS} = \boldsymbol{0}$ .

$$\frac{s\boldsymbol{X}^T\boldsymbol{R}\boldsymbol{X}}{r\boldsymbol{X}^T\boldsymbol{S}\boldsymbol{X}}\sim F_{r,s}$$

5/17

# Acid rain in Norwegian lakes

Measured pH in Norwegian lakes explained by content of

- ▶ x1: SO<sub>4</sub>: sulfate (the salt of sulfuric acid),
- ▶ x2: N0<sub>3</sub>: nitrate (the conjugate base of nitric acid),
- x3: Ca: calsium.
- ► x4: latent AI: aluminium,
- x5: organic substance,
- x6: area of lake.
- ▶ x7: position of lake (Telemark or Trøndelag),

Random sample of n = 26 lakes.

7 / 17

# Properties: $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2$

Least squares and maximum likelihood estimator for *β*:

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

has mean  $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$  and  $Cov(\hat{\boldsymbol{\beta}}) = \sigma^2(\boldsymbol{X}^T\boldsymbol{X})^{-1}$ .

- ▶ In addition  $\hat{\beta}$  is best linear unbiased estimator (BLUE), that is, among all unbiased estimator it has minimum variance in each component. (More in TMA4295 Statistical Inference.)
- ▶ For the normal model:  $\hat{\boldsymbol{\beta}} \sim N_p(\boldsymbol{\beta}, \sigma^2(\boldsymbol{X}^T\boldsymbol{X})^{-1})$ .
- ▶ Restricted maximum likelihood estimator for  $\sigma^2$ :

$$\hat{\sigma^2} = \frac{1}{n-p} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \frac{\mathsf{SSE}}{n-p}$$

► For the normal model

> fit=lm(v~.,data=ds)

$$\frac{(n-p)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p}^2$$

6 / 17

# Output from fitting the full model in ${\sf R}$

```
> summary(fit)
Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 5.6764334 0.1389162 40.862 < 2e-16 ***
           -0.3150444 0.0587512 -5.362 4.27e-05 ***
x2
           -0.0018533 0.0012587 -1.472
x3
            0.9751745 0.1449075 6.730 2.62e-06 ***
           -0.0002268 0.0010038 -0.226
x5
           -0.0334242 0.0225009 -1.485
                                          0.155
           -0.0039399 0.0724339 -0.054
                                          0.957
x7
            0.0888722 0.1025724
                                0.866
                                          0.398
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 ' 1
```

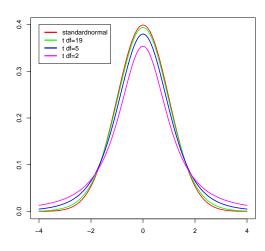
Residual standard error: 0.1165 on 18 degrees of freedom Multiple R-squared: 0.93,Adjusted R-squared: 0.9027 F-statistic: 34.15 on 7 and 18 DF, p-value: 3.904e-09

#### W. S. Gosset alias Student



9 / 17

### *t*-distribution



11 / 17

# Historisk: Student-t fordelingen

- W.S. Gosset (1876-1937) was employed by the Guinness Brewing Company of Dublin.
- Sample sizes available for experimentation in brewing were necessarily small, and Gosset knew that a correct way of dealing with small samples was needed.
- He consulted Karl Pearson (1857-1936) of University College in London about the problem. Pearson told him the current state of knowledge was unsatisfactory.
- ► The following year Gosset undertook a course of study under Pearson. An outcome of his study was the publication in 1908 of Gosset's paper on "The Probable Error of a Mean," which introduced a form of what later became known as Student's t-distribution.
- ► Gosset's paper was published under the pseudonym "Student."
- The modern form of Student's t-distribution was derived by R.A. Fisher and first published in 1925.

10 / 17

### DEF: t-distribution

Let Z be a standard normal random variable and V a chi-squared random variable with parameter  $\nu$  (degrees of freedom). If Z and V are independent, the distribution of the random variable T

$$T = \frac{Z}{\sqrt{V/\nu}}$$

has probability density function

$$h(t) = \frac{\Gamma[(\nu+1)/2]}{\Gamma(\nu/2)\sqrt{\pi\nu}} (1 + \frac{t^2}{\nu})^{-(\nu+1)/2}$$

for  $-\infty < t < \infty$ . This distribution is called the (Student) t-distribution with  $\nu$  degrees of freedom.

- ▶  $E(T) = 0 \text{ if } \nu \ge 2.$

# Are $\hat{\beta}$ and SSE are independent?

Independence – from Part 1:

Let  $X_{(p\times 1)}$  be a random vector from  $N_p(\mu, \Sigma)$ . Then AX and BX are independent iff  $A\Sigma B^T = 0$ .

#### We have:

- $ightharpoonup Y \sim N_n(X\beta, \sigma^2 I)$
- $ightharpoonup AY = \hat{eta} = (X^TX)^{-1}X^TY$ , and
- BY = (I H)Y.
- Now  $\mathbf{A}\sigma^2 \mathbf{I} \mathbf{B}^T = \sigma^2 \mathbf{A} \mathbf{B}^T = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{I} \mathbf{H}) = \mathbf{0}$
- ► since X(I H) = X HX = X X = 0.
- We conclude that  $\hat{\boldsymbol{\beta}}$  is independent of  $(\boldsymbol{I} \boldsymbol{H})\boldsymbol{Y}$ ,
- ▶ and, since SSE=function of (I H)Y: SSE= $Y^T(I H)Y$ ,
- **b** then  $\hat{\beta}$  and SSE are independent.

13 / 17

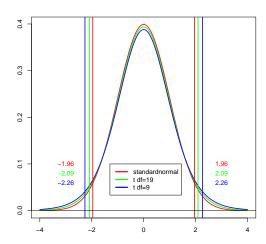
#### Kritiske verdier i t-fordelingen

 $P(T > t_{\alpha,\nu}) = \alpha$ 

$\nu \backslash \alpha$	.150	.100	.075	.050	.025	.010	.005	.001	.0005
- 1	1.963	3.078	4.165	6.314	12.706	31.821	63.657	318.309	636.619
2	1.386	1.886	2.282	2.920	4.303	6.965	9.925	22.327	31.599
3	1.250	1.638	1.924	2.353	3.182	4.541	5.841	10.215	12.924
4	1.190	1.533	1.778	2.132	2.776	3.747	4.604	7.173	8.610
5	1.156	1.476	1.699	2.015	2.571	3.365	4.032	5.893	6.869
6	1.134	1.440	1.650	1.943	2.447	3.143	3.707	5.208	5.959
7	1.119	1.415	1.617	1.895	2.365	2.998	3.499	4.785	5.408
8	1.108	1.397	1.592	1.860	2.306	2.896	3.355	4.501	5.041
9	1.100	1.383	1.574	1.833	2.262	2.821	3.250	4.297	4.781
10	1.093	1.372	1.559	1.812	2.228	2.764	3.169	4.144	4.587
11	1.088	1.363	1.548	1.796	2.201	2.718	3.106	4.025	4.437
12	1.083	1.356	1.538	1.782	2.179	2.681	3.055	3.930	4.318
13	1.079	1.350	1.530	1.771	2.160	2.650	3.012	3.852	4.221
14	1.076	1.345	1.523	1.761	2.145	2.624	2.977	3.787	4.140
15	1.074	1.341	1.517	1.753	2.131	2.602	2.947	3.733	4.073
16	1.071	1.337	1.512	1.746	2.120	2.583	2.921	3.686	4.015
17	1.069	1.333	1.508	1.740	2.110	2.567	2.898	3.646	3.965
18	1.067	1.330	1.504	1.734	2.101	2.552	2.878	3.610	3.922
19	1.066	1.328	1.500	1.729	2.093	2.539	2.861	3.579	3.883
20	1.064	1.325	1.497	1.725	2.086	2.528	2.845	3.552	3.850
21	1.063	1.323	1.494	1.721	2.080	2.518	2.831	3.527	3.819
22	1.061	1.321	1.492	1.717	2.074	2.508	2.819	3.505	3.792
23	1.060	1.319	1.489	1.714	2.069	2.500	2.807	3.485	3.768
24	1.059	1.318	1.487	1.711	2.064	2.492	2.797	3.467	3.745
25	1.058	1.316	1.485	1.708	2.060	2.485	2.787	3.450	3.725
26	1.058	1.315	1.483	1.706	2.056	2.479	2.779	3.435	3.707
27	1.057	1.314	1.482	1.703	2.052	2.473	2.771	3.421	3.690
28	1.056	1.313	1.480	1.701	2.048	2.467	2.763	3.408	3.674
29	1.055	1.311	1.479	1.699	2.045	2.462	2.756	3.396	3.659
30	1.055	1.310	1.477	1.697	2.042	2.457	2.750	3.385	3.646
35	1.052	1.306	1.472	1.690	2.030	2.438	2.724	3.340	3.591
40	1.050	1.303	1.468	1.684	2.021	2.423	2.704	3.307	3.551
50	1.047	1.299	1.462	1.676	2.009	2.403	2.678	3.261	3.496
60	1.045	1.296	1.458	1.671	2.000	2.390	2.660	3.232	3.460
80	1.043	1.292	1.453	1.664	1.990	2.374	2.639	3.195	3.416
100	1.042	1.290	1.451	1.660	1.984	2.364	2.626	3.174	3.390
120	1.041	1.289	1.449	1.658	1.980	2.358	2.617	3.160	3.373
$\infty$	1.036	1.282	1.440	1.645	1.960	2.326	2.576	3.090	3.291

15 / 17

# Quantiles and critical values: N og t: $\alpha/2 = 0.025$



14 / 17

# Acid rain in R

```
\label{lem:ds} $$ds=read.table("https://www.math.ntnu.no/emner/TMA4267/2017v/acidrain.txt",header=TRUE)$$fit=lm(y~.,data=ds)
```

> confint(fit)

	/		
	2.5 %	97.5 %	
(Intercept)	5.384581378	5.9682854281	
x1	-0.438476153	-0.1916126966	
x2	-0.004497716	0.0007911594	
x3	0.670735075	1.2796138706	
x4	-0.002335625	0.0018820903	
x5	-0.080696921	0.0138484550	
x6	-0.156117992	0.1482381575	
x7	-0.126624544	0.3043688780	

P-values: http://www.statistrikk.no/wp-content/uploads/2017/02/nerdekort.jpg

### Today

- ▶ Distribution of SSE/ $\sigma^2$  is chisquared (n-p).
- ▶ Independence of  $\hat{\beta}$  and SSE.
- Inference about β components can be performed using the t-distribution

17 / 17

Drishibution of SSE and 
$$\hat{\sigma}^2$$

SSE =  $\hat{\mathcal{E}} \uparrow \hat{\mathcal{E}} = Y^T (I-H)(I-H)Y$ 
 $\hat{\mathcal{I}} (y_i - \hat{\mathcal{G}})^e$ 
 $\hat{\mathcal{E}} = (I-H)Y$ 

SSE =  $Y^T (I-H)Y$ 

Temanber Fenn  $(I-H) = v_1 - p$ .

RecEx3.83 boloo at the distribution of  $\hat{\mathcal{E}} Y^T (I-H)Y = \hat{\mathcal{O}}^2$ 
by using the result on questric form from the A

(see state)

 $Y \sim N_n(X_i^n, \sigma^* I)$ 
 $Y^* = \hat{\mathcal{O}} (Y-X_i^n) \sim N_n(0, I)$ 

OIL FORFAIT PART 2: LINEAR REGRESSION 14.02.2017 [FB2] Distribution of & (residuals) Residualo: ê= Y-Y=Y-Xp=Y-Xp=Y-XXXXY Y(H-I) = YH -Y =  $E(\hat{\mathcal{E}}) = E(CI-H)Y = (I-H)E(Y)$ Y= XB+E = (I-H) XB = 0 project and the spece orth. to column open of the or (I-H) XB = (X - HX) = 0  $Cov(\hat{\mathcal{E}}) = Gv(fT-H)Y) = (T-H)Gv(Y)(T-H)^T$ = 02 ( I-H) I (I-H) = 02 (I-H) ABSUME C~ Nn(0, 02 I) => E~Nn (0, 02 (I-H)) 16: renk(H)=p, renk(I-H)= n-p, which meens (TH) - does not exist end we use the singular version of the normal poly.

$$E(V) = n-p \quad Ver(V) \cdot 2(n-p)$$

$$E(\hat{\sigma}^2) = E(\frac{1}{n-p}, SSE) = \frac{1}{n-p} E(\sigma^2V)$$

$$= \frac{\sigma^2}{n-p} \frac{E(V)}{n-p} = \underline{\sigma}^2 \quad \text{outbrased.}$$
This is true when we assume  $c \sim V$ .

If we do not assume  $e \sim V$ , then we cen use the freee-framely
$$E(SSE) = E(Y^*(I-H)Y) \quad E(Y) = K_B$$

$$= (n-p)\sigma^2 + 0$$

$$E(\delta^2) - E(\frac{SSE}{n-p}) = \underline{\sigma}^2$$

3

#### Inference about one By

Ex: Acid coin Bi-effect of soy on pH of lake

$$SD(\hat{\beta}_1) = \sqrt{6^2 \left[ \left( \frac{1}{2} \left( \frac{1}{2} \right)^4 \right]} \left[ Colony, ho Soy \right]$$

$$SD(\hat{\beta}_1) = \sqrt{6^2 \left[ \left( \frac{1}{2} \left( \frac{1}{2} \right)^4 \right]} \right]} Collect, ho Soy 1$$

1820.0 € rong.18

6: "Residual shenderderror" = 0.1165 n=26 3 n-p= 68

on 18 degrees of beedom".

To find a confidence inhead for B', - or to test hypotheses about p; we need to know the distribution of a stehslic involving fij end pjwith no other lunknown peremetes.

We have: 
$$\frac{\beta_j - \beta_j}{\sqrt{Q_{ij}^{\gamma}}} \sim N(0,1)$$
and  $\frac{(n-p)\delta^2}{\delta^2} \sim \chi_{n-p}^2$ 

B, and ô2 need to be independent incrow that this holds. => Rec & 3 P3 + sides

Use 
$$T_j = \frac{\hat{\beta} - \hat{\beta}}{\sqrt{c_j}} \sim t_{mp}$$
 for infrares.

a) Find a 95% confidence injurial (CI)

b

and  $\sigma^2 = \frac{SSE}{n-p}$  when  $\frac{SSE}{\sigma^2} \wedge \chi^2_{n-p}$ .

 $SD(\hat{\beta}_j) = \overline{C_{jj}} \cdot \hat{G}$  and  $\hat{\beta}_j$  and  $\hat{G}_j$  ere independent (to be shown)

General result:

pervel result: 
$$\frac{N(o_i 1)}{\sqrt{\frac{X_q^i}{q_i^4}}} \sim t_q$$

5

b) Test:

Ex: Acid rain por effect of Soy on pH.

$$P(-t_{\underline{x},n-p} < T_j < t_{\underline{x}_j,n-p}) = 1-\alpha$$

How do you interpret this interval?

- strong belief (75%) that p is in Imenal:

We see that 0 is not in the intend-what does this been? -> Reject Ho.g.=0 us thing +0 at sign Lind 5%

Ę

# Today

- 1. Hypothesis testing for  $\beta_i$ .
- 2. Residuals: standardized (or studentized) preferred.
- 3. Decomposition of variability: SST=SSR+SSE, and significance of regression.
- 4.  $R^2$  gives the proportion of variability explained by the regression model. and will never decrease if new covariates are added to the model.
- 5. Model choice considerations.
- 6. SPSE: Expected squared prediction error.

TMA4267 Linear Statistical Models V2017 (L11)

Part 2: Linear regression:

Parameter estimation [F:3.2] and model selection [F:3.4]
Hypothesis test for one regression coefficient
Studentized and standardized residuals
decomposition of variability and signficance of regression
R<sup>2</sup>, SPSE=Expected squared prediction error

# Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 21, 2017

1/30

## The classical linear model

The model

$$Y = X\beta + \varepsilon$$

is called a classical linear model if the following is true:

- 1.  $E(\varepsilon) = 0$ .
- 2.  $Cov(\varepsilon) = E(\varepsilon \varepsilon^T) = \sigma^2 I$ .
- 3. The design matrix has full rank rank(X) = k + 1 = p.

The classical *normal* linear regression model is obtained if additionally

1. 
$$\varepsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$$

holds. For random covariates these assumptions are to be understood conditionally on  $\boldsymbol{X}$ .

# Properties for the normal linear model

Least squares and maximum likelihood estimator for β:

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

with 
$$\hat{\boldsymbol{\beta}} \sim N_p(\boldsymbol{\beta}, \sigma^2(\boldsymbol{X}^T\boldsymbol{X})^{-1})$$
.

▶ Restricted maximum likelihood estimator for  $\sigma^2$ :

$$\hat{\sigma^2} = \frac{1}{n-p} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \frac{\mathsf{SSE}}{n-p}$$

with 
$$\frac{(n-p)\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-p}$$
.

Statistic for inference about  $\beta_j$ ,  $c_{jj}$  is diagonal element j of  $(\boldsymbol{X}^T\boldsymbol{X})^{-1}$ .

$$T_j = rac{\hat{eta}_j - eta_j}{\sqrt{c_{jj}}\hat{\sigma}} \sim t_{n-p}$$

3/30

5 / 30

# Output from fitting the full model in R

```
> fit=lm(y~.,data=ds)
> summary(fit)
Coefficients:
```

Estimate Std. Error t value Pr(>|t|) (Intercept) 5.6764334 0.1389162 40.862 < 2e-16 \*\*\* -0.3150444 0.0587512 -5.362 4.27e-05 \*\*\* x2 -0.0018533 0.0012587 -1.472 0.158 x3 0.9751745 0.1449075 6.730 2.62e-06 \*\*\* x4 -0.0002268 0.0010038 -0.226 x5 -0.0334242 0.0225009 -1.485 0.155 x6 -0.0039399 0.0724339 -0.054 0.957 0.398 x7 0.0888722 0.1025724 0.866

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1165 on 18 degrees of freedom Multiple R-squared: 0.93,Adjusted R-squared: 0.9027 F-statistic: 34.15 on 7 and 18 DF, p-value: 3.904e-09

Acid rain in Norwegian lakes

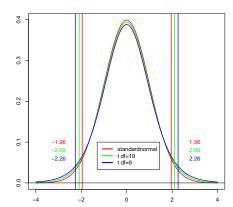
Measured pH in Norwegian lakes explained by content of

- ▶ x1: SO<sub>4</sub>: sulfate (the salt of sulfuric acid),
- ▶ x2: N0<sub>3</sub>: nitrate (the conjugate base of nitric acid),
- ▶ x3: Ca: calsium,
- x4: latent AI: aluminium,
- x5: organic substance,
- ▶ x6: area of lake,
- x7: position of lake (Telemark or Trøndelag),

Random sample of n = 26 lakes.

4 / 30

# Quantiles and critical values: N og t: $\alpha/2 = 0.025$



In R: specify area to the left, but our notation gives area to the right. Fahrmeir et al: notation with area to the left.

# Properties of the residuals

- ▶ Residuals (raw):  $\hat{\boldsymbol{\varepsilon}} = \boldsymbol{Y} \hat{\boldsymbol{Y}}$ .
- with mean  $E(\hat{\varepsilon}) = \mathbf{0}$  and covariance matrix  $Cov(\hat{\varepsilon}) = \sigma^2 (\mathbf{I} \mathbf{H})$  where  $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ .
- ▶ In the normal model  $\varepsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$  and then also the vector of residuals are normal, but with heteroscedastic variances and non-zero covariances.
- Standardized residuals: divide (raw) residuals by estimated standard deviation.
- Studentized residuals: leave-one-out version.
- ► Studentized residuals are compared with the normal distribution to assess normality of the error term.

7/30

# Simulating data and checking residuals

```
n=1000
beta=matrix(c(0,1,1/2,1/3),ncol=1)
set.seed(123)
x1=rnorm(n,0,1); x2=rnorm(n,0,2); x3=rnorm(n,0,3)
X=cbind(rep(1,n),x1,x2,x3)

y=X%*%beta+rnorm(n,0,2)
fit=lm(y~x1+x2+x3)
yhat=predict(fit)
summary(fit)
ehat=residuals(fit); estand=rstandard(fit); estud=rstudent(fit)
plot(yhat,ehat,pch=20)
points(yhat,estand,pch=20,col=2)
#points(yhat,estud,pch=20,col=5)
```

9 / 30

#### 3.12 Overview of Residuals

#### **Ordinary Residuals**

The residuals are given by

$$\hat{\varepsilon}_i = y_i - \hat{y}_i = y_i - x_i' \hat{\beta}$$
  $i = 1, ..., n$ 

#### Standardized Residuals

The standardized residuals are defined by

$$r_i = \frac{\hat{\varepsilon}_i}{\hat{\sigma}\sqrt{1 - h_{ii}}}$$

where  $h_{ii}$  is the ith diagonal element of the hat matrix.

#### Studentized Residuals

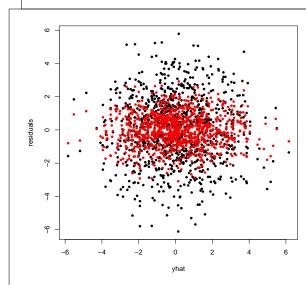
The studentized residuals are defined by

$$r_i^* = \frac{\hat{\varepsilon}_{(i)}}{\hat{\sigma}_{(i)}(1 + \mathbf{x}_i'(\mathbf{X}_{(i)}'\mathbf{X}_{(i)})^{-1}\mathbf{x}_i)^{1/2}} = \frac{\hat{\varepsilon}_i}{\hat{\sigma}_{(i)}\sqrt{1 - h_{ii}}} = r_i \left(\frac{n - p - 1}{n - p - r_i^2}\right)^{1/2}$$

The studentized residuals are used to verify model assumptions and to discover outliers (see Sect. 3.4.4).

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.126)

8/30



Black: raw residuals, red: standardized residuals (identical to studentized here).

# Examination of model assumptions

1. Linearity of covariates:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 

2. Homoscedastic error variance:  $Cov(\varepsilon) = \sigma^2 I$ .

3. Uncorrelated errors:  $Cov(\varepsilon_i, \varepsilon_i) = 0$ .

4. Additivity of errors:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 

5. Assumption of normality:  $\varepsilon \sim N_n(0, \sigma^2 I)$ 

11/30

# Volume of a tree

Data for 31 trees of a certain kind in a national park in the US are given below. Three variables are measured for each tree. These are:

- ➤ D: The diameter of the tree measured in inches 1.5 m above ground level
- ▶ H: The height of the tree measured in feet.
- V: The volume of the tree measured in cubic feet.

Obs.	D	Н	V	Obs.	D	Н	V
1	8.3	70	10.3	17	12.9	85	33.8
2	8.6	65	10.3	18	13.3	86	27.4
3	8.8	63	10.2	19	13.7	71	25.7
4	10.5	72	16.4	20	13.8	64	24.9
5	10.7	81	18.8	21	14.0	78	34.5
6	10.8	83	19.7	22	14.2	80	31.7
7	11.0	66	15.6	23	14.5	74	36.3
8	11.0	75	18.2	24	16.0	72	38.3
9	11.1	80	22.6	25	16.3	77	42.6
10	11.2	75	19.9	26	17.3	81	55.4
11	11.3	79	24.2	27	17.5	82	55.7
12	11.4	76	21.0	28	17.9	80	58.3
13	11.4	76	21.4	29	18.0	80	51.5
14	11.7	69	21.3	30	18.0	80	51.0
15	12.0	75	19.1	31	20.6	87	77.0
16	12.9	74	22.2				

13 / 30

# Plotting residuals

- 1. Plot the residuals,  $r_i^*$  against the predicted values,  $\hat{y}_i$ .
  - Dependence of the residuals on the predicted value: wrong regression model?
  - Nonconstant variance: transformation or weighted least squares is needed?
- 2. Plot the residuals,  $r_i^*$ , against predictor variable or functions of predictor variables. Trend suggest that transformation of the predictors or more terms are needed in the regression.
- 3. Assessing normality of errors: QQ-plots and histograms of residuals. As an additional aid a test for normality can be used, but must be interpreted with caution since for small sample sizes the test is not very powerful and for large sample sizes even very small deviances from normality will be labelled as significant.
- 4. Plot the residuals,  $r_i^*$ , versus time or collection order (if possible). Look for dependence or autocorrelation.

12 / 30

# Volume of a tree

- ▶ If one wants to measure the volume of a tree the tree has to be cut down.
- But, height and diameter can be measured without cutting down the tree.
- ► Of interest: develop a model that can be used to estimate the tree volume from the height and diameter.

As an illustration assume we want to fit a linear model with V as response and D and H as covariates. What is the  $R^2$  of this model?

Comment: if we start with the volume of a cylinder (area of circle times height) we may suggest a different regression model (on the log scale). Which model?

# Volume: height and diameter

```
fit <- lm(Volume~.,data=ds)</pre>
summary(fit)
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -57.9877
                        8.6382 -6.713 2.75e-07 ***
                        0.2643 17.816 < 2e-16 ***
Diameter
             4.7082
Height
             0.3393
                        0.1302 2.607 0.0145 *
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 '
Residual standard error: 3.882 on 28 degrees of freedom
Multiple R-squared: 0.948, Adjusted R-squared: 0.9442
F-statistic: 255 on 2 and 28 DF. p-value: < 2.2e-16
```

15 / 30

# Volume: height and diameter – and IQ of lumberjack

```
set.seed(123) # reproducible results
iq < rnorm(31,100,16)
fit2 <- lm(Volume~Height+Diameter+iq,data=ds)</pre>
summary(fit2)
```

```
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -61.03399 10.20868 -5.979 2.24e-06 ***
Height
             0.34099
                        0.13176 2.588 0.0154 *
Diameter
             4.72507
                        0.26906 17.561 2.68e-16 ***
             0.02704
                        0.04678 0.578 0.5681
iq
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 '
Residual standard error: 3.929 on 27 degrees of freedom
Multiple R-squared: 0.9486, Adjusted R-squared: 0.9429
```

F-statistic: 166.1 on 3 and 27 DF, p-value: < 2.2e-16

Volume of a tree: IQ of lumberjack added

- ▶ We want to add the IQ of the lumberjack that cut down the tree as a covariate in the model
- ► This should for obvious reasons not be a good predictor for the volume of the tree.
- ▶ To mimic this situation we simulate new data to resemble the IQ of different lumberjacks by drawing data from the normal distribution with mean 100 and standard deviation 16. and since we have 31 trees we simulate 31 observations.
- $\triangleright$  Q: will the  $R^2$  of this new model be higher than the  $R^2$  of the previous model?

16 / 30

# Acid rain in Norwegian lakes

Data on n = 26 lakes, with

- y: measured pH in lake,
- ► x1: SO<sub>4</sub>: sulfate (the salt of sulfuric acid).
- ► x2: N0<sub>3</sub>: nitrate (the conjugate base of nitric acid),
- x3: Ca: calsium.
- ▶ x4: latent *AI*: aluminium.
- ► x5: organic substance,
- ▶ x6: area of lake.
- x7: position of lake (Telemark or Trøndelag).

We would like to use a regression model with pH of the lake as the response. Should we fit a model will all 7 covariates, or choose a subset?

### Simulated data (Fahrmeir et al: Fig 3.17)

True model:

$$Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \varepsilon_i$$

Known that the model is polynomial in nature, but not up to which degree.

Try to fit polynomial also with higher order terms.

New: in addition to the data set to be used to fit the regression (called *training set*) also a data set to assess the model fit is present (called a *validation* set).

Mean Squared Error (MSE) is a scaled version of the SSE, that is  $\frac{1}{n}\sum_{i=1}^{n}(Y_i-\hat{Y}_i)^2$ .

19 / 30

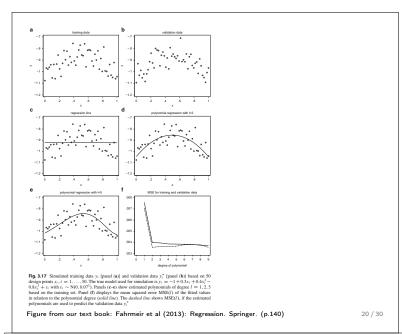
### Simulated data (Fahrmeir et al: Fig 3.18, Tab3.3, Tab3.4)

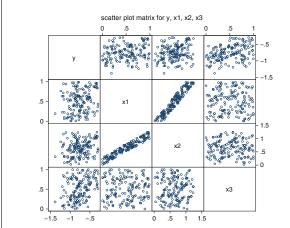
True model:

$$Y \sim N(-1 + 0.3x_1 + 0.2x_3, 0.2^2)$$

where also  $x_2 = x_1 + u$  is observed ( $u \sim$  uniform in 0,1). The variables  $x_1$  and  $x_3$  are uncorrelated.

21 / 30





ig. 3.18 Scatter plot matrix for the variables y,  $x_1$ ,  $x_2$ , and  $x_3$ 

Figure from our text book: Fahrmeir et al (2013): Regression. Springer. (p.141)

**Table 3.3** Results for the model based on covariates  $x_1$ ,  $x_2$ , and  $x_3$ 

Variable	Coefficient	Standard error	t-value	p-value	95 % Con	fidence interval
intercept	-0.970	0.047	-20.46	< 0.001	-1.064	-0.877
$\overline{x_1}$	0.146	0.187	0.78	0.436	-0.224	0.516
<i>x</i> <sub>2</sub>	0.027	0.177	0.15	0.880	-0.323	0.377
<i>x</i> <sub>3</sub>	0.227	0.052	4.32	< 0.001	0.123	0.331

**Table 3.4** Results for the correctly specified model based on covariates  $x_1$  and  $x_3$ 

Variable	Coefficient	Standard error	t-value	p-value	95 % Con	fidence interval
intercept	-0.967	0.039	-24.91	< 0.001	-1.042	-0.889
$\overline{x_1}$	0.173	0.055	3.17	0.002	0.065	0.281
<i>x</i> <sub>3</sub>	0.226	0.052	4.33	< 0.001	0.123	0.330

Table from our text book: Fahrmeir et al (2013): Regression. Springer. (p.142)

23 / 30

### Two subsets of covariates (Exam V2014 Problem 4b)

Classical linear model with identically normally distributed random errors,  $Cov(\varepsilon) = \sigma^2 I$ , but now look at misspecification of E(Y). Suppose that the true model is

$$Y = X_1 \beta_1 + X_2 \beta_2 + \varepsilon,$$

$$\varepsilon \sim N_0(0, \sigma^2 I),$$
(1)

where we have partitioned the design matrix into two parts  $X_1$   $(n \times p_1)$  and  $X_2$   $(n \times p_2)$  and  $\beta_1$  and  $\beta_2$  are unknown  $p_1$ - and  $p_2$ -dimensional vectors of regression coefficients  $(p = p_1 + p_2)$ .

25 / 30

### Irrelevant and/or missing covariates in the regression

Irrelevant: variables that are included in the regression but should not have been.

missing : variables that are not included, but should have been.

24 / 30

### Two subsets of covariates (cont.)

Assume that we ignore the covariates in  $\boldsymbol{X}_2$  and fit the model

$$\mathbf{Y} = \mathbf{X}_1 \alpha_1 + \delta, \delta \sim N_n(\mathbf{0}, \tau^2 \mathbf{I}).$$
 (2)

Here  $\alpha_1$  is used in place of  $\beta_1$  to emphasize that  $\alpha_1$  (and estimates thereof) will in general be different from  $\beta_1$  in the true model. The least squares estimator for model (2) is  $\hat{\alpha_1} = (\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1} \boldsymbol{X}_1^T \boldsymbol{Y}$ .

### Two subsets of covariates (cont.)

Find the expected value and covariance matrix of  $\hat{\alpha_1}$  under the true model.

$$E(\hat{\alpha_1}) = \beta_1 + (\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1} \boldsymbol{X}_1^T \boldsymbol{X}_2 \beta_2$$

We see that the bias term for  $\hat{\alpha}_1$  is  $(\boldsymbol{X}_1^T\boldsymbol{X}_1)^{-1}\boldsymbol{X}_1^T\boldsymbol{X}_2\boldsymbol{\beta}_2$ . When is the bias term equal to zero?

$$\operatorname{Cov}(\hat{\boldsymbol{\alpha}_1}) = \sigma^2(\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1}$$

Observe,  $Cov(\hat{\alpha_1})$  is not dependent on  $\beta_2$ .

27 / 30

### Irrelevant covariates included: findings

Bias: The estimator for the true covariates are unbiased, also if irrelevant covariates are included.

Variance: The model with the irrelevant covariants have larger variance for the true covariates, compared with the model without the irrelevant covariates. So, again sparse model is the best.

Missing covariates: findings

Bias: The estimator for the (true) covariates (in the model) is only unbiased if the true and missing covariates are uncorrelated (orthogonal design) in the data

Variance: The variance of the estimator for the true covariates may be smaller based on the model with the missing covariates (than for the correctly specified model), and even the sum of the bias<sup>2</sup> and the variance may better for the model with the missing variables. So the sparse model may be better on overall (even though it is biased).

28 / 30

### Irrelevant and/or missing covariates in the regression

Irrelevant: variables that are included in the regression but should not have been.

missing : variables that are not included, but should have been.

Conclusion in book: the model should not contain irrelevant covariates, and we should aim for a sparse model.

### Law of parsimony

If two models are not very different - then always choose the simplest one

31/30

### Hypothesis toling about \$1

21.02.2017 11

By is the 1th cleaner of B= (XTX) - XTY \$~ Np(B, (XIX)-101)

Text for association (linear) between surposes Y





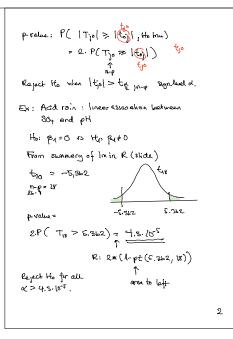
### Today

- ▶ T-test for significance of one regression coefficient.
- ▶ Residuals: standardized (or studentized) preferred.
- ▶ Significance of regression based on F-test with SSR/(p-1) divided by SST/(n-1).
- $ightharpoonup R^2$  gives the proportion of variability explained by the regression model.

$$R^2 = \frac{\text{SSR}}{\text{SST}} = 1 - \frac{\text{SSE}}{\text{SST}}$$

and will never decrease if new covariates are added to the model.

▶ Model selection: want to choose the model that minimize the expected squared prediction error.



### Residuals (again)

The residuan have heteroscedestic resience  $Var(\hat{\mathcal{C}}) = 0^{\alpha} (1-hij)$  and  $Gv(\hat{\mathcal{C}}_i,\hat{\mathcal{C}}_j) = 0^{\alpha} (0-hij)$  can in general lead to plant in most common and the common and the

### Standardrzed residuals:

### R: rstanderd (tit)

Stadentized residuals: Fitting the woodet to all obs. except i to make  $r_i^*$ . — see slide

hee studentized! R: (Student (pt)

See exemple on slide for (i s &:

3

Ex: Acid rain, full model (all coverses available) R2= 0.73, 93%

Is the regression significent?

the at least one py +0 j=1,..., k

Text statistics:  $f = \frac{3SR/k}{8SE/(n-p)} \sim F_k$ , where

prove this in Port 3 in a general setting

### Ex: Acid rain:

2

### Analysis of varience decomposition and R2

$$\frac{\sum_{i=1}^{n} (y_i - \overline{y})^2}{\text{Supp of squared}} = \frac{\sum_{i=1}^{n} (y_i - \overline{y}_i^2 + \overline{y}_i^2 - \overline{y}_i^2)^2}{\sum_{i=1}^{n} (y_i - \overline{y}_i^2)^2} + \frac{\sum_{i=1}^{n} (y_i^2 - \overline{y}_i^2)^2}{\sum_{i=1}^{n} (y_i^2 - \overline{y}_i^2)^2} + \frac{\sum_{i=1}^{n} (y_i^2 - \overline{y}_i^$$

With vector and matrices. Atc

This is used to define:

coefficient of determination

relative propertien of total recreatility explaned by the regression

4

Ex: Volume of tree and the lumber jack

Big model: Yi=po+paxi+ paxi+paxi+ e.

Small model: Yi=po+paxi+ paxi+ fixi+ fix

R'ony ≥ R'onel ← since fis is found

to minimize SSE-ÉtÉ, thus maximize

Re= 1- 85€ T28

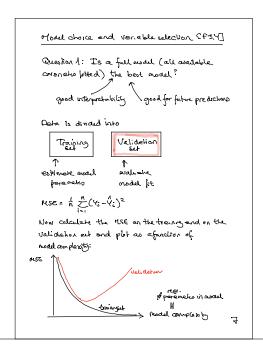
 $R^2$ 819 =  $R^2$ 811011 if  $\frac{1}{8}$ 3 = 0 if  $\frac{1}{8}$ 5 + 0 the SSE<sub>819</sub> < SSE<sub>81111</sub> and  $R^2$ 819 >  $R^2$ 8111.

Re will always increase (or stay unchanged) when a new correlate is added to the model.

Next: next on choosing a good model, and then  $R^{E}ady = 1 - \frac{33E/(n-p)}{85T/(n-1)} \stackrel{\text{penoticus}}{\text{onens}}$ 

is one criterion to one instead of 2° for model selection.

6



# TMA4267 Linear Statistical Models V2017 (L12)

Part 2: Linear regression:
Model selection [F:3.4]
Transformation and Taylor expansion
Quiz

### Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 24, 2017

Answer 1: No, this may lead to overhing = fitting the trend + the norms:

=> so, what can we do instead?

8

### What is the "best" model?

Acid rain in Norwegian lakes, data on n=26 lakes, with

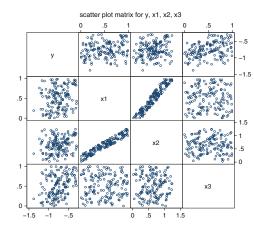
- ▶ y: measured pH in lake,
- ► x1: SO<sub>4</sub>: sulfate (the salt of sulfuric acid),
- ▶ x2: N0<sub>3</sub>: nitrate (the conjugate base of nitric acid),
- x3: Ca: calsium,
- ▶ x4: latent AI: aluminium,
- x5: organic substance,
- ► x6: area of lake,
- ▶ x7: position of lake (Telemark or Trøndelag),

### Topic: choosing the "best" linear regression model!

- ► First, debunk popular strategies (based on simulations studies were we knew the "true" model):
  - Popular 1: fit all available covariates.
     Problem: overfitting (=fitting trends and noise).
  - Popular 2: fit all available covariates, then remove the insignificant ones (=those  $\beta_j$  where  $H_0: \beta_j = 0$  is not rejected).

2 / 47

4 / 47



ig. 3.18 Scatter plot matrix for the variables y,  $x_1$ ,  $x_2$ , and  $x_3$ 

Figure from our text book: Fahrmeir et al (2013): Regression. Springer. (p.141)

Simulated data (Fahrmeir et al: Fig 3.18, Tab3.3, Tab3.4)

True model:

$$Y \sim N(-1 + 0.3x_1 + 0.2x_3, 0.2^2)$$

where also  $x_2 = x_1 + u$  is observed ( $u \sim$  uniform in 0,1). The variables  $x_1$  and  $x_3$  are uncorrelated.

3 / 47

**Table 3.3** Results for the model based on covariates  $x_1$ ,  $x_2$ , and  $x_3$ 

Variable	Coefficient	Standard error	t-value	p-value	95 % Con	fidence interval
intercept	-0.970	0.047	-20.46	< 0.001	-1.064	-0.877
$\overline{x_1}$	0.146	0.187	0.78	0.436	-0.224	0.516
$x_2$	0.027	0.177	0.15	0.880	-0.323	0.377
$\overline{x_3}$	0.227	0.052	4.32	< 0.001	0.123	0.331

**Table 3.4** Results for the correctly specified model based on covariates  $x_1$  and  $x_3$ 

Variable	Coefficient	Standard error	t-value	p-value	95 % Confidence interva	
intercept	-0.967	0.039	-24.91	< 0.001	-1.042	-0.889
$x_1$	0.173	0.055	3.17	0.002	0.065	0.281
<i>x</i> <sub>3</sub>	0.226	0.052	4.33	< 0.001	0.123	0.330

Table from our text book: Fahrmeir et al (2013): Regression. Springer. (p.142)

### Topic: choosing the "best" linear regression model!

- ► First, debunk popular strategies (based on simulations studies were we knew the "true" model):
  - Popular 1: fit all available covariates.
     Problem: overfitting (=fitting trends and noise).
  - Popular 2: fit all available covariates, then remove the insignificant ones (=those  $\beta_j$  where  $H_0: \beta_j = 0$  is rejected). Problem: may also remove important covariates that are correlated with unimportant ones but insignificant because being masked by the unimportant ones.
- Study of irrelevant and missing covariates:

Irrelevant: variables that are included in the regression but should not have been (IQ of lumberjack)

missing: variables that are not included, but should have been (omitting height in the tree volum example)

Conclusion in book: the model should not contain irrelevant covariates, and we should aim for a sparse model.

Take home message is the "Law of parsimony": *If two models* are not very different – then always choose the simplest one.

6 / 47

### Expected squared prediction error (SPSE)

Possible criterion we want to minimize: SPSE. Definition (j, M, ... given in classnotes)

$$\mathsf{SPSE} = \sum_{j=1}^J \mathrm{E}((Y_j - \hat{Y}_{jM})^2)$$

can be written as:

$$\mathsf{SPSE} = \sum_{j=1}^J \mathrm{E}((Y_j - \hat{Y}_{jM})^2) = n\sigma^2 + |M|\sigma^2 + \sum_{j=1}^J (\mu_{jM} - \mu_j)^2$$

Problem: Not useful on practise since  $\mu_j$  and  $\sigma^2$  are unknown. Plan: Find a way to estimate SPSE and then choose the model M with the minimum SPSE!

8 / 47

### All models are wrong?

A model is a simplification or approximation of reality and hence will not reflect all of reality.

George Box noted that "all models are wrong, but some are useful". While a model can never be "truth" a model might be ranked from very useful, to useful, to somewhat useful to, finally, essentially useless.

Burnham, K. P.; Anderson, D. R. (2002), Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach.

7 / 47

### How to estimate SPSE?

$$\mathsf{SPSE} = \sum_{i=1}^J \mathrm{E}((Y_j - \hat{Y}_{jM})^2)$$

Assume we have fitted a model M with |M| regression parameters.

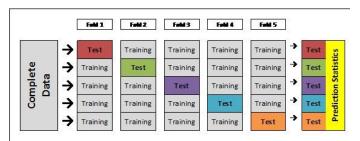
1. Use new (independent) data – if available (seldom the case):

$$\widehat{SPSE} = \sum_{j=1}^{J} (Y_j - \hat{Y}_{jM})^2$$

2. Cross-validation: mimic new data by dividing data into k folds (popular is k=n and k=10). In a for-loop let  $j=1,\ldots,k$ , and use all folds except fold j to estimate regression parameter, and use the jth fold to calculated the  $\widehat{SPSE}$ . Sum across folds.

Choose the model M that minimizes the  $\widehat{SPSE}$ .

### Cross-validation (5-fold)



Will be taught in TMA4300 Computational statistics and will be a backbone in TMA4268 Statistical Learning.

http://blog-test.goldenhelix.com/wp-content/uploads/2015/04/B-fig-1.jpg

10 / 47

# For models with the same model complexity – easy solution: SSF

Estimators for SPSE to be used on the same data as to be used for estimating the model parameters have the same form; a first term based on SSE (or  $R^2$ ) for model M, and a second term penalizing the model complexity.

If we consider two models with the same model complexity then SSE can be used to choose between these models.

12 / 47

### How to estimate SPSE?

$$SPSE = \sum_{j=1}^{J} E((Y_j - \hat{Y}_{jM})^2)$$

Assume we have fitted a model M with |M| regression parameters.

3. Use existing data (only): It can be shown that  $E(\widehat{SPSE}) = SPSE - 2 \mid M \mid \sigma^2$  when used on the same data that was used to make the prediction, so a better estimate for existing data is

$$\widehat{SPSE} = \sum_{i=1}^{n} (Y_i - \hat{Y}_{iM})^2 + 2|M|\hat{\sigma}^2 = SSE + 2|M|\hat{\sigma}^2$$

where  $\hat{\sigma}^2$  is the same for all models M, and is often estimated using the most complex model under study.

4. Other criteria: all have the same form; a first term based on SSE (or  $R^2$ ) for model M, and a second term penalizing the model complexity.

Choose the model M that minimizes the  $\widehat{SPSE}$ .

11 / 47

### Acid rain (1). Best subset

For 1,...,7 covariates: fit all possible models, and report the model with the smallest SSE (given below) for each value for the model complexity. Explain what you see! How many models have been searched for each model complexity?

Names: x1:  $SO_4$ , x2:  $NO_3$ , x3: Ca, x4: latent AI, x5: organic substance, x6: area of lake, x7: position of lake (Telemark or Trøndelag).

### Popular model choice criteria

 $R^2$  adjusted (corrected) Mallows'  $C_p$ Akaike Information Criterion (AIC) Bayesian Information Criterion (BIC)

NB: there is no overall best choice for criterion - all of these are used.

14 / 47

### Happiness (n = 39)

Are love and work the important factors determining happiness?

- y, happiness. 10-point scale, with 1 representing a suicidal state, 5 representing a feeling of «just muddling along», and 10 representing a euphoric state.
- ▶ x<sub>1</sub>, money. Annual family income in thousands of dollars.
- x<sub>2</sub>, sex. Sex was measured as the values 0 or 1, with 1 indicating a satisfactory level of sexual activity.
- x<sub>3</sub>, love. 3-point scale, with 1 representing loneliness and isolation, 2 representing a set of secure relationships, and 3 representing a deep feeling of belonging and caring in the context of some family or community.
- x<sub>4</sub>, work. 5-point scale, with 1 indicating that an individual is seeking other employment, 3 indicating the job is OK, and 5 indicating that the job is enjoyable.

Data taken from library faraway, data set happy.

16 / 47

### $R^2$ adjusted (corrected)

 $\hat{Y}_i$  is from fitting the regression model M. Remember, for a regression model (with intercept) we have the SST=SSR+SSE.

$$SST = \sum_{i=1}^{n} (Y_i - \bar{Y})^2$$

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

$$R^2 = 1 - \frac{SSE}{SST}$$

$$R_{adj}^2 = 1 - \frac{\frac{SSE}{n-p}}{\frac{SST}{n-p}} = 1 - \frac{n-1}{n-p} (1 - R^2)$$

Choose the model with the largest  $R_{\text{adj}}^2$ 

"All" statistical software outputs this automatically! However, Fahrmeir et al (2013) believes that the penalty n-p is too small.

15 / 47

### Нарру

	money	sex	love	work	Ν	р	$R^2$	$R_{\rm adj}^2$
1	0.014				1	0.000747	7.3	4.8
2		-0.130			1	1	0.1	-2.6
3			2.270		1	8.35e-24	61.5	60.5
4				0.990	1	1.36e-13	29.1	27.2
5	0.016	-0.508			2	0.0504	8.8	3.8
6	0.009		2.206		2	8.77e-19	64.5	62.5
7	0.012			0.961	2	3.68e-10	34.6	31.0
8		-0.277	2.279		2	5.55e-18	62.0	59.9
9		0.610		1.079	2	3.48e-09	31.2	27.4
10			1.959	0.511	2	5.75e-20	68.1	66.3
11	0.011	-0.536	2.209		3	9.49e-16	66.2	63.3
12	0.011	0.305		1.009	3	1.84e-07	35.1	29.5
13	0.009		1.902	0.504	3	2.63e-17	70.9	68.4
14		0.108	1.944	0.530	3	2.22e-16	68.1	65.4
15	0.010	-0.149	1.919	0.476	4	9.89e-15	71.0	67.6

Intercept included, N=p-1, p-value for significance of regression.  $R^2=1-\frac{SSE}{SST}$ ,  $R^2_{adj}=1-\frac{SSE}{\frac{SSE}{SST}}$ . Which model to prefer?

18 / 47

### AIC

Akaike information criterion – one of the most widely used. Designed for likelihood-based inference.

For a normal regression model:

$$AIC = n \ln(\hat{\sigma}^2) + 2(|M| + 1)$$

Choose the model with the minimum AIC.

Mallows'  $C_p$ 

 $\hat{Y}_i$  is from fitting regression model M. Mallows is the name of a person.

$$C_p = \frac{\sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2}{\hat{\sigma}^2} - n + 2|M|$$

Minimizing Cp gives the same optimal model as minimizing  $\widehat{SPSE}$ .

See Exam V2015 Problem 3 for an in depth explanation of the theory behind Mallow's Cp.

19 / 47

BIC

Bayesian information criterion.

For a normal regression model:

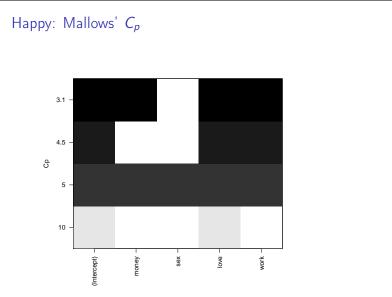
$$\mathsf{BIC} = n \ln(\hat{\sigma}^2) + \ln(n)(|M| + 1)$$

Choose the model with the minimum BIC.

AIC and BIC are motivated in very different ways, but the final result for the normal regression model is very similar.

BIC has a larger penalty than AIC (log(n)vs.2), and will often give a smaller model (=more parsimonious models) than AIC.

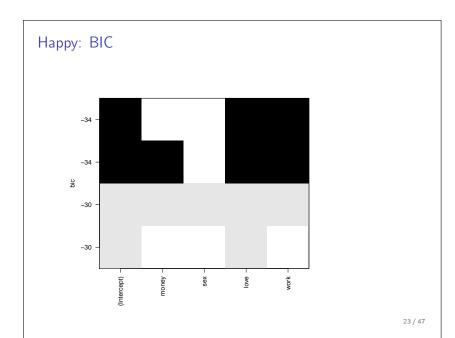
21 / 47

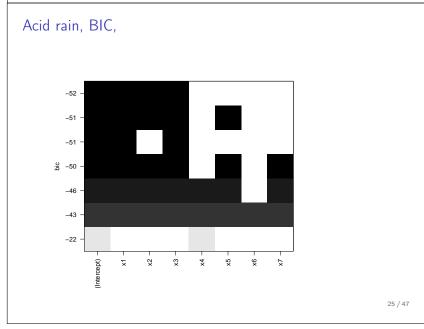


22 / 47

24 / 47

## Acid rain (2)





### Practical use of the model criteria

- ► All subset selection: use smart "leaps and bounds" algorithm, works fine for number of covariates in the order of 40.
- ► Forward selection: choose starting model (only intercept), then add one new variable at each step selected to make the best improvement in the model selection criteria. End when no improvement is made.
- Backward elimination: : choose starting model (full model), then remove one new variable at each step - selected to make the best improvement in the model selection criteria. End when no improvement is made.
- ▶ Stepwise selection: combine forward and backward.

26 / 47

### Acid rain (4): Forward

28 / 47

### Acid rain (3): stepAIC

```
> all=lm(happy~.,data=happy)
> stepAIC(all)
Start: AIC=9.08
happy ~ money + sex + love + work
       Df Sum of Sq RSS AIC
<none>
                   38.087 9.076
             3.782 41.869 10.768
- money 1
             6.386 44.473 13.122
- work 1
- love 1 47.272 85.359 38.549
Step: AIC=7.22
happy ~ money + love + work
       Df Sum of Sq RSS AIC 38.229 7.221
<none>
- money 1 3.723 41.952 8.846
             8.410 46.639 12.976
- work 1
- love 1 47.742 85.971 36.828
lm(formula = happy ~ money + love + work, data = happy)
(Intercept)
                             love
  -0.185936
             0.008959
                          1.901709
                                      0.503602
```

27 / 47

### Acid rain (5): Backward

### Model diagnosis

- Influential observations and outliers: impact of specific observations on model fit.
- ▶ Collinearity analysis: Highly correlated variables cause imprecise estimation of the regression parameters. (Why? Look at diagonal elements of  $Cov(\hat{\beta}) = \sigma^2(\boldsymbol{X}^T\boldsymbol{X})^{-1}$ , and look back to Problem 2 in the start of this lecture.)
- ► Examination of model assumptions: residual plots!

30 / 47

### Transformations

- ► Multiplicative or additive model?
- ▶ Box–Cox transform with profile likelihood.
- ► Stabilizing the variance.

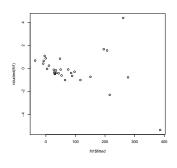
### Influential observations— and outliers

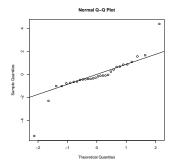
- ► Observations that significantly affect inferences drawn from the data are said to be influential.
- ▶ The leverage,  $h_{ii}$ , associated with the *i*th datapoint measures "how far the *i*th observation is from the other n-1 observations".
- Methods for assessing influential observations may be be based on change in β estimate when observations are deleted.
- ► Always investigate possible causes of an influential observation (if possible).
- ► Cook's distance can be used to identify influential observations.
- ▶ Robust methods (median, quantile regression) can be useful.

Want to understand more? Read for yourself in Fahrmeir et al (2013): p 160-166.

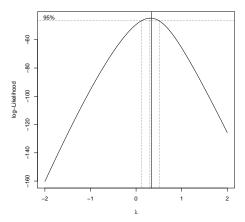
31 / 47

### Galapagos islands, Model A, Exam V2014 Problem 2





### Box-Cox plot



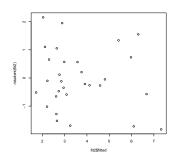
Box–Cox transformation plot based on Model A for the Galapagos data set, RecEx4. Line at x=1/3.

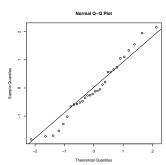
34 / 47

# Approximation of E and Var for nonlinear functions

- ▶ Have RV X, with mean  $E(X) = \mu$  and some variance Var(X).
- ▶ Want to look at a nonlinear function of X, called g(X).
- ▶ Aim: find an approximation to E(g(X)) and Var(g(X)).
- ▶ And, the same for two RVs  $X_1$  and  $X_2$  with  $g(X_1, X_2)$ .

### Galapagos islands, Model B, Exam V2014 Problem 2





35 / 47

### Example In of BMI

Looking at residual plots from a regression model the conclusion was to analyse data of BMI on the natural logarithmic scale. After a regression model was fitted the predicted value for the ln(BMI) for a specific combination of the covariates was found to be 3.2151 with an estimated standard deviation of 0.1656. Use approximate methods to arrive at an estimate of the predicted value and estimated standard deviation on the original scale,  $kg/m^2$ , and not on the logarithmic scale.

36 / 47

### E(g(X)) and Var(g(X))

- Let g(X) be a general function. When is E(g(X)) = g(E(X))?
  - ▶ When g(X) is a linear function of X.
- ▶ What can we do if this is not the case?
  - ▶ We can calculate  $E(g(X)) = \int_{-\infty}^{\infty} g(x)f(x)dx$  when X is continuous, or a version thereof in the discrete case,
  - or if g is monotone we can use the transformations formula to find the distribution of Y = g(X) and then calculate E(Y) and Var(Y), if possible.
- ▶ What if we only know  $E(X) = \mu$  and  $Var(X) = \sigma^2$  and not f(x)?
  - ▶ Use a Taylor series approximation of g(X) around  $g(\mu)$ . g need to be differentiable.

38 / 47

# Treatment of tennis elbow (exam TMA4255 V2012, 3b)

The term *tennis elbow* is used to describe a state of inflammation in the elbow, causing pain. This injury is common in people who play racquet sports, however, any activity that involves repetitive twisting of the wrist (like using a screwdriver) can lead to this condition. The condition may also be due to constant computer keyboard and mouse use.

In a randomized clinical study the aim was to compare three different methods for treatment of tennis elbow.

- ► A: physiotherapy intervention,
- ▶ B: corticosteroid injections and
- ► C: wait-and-see (the patients in the wait-and-see group did not get any treatment but was told to use the elbow as little as possible).

Univariate function

First order Taylor approximation of g(X) around  $\mu$ .

$$g(X) \approx g(\mu) + g'(\mu)(X - \mu)$$

This leads to the following approximations:

$$\mathrm{E}(g(X)) \approx g(\mu)$$
  
 $\mathrm{Var}(g(X)) \approx [g'(\mu)]^2 \mathrm{Var}(X)$ 

39 / 47

### Treatment of tennis elbow (cont.)

We will look at the short-term effect of treatment by studying measurements at 6 weeks. All patients participating in the study only had one affected arm.

We will look at the outcome measure called *pain-free grip force*. This was measured by a digital grip dynamometer and normalized to the grip force of the unaffected arm. A pain-free grip force of 100 would mean that the affected and the unaffected arm performed equally good.

Summary statistics for each of the treatment groups.

Treatment	Sample size	Average	Standard deviation
A (physiotherapy)	63	70.2	25.4
B (injection)	65	83.6	22.9
C (wait-and-see)	60	51.8	23.0
Total	188	69.0	

42 / 47

### Bivariate function: first order Taylor

 $X_1$  is a RV with  $\mu = E(X_2)$  and  $X_2$  is a RV with  $\mu_2 = E(X_2)$ . Let g be a bivariate function of  $X_1$  and  $X_2$ , and define

$$g_1'(\mu_1, \mu_2) = \frac{\partial g(x_1, x_2)}{\partial x_1} \mid_{x_1 = \mu_1, x_2 = \mu_2}$$
$$g_2'(\mu_1, \mu_2) = \frac{\partial g(x_1, x_2)}{\partial x_2} \mid_{x_1 = \mu_1, x_2 = \mu_2}$$

First order Taylor approximation:

$$g(X_1, X_2) \approx g(\mu_1, \mu_2) + g'_1(\mu_1, \mu_2)(X_1 - \mu_1) + g'_2(\mu_1, \mu_2)(X_2 - \mu_2)$$

44 / 47

### Example 2: Exam TMA4255 V2012 3d (fraction)

Let  $\mu_A$  be the expected pain-free grip force for a population where the physiotherapy intervention treatment is used to treat tennis elbow, and  $\mu_C$  be the expected pain-free grip force for a population where the wait-and-see treatment is used. Define the relative difference between these two expected values as

$$\gamma = \frac{\mu_{\mathsf{A}} - \mu_{\mathsf{C}}}{\mu_{\mathsf{C}}}.$$

This can be interpreted as the expected relative gain by using physiotherapy instead of wait-and-see. Based on two independent random samples of size  $n_A$  and  $n_C$  from the physiotherapy and wait-and-see treatment groups, respectively, suggest an estimator,  $\hat{\gamma}$ , for  $\gamma$ .

Use approximate methods to find the expected value and variance of this estimator, that is,  $E(\hat{\gamma})$  and  $Var(\hat{\gamma})$ .

43 / 47

### Bivariate function: first order Taylor

$$\mathrm{E}(g(X_1, X_2)) \approx g(\mu_1, \mu_2)$$
 $\mathrm{Var}(g(X_1, X_2)) \approx [g_1'(\mu_1, \mu_2)]^2 \mathrm{Var}(X_1) + [g_2'(\mu_1, \mu_2)]^2 \mathrm{Var}(X_2) + 2 \cdot g_1'(\mu_1, \mu_2) \cdot g_2'(\mu_1, \mu_2) \mathrm{Cov}(X_1, X_2)$ 

### Multivariate version

### From Tabeller og formler i statistikk.

### Rekkeutvikling

En første ordens Taylorutvikling av funksjonen  $g(X_1,\dots,X_n)$  omkring  $g(\mu_1,\dots,\mu_n)$ , der  $\mathrm{E}(X_i)=\mu_i,\ i=1,\dots,n$ , gir approksimasjonene

$$E[g(X_1, ..., X_n)] \approx g(\mu_1, ..., \mu_n),$$

$$\operatorname{Var}[g(X_1,\dots,X_n)] \approx \sum_{i=1}^n \left(\frac{\partial g(\mu_1,\dots,\mu_n)}{\partial \mu_i}\right)^2 \operatorname{Var}(X_i) + 2 \sum_{i>j} \frac{\partial g}{\partial \mu_i} \frac{\partial g}{\partial \mu_j} \operatorname{Cov}(X_i,X_j).$$

46 / 47

Model selection [F.3.4]

L12. 24,02.2017

SPSE = Expected squared prediction error

 $E(Y) = \mu$  is the truth, but we modely as  $\mu = X_{FM}$  and we assume

nodel to X11... ) ha based on a subset of all the available Corsisha

TRAINING: i=1,-1,0 on observations, available Yi, Xi and fire is the advanction from the training into for our model M.

UALIDATION: j=1,...,T new observation, available as  $Y_j$  and  $X_jT$ .

$$\sum_{j=1}^{T} E\left( (Y_j - Y_j n_j)^2 \right) = SPSE$$

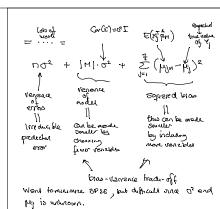
$$\sum_{j=1}^{n_{out}} \int_{Olso}^{n_{out}} SPSE + S$$

### Today

- Choosing between models of equal model complexity: choose the model with the minimum SSE.
- ► Choosing between models of *different model complexity*: Model selection based on penalized criteria (Mallows Cp,  $R_{\text{adj}}^2$ ,AIC and BIC). Try out on RecEx4 and Compulsory Exercise 2.
- ▶ BoxCox transformation: see RecEx4.
- Work for for yourself: Taylor solution to E and Var of nonlinear function, useful when you want to look at transformations of the data or functions of parameter estimates.

Summary of Part 2 in Kahoot!

47 / 47



PLAN: elimate SPSE and choose the of that minimizes this epipovale.

2

# Finding the best model all subsels onethod

1) Have k coveriets that night be used

How many possible module can I make (went to have increase)?

2-2-2-2 = 24 possible module

Fit all possible 2h models.

2) For M = 10,1,2,.., k3 choose the model with the smallest SSE.

Ex: Acid rain: k=7 > total 22= 128 possible reado

Complexity  [M = 1 2:	seeran d	but nodel	
M  = 1 2:	7	X4 (AL)	
MI = 253:	(2)=21	x, 123	
(n) =954:	(Z) = 33	X1, X21 X3	
orte 45.	(t) = 34	χ <sub>1</sub> , χ <sub>2</sub> , χ <sub>3</sub> , χ <sub>5</sub>	
H1= 86	(§)	X1, X2, X3, K5, X7	
M= 52	(3)	X1, X1, X3, X4, X5, X+	
n = 2 8	(2) <u>1</u>	X1, X2, X3, X4, X5, X6, X4	
	94		-

Ex: Happy: BIC best award = love + were

Homework : slide 24 Acidrain

Transfermation of response and predictes might improve the lit of the regionsion model.

The BoxCox trensform

$$g_{\lambda}(\alpha) = \begin{cases} y^{\lambda-1} & \lambda \neq 0 \\ \chi & \lambda \neq 0 \end{cases}$$
Close of function

For Yo XB+C, en N(0,0°F) the but value of A on board on maximizing the likelihood

SSEA (5 the SSE When Ja(4) is the response

R: boxcon (fit), see plat.

5

- 3) Now we need to choose between these kell mouth found in 2). Which criterian should I was?

  835 [17] op-kel
  - i)  $R^2$  and  $= 1 \frac{R \times E}{(n-y)}$

Ex: Happiness:

- ii) Hellows' Cp = SSE n+ 2|n|

  No SPSE c SSE + 2|n| · Fine. Type some coult.
- iii) AIC = n. ln(3°) + 2(1111+1)
- iv) BK = n. (n(ô2)+ (n(n) (1+1+1)
  BK gives own penally than Aic to large roads.

4

### TMA4267 Linear statistical models

Part 2: Linear regression

February 20, 2017

### Normal equations

$$\mathbf{Y} = \mathbf{X}\mathbf{\beta} + \mathbf{\epsilon}$$
 where  $E(\mathbf{\epsilon}) = \mathbf{0}$  and  $Cov(\mathbf{\epsilon}) = \sigma^2 \mathbf{I}$ 

Which of the following are the normal equations?

$$\mathbf{A} \quad \boldsymbol{X}\hat{\boldsymbol{\beta}} = \boldsymbol{H}\boldsymbol{Y}$$

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

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

$$\mathbf{D} \quad (\mathbf{X}^T \mathbf{X}) \mathbf{Y} = \mathbf{X}^T \hat{\boldsymbol{\beta}}$$

### Estimator for $\sigma^2$

$$m{Y} = m{X}m{\beta} + m{\epsilon}$$
 where  $E(m{\epsilon}) = m{0}$  and  $Cov(m{\epsilon}) = \sigma^2 m{I}$   
 $m{H} = m{X}(m{X}^Tm{X})^{-1}m{X}^T$ 

An unbiased estimator for  $\sigma^2$  is:

**A** SSE/
$$n$$
 **B**  $\mathbf{Y}^T(\mathbf{I}-\mathbf{H})\mathbf{Y}/(n-p)$ 

$$(\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{Y} / (n - p)$$
  $\boldsymbol{D}$   $(\boldsymbol{X}^T \boldsymbol{X})^{-1} SSE / n$ 

$$(\boldsymbol{X}^T\boldsymbol{X})^{-1}SSE/r$$

### The hat matrix

Design matrix X has n rows and p linearly independent columns.  $\boldsymbol{H} = \boldsymbol{X}(\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T$ is called the hat-matrix.

Which of the following statements are NOT true?

$$\mathbf{A} \ \mathbf{H} = \mathbf{H}^T = \mathbf{H}^2 \quad \mathbf{B} \ \operatorname{rank}(\mathbf{H}) = \mathbf{p}$$

$$\mathbf{B} \quad \operatorname{rank}(\mathbf{H}) = \mathbf{p}$$

$$C HY = Y$$

$$\mathbf{C} \ \mathbf{H} \mathbf{Y} = \mathbf{Y} \qquad \mathbf{D} \ \mathbf{H} (\mathbf{I} - \mathbf{H}) = \mathbf{0}$$

### Inference about B

$$m{Y} = m{X}m{\beta} + m{\epsilon} \text{ where } m{\epsilon} \sim N_n(m{0}, \sigma^2 m{I})$$
 and  $\hat{m{\beta}} = (m{X}^Tm{X})^{-1}m{X}^Tm{Y}$ .

What are the properties of  $\hat{\beta}$ ?

- A Chi-squared distributed with n - pdegrees of freedom.
- C Multivariate normal with covariance matrix  $(\mathbf{I} - \mathbf{H})\sigma^2$ .
- B Chi-squared distributed with p degrees of freedom.
- D Multivariate normal with covariance matrix  $(\boldsymbol{X}^T\boldsymbol{X})^{-1}\sigma^2$ .

### Happiness=money+sex+love+work

```
Estimate Std. Error t value Pr(>|t|)
             0.009578
                       0.005213
                                   1.837
                                          0.0749
money
            -0.149008
                     0.418525
                                 -0.356
                                          0.7240
sex
                      0.295451
                                   6.496 1.97e-07
love
            1.919279
             0.476079 0.199389
                                  2.388
                                          0.0227
work
```

Which of the regression coefficient estimates has the largest estimated variance?

- A money B sex
- C love D work

### **Happiness**

```
Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.072081
                       0.852543
            0.009578
                                  1.837
                       0.005213
                                          0.0749
money
           -0.149008
                       0.418525
                                 -0.356
                                          0.7240
sex
            1.919279
love
                      0.295451
                                  6.496 1.97e-07
work
            0.476079 0.199389
                                  2.388
                                        0.0227
```

For which  $\beta_j$  would we reject the null hypothesis  $\beta_j = 0$  at significance level 1%?

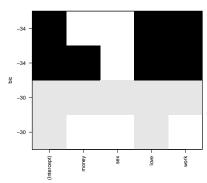
- A money B sex
- C love D work

### Happiness=money+sex+love+work

The  $R^2$  for the happiness-regression model is 71%. What does that mean?

- A The regression is significant for significance level 71%
- B The regression explains 71% of the variability in the data
- c The estimate for the variance  $\sigma^2$  is 0.71
- D The covariates have a correlation of 0.71

### Best model



Which model does the BIC criterion report to be the best?

A love+work

- B love
- C money+love+work
- **D** money+sex+love+work

# What is this plot used for? A Check residuals C Assess linearity D Find transform of response

### Answers

- 1. C: The normal equation  $(\boldsymbol{X}^T\boldsymbol{X})\hat{\boldsymbol{\beta}} = \boldsymbol{X}^T\boldsymbol{Y}$  is before you solve for  $\hat{\boldsymbol{\beta}}$ .
- 2. C: The hat matrix is symmetric and idempotent (so A is ok), and has rank p, but the reason for the name of the hat matrix is that is puts the hat on the  $\mathbf{Y}$  so  $\mathbf{H}\mathbf{Y} = \hat{\mathbf{Y}}$ . We know that for symmetric projection matrices the two matrices  $\mathbf{H}$  and  $(\mathbf{I} \mathbf{H})$  are orthogonal so the product must be zero.

### Correct?

Are you sure you want to read the correct answers? Maybe try first? The answers are explained on the next two slides.

### Answers

- 3. B: Since SSE has mean  $(n-p)\sigma^2$ , then SSE/(n-p) must be an unbiased estimator for  $\sigma^2$ . We know that  $(\mathbf{I} \mathbf{H})$  projects onto the space othogonal to the column space of the designmatrix, so that must have to do with SSE.
- 4. D: We know that linear combinations of multivariate normal random vectors are also multivariate normal (so the chisquare is not suitable). The residuals have (I H) as part of their covariance matrix, but  $\hat{\beta}$  has not.

### **Answers**

- 5. B: Sex has the largest estimated variance for regression estimate.
- 6. B:  $R^2$  gives the percent of variability explained.
- 7. C: only love is significant on level 1%, since this is the only p-value below 0.01 (last column).
- 8. A: love+work has smallest BIC.
- 9. D: Box-Cox plot used to find transformation of response.