**Multiple linear regression 1**
Morten Frydenberg ©
Section of Biostatistics, Aarhus Univ, Denmark

**Why do we need multiple linear regression.**

**An example**
  Interpretation of the parameters

**The general model**
  The assumptions.
  The parameters.
  Estimation.
  The distribution of the estimates
  Confidence intervals
  The F-test , R-squared

**Checking the model**
  Fitted values, residuals and leverage
  Extending the model

1

---

**Why do we need a multiple regression**

The simple linear regression model only models how the dependent variable, $y$, depend on **one** independent variable (covariate) , $x_1$.

We are often interested in **how** several independent variables, $x_1$ , $x_2$ ,..., $x_k$ , influence the dependent variable , $y$.

Sometimes we want to **adjust** the influence of some of the information, such as age and sex, before we look at the 'effect' of other variables.

2

---

**A multiple linear regression model**

We will here start by considering a **random** subsample consisting of 200 persons from the Frammingham study with focus on the baseline characteristics:
A **multiple** linear regression  model:

$$\ln(sbp) = \beta_0 + \beta_1 \cdot age + \beta_2 \cdot woman + \beta_3 \cdot \ln(bmi) + E$$

Where the **errors, $E$,** are assumed to be **independent** and **normal** with mean zero and standard deviation $\sigma$.

Note, that the variable *woman* is a **indicator** variable, that it is
  **one**   if the person is a **woman**
and
  **zero**  if the person  is a **man**.

3

---

**Interpretation of the coefficients 0 – the constant**

$$\ln(sbp) = \boxed{\beta_0} + \beta_1 \cdot age + \beta_2 \cdot woman + \beta_3 \cdot \ln(bmi) + E$$

The first coefficient (the constant term) is the **expected** $\ln(sbp)$ for

  a **man**          (that is ok!)

  $age$=0            ??????

  $bmi$=1 kg/m²      ??????       ( ln(1)=0 ).

As in the simple linear regression this is not of any interest.

But again we can control the interpretation, by choosing **relevant reference** values for $age$ and $bmi$.  E.g.

$$\ln(sbp) = \alpha_0 + \beta_1 \cdot (age - 45) + \beta_2 \cdot woman + \beta_3 \cdot \ln\left(\frac{bmi}{25}\right) + E$$

`age45`                          `lnBMI25`

4

**Interpretation of the coefficients 1**

$$\ln(sbp) = \beta_0 + \boxed{\beta_1} \cdot age + \beta_2 \cdot woman + \beta_3 \cdot \ln(bmi) + E$$

The **expected** $\ln(sbp)$ for a **man** with $bmi$=27 kg/m² is:

$$\beta_0 + \beta_1 \cdot age + \beta_3 \cdot \ln(27)$$

The **expected** $\ln(sbp)$ for another **man** with the same $bmi$, but 1.7 **year older**:

$$\beta_0 + \beta_1 \cdot (age + 1.7) + \beta_3 \cdot \ln(27)$$

The difference is: $1.7\beta_1$

We see that this difference

·**does not** depend on the $age$ of the first man.

·**does not** depend on the $bmi$ as long as it is the same for the two men.

·would be the same if the two persons were women.

5

**Interpretation of the coefficients 2**

$$\ln(sbp) = \beta_0 + \beta_1 \cdot age + \boxed{\beta_2} \cdot woman + \beta_3 \cdot \ln(bmi) + E$$

The **expected** $\ln(sbp)$ for a 50 year old **man** with $bmi$=27 kg/m² is:    $\beta_0 + \beta_1 \cdot 50 \qquad + \beta_3 \cdot \ln(27)$

The **expected** $\ln(sbp)$ for **woman** with the same $age$ and $bmi$

$$\beta_0 + \beta_1 \cdot 50 + \beta_2 \qquad + \beta_3 \cdot \ln(27)$$

The difference is: $\beta_2$

We see that this difference

·**does not** depend on the $age$ as long as it is the same for the two persons.

·**does not** depend on the $bmi$ as long as it is the same for the two persons.

6

**Interpretation of the coefficients 3**

$$\ln(sbp) = \beta_0 + \beta_1 \cdot age + \beta_2 \cdot woman + \boxed{\beta_3} \cdot \ln(bmi) + E$$

The **expected** $\ln(sbp)$ for a **woman** who is 50 year old:

$$\beta_0 + \beta_1 \cdot 50 + \beta_2 + \beta_3 \cdot \ln(bmi)$$

The **expected** $\ln(sbp)$ for another **woman** with the same age, but with a $bmi$ which is 10% higher:

$$\beta_0 + \beta_1 \cdot 50 + \beta_2 + \beta_3 \cdot \ln(1.1 \cdot bmi)$$

The difference $\quad \beta_3 \cdot \left[\ln(1.1 \cdot bmi) - \ln(bmi)\right] = \beta_3 \cdot \ln(1.1)$

We see that this difference

·**does not** depend on the $bmi$ of the first woman.

·**does not** depend on the $age$ as long as it is the same for the two women.

·would be the same if the two persons were **men**.

7

**Interpretation of the coefficients 4**

$$\ln(sbp) = \beta_0 + \beta_1 \cdot age + \beta_2 \cdot woman + \boxed{\beta_3} \cdot \ln(bmi) + E$$

$$\beta_3 \cdot \left[\ln(1.1 \cdot bmi) - \ln(bmi)\right] = \beta_3 \cdot \ln(1.1)$$

As the $bmi$ is introduced on the **log-scale**, then "differences " of this variable is measured **relatively**.

So comparing a pair of persons who **only differ** in bmi . One having $bmi$=25 kg/m² and the other $bmi$=27 kg/m² .

Then the expected difference in $\ln(sbp)$ is:

$$\beta_3 \cdot \ln\left(\frac{27}{25}\right) = \beta_3 \cdot 0.077$$

If the bmi's were 21 kg/m² and 23 kg/m² , then the expected difference in $\ln(sbp)$ would be:

$$\beta_3 \cdot \ln\left(\frac{23}{21}\right) = \beta_3 \cdot 0.091$$

8

## Interpretation of the coefficients 5

$$\boxed{\ln\left(sbp\right)} = \beta_0 + \beta_1 \cdot age + \beta_2 \cdot woman + \beta_3 \cdot \ln\left(bmi\right) + E$$

Taking the **expone
ntial** we get:

$$sbp = \gamma_0 \cdot \gamma_1^{age} \cdot \gamma_2^{woman} \cdot bmi^{\beta_3} \cdot \exp\left(E\right)$$

where $\gamma_0 = \exp\left(\beta_0\right)$, $\gamma_1 = \exp\left(\beta_1\right)$ and $\gamma_2 = \exp\left(\beta_2\right)$

That is a non-linear model on the *sbp* scale!

The error is **multiplicative**.

As **medians** are preserved by the exponential transformation then the estimates are measuring the **effects on the median** *sbp*.

**An example**: The age and bmi adjusted median sbp is a factor $\gamma_2$ higher for women compared to men.

9

## The multiple linear regression in general

$Y$                the **dependent** variable

$(x_1, x_2, \ldots, x_k)$      the **independent** variables.

$$Y = \beta_0 + \sum_{p=1}^{k} \beta_p \cdot x_p + E \qquad E \sim N\left(0, \sigma^2\right)$$

This model is based on the **assumptions**:

1. The **expected** value of $Y$ is   $\beta_0 + \sum_{p=1}^{k} \beta_p \cdot x_p$

2. The **unexplained** random deviations are **independent**.

3. The unexplained random deviations have the **same distributions**.

4. This distribution is **normal**.

10

## The multiple linear regression in general

$$Y = \beta_0 + \sum_{p=1}^{k} \beta_p \cdot x_p + E \qquad E \sim N\left(0, \sigma^2\right)$$

We see that the assumptions fall in **two parts**:

The **first concerning** the systematic part

and the three other which focus on the error, the unexplained random variation.

Before we turn to how one can check some of the assumptions, we will take a closer look at the first assumption.

The **expected** value of $Y$ is   $\beta_0 + \sum_{p=1}^{k} \beta_p \cdot x_p$

11

## The assumption of linearity

The **expected** value of $Y$ is   $\beta_0 + \sum_{p=1}^{k} \beta_p \cdot x_p$

This is based on three (sub) assumptions:

a. **Additivity**: The contribution from each of the independent variables are **added**.

b. **Proportionality**: The contribution from a independent variable is **proportional** to its value (with a factor $\beta$ )

c. **No effectmodification**: The contribution from one independent variables **is the same** whatever the values are for the other.

12

---

**The assumption of linearity**

The **expected** value of $Y$ is $\beta_0 + \sum_{p=1}^{k} \beta_p \cdot x_p$

If one consider two persons who differ with

$$\Delta x_1 \text{ in } x_1 , \Delta x_2 \text{ in } x_2 \dots \text{ and } \Delta x_k \text{ in } x_k$$

then the difference in the **expected** value of $Y$ is :

$$\sum_{p=1}^{k} \beta_p \cdot \Delta x_p$$

Again we see that the **contribution** for each of the explanatory variables:
  are **added**,
  are **proportional** to the difference
  and **does not dependent** of the differences in the other

13

---

**Estimation**

It is almost impossible to find the estimates by hand, but easy if you use a computer.

In Stata: `regress lnSBP age45 woman lnBMI25`

(Note first we have to generate `lnSBP`, `age45`, `woman` and `lnBMI25`)

```
  Source |       SS       df       MS              Number of obs =     200
---------+------------------------------           F(  3,   196) =   16.46
   Model | 1.05572698      3  .351908994           Prob > F      =  0.0000
Residual | 4.18969066    196  .021375973           R-squared     =  0.2013
---------+------------------------------           Adj R-squared =  0.1890
   Total | 5.24541764    199  .026358883           Root MSE      =  .14621

-------------------------------------------------------------------------
   lnSBP |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
---------+---------------------------------------------------------------
   woman |   .0036329   .0208905     0.17   0.862    -.0375662    .0448319
   age45 |   .0065384   .0012844     5.09   0.000     .0040053    .0090715
 lnBMI25 |   .2583399   .0758295     3.41   0.001     .1087934    .4078864
   _cons |   4.856592   .0154266   314.82   0.000     4.826169    4.887016
-------------------------------------------------------------------------
```

14

---

**Estimation**

The last part of the output: No CI for $\sigma$!
It should be calculated "by hand"

$\hat{\sigma}$
Root MSE = .14621

```
   lnSBP |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
---------+---------------------------------------------------------------
   woman |   .0036329   .0208905     0.17   0.862    -.0375662    .0448319
   age45 |   .0065384   .0012844     5.09   0.000     .0040053    .0090715
 lnBMI25 |   .2583399   .0758295     3.41   0.001     .1087934    .4078864
   _cons |   4.856592   .0154266   314.82   0.000     4.826169    4.887016
```

the $\hat{\beta}'s$   the se's        The CI 's
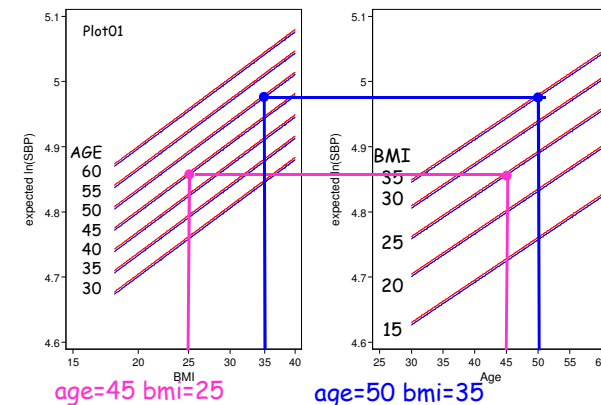
Test for $\beta_2 = 0$

The hypothesis: "no difference in $\ln(sbp)$ between men and women **adjusted** for age and bmi"

15

---

**Estimated systematic part**

$$\ln(sbp) = 4.8566 + 0.0065 \cdot (age - 45) + 0.0036 \cdot woman + 0.2583 \cdot \ln\left(\frac{bmi}{25}\right)$$



age=45 bmi=25     age=50 bmi=35

16

---

---

**Stata special – plotting response curves**

```
regress lnSBP age45 woman lnBMI25

------------------------------------------------------------------------------
  lnSBP |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
--------+---------------------------------------------------------------------
  woman |  .0036329    .0208905     0.17   0.862    -.0375662     .0448319
  age45 |  .0065384    .0012844     5.09   0.000     .0040053     .0090715
lnBMI25 |  .2583399    .0758295     3.41   0.001     .1087934     .4078864
  _cons |  4.856592    .0154266   314.82   0.000     4.826169     4.887016
------------------------------------------------------------------------------
```

After a regression commando, Stata leaves several
information in the memory of the computer for later use.

You can get a list by writing "`ereturn list`".
We have already used this feature in the calculation of the
confidence interval for $\sigma$.

Another example:

```
. display %12.7f _b[woman] %12.7f _se[woman]
0.0036329    0.0208905
```

17

---

**Stata special – plotting "response" curves**

I have made a Stata command that extracts the estimated
equations and the coefficients for later use.
The command file
    `regeq.ado`
and the small help file
    `regeq.sthlp`
should be place in your ado folder typically
    `c:\ado\personal`.

You can run the `regeq` command after any linear or logistic
regression estimation.
Here you get the output :
```
estimated equation
4.85659  +0.003632 * woman  +0.006538 * age45 +0.25834* lnBMI25
equation
b0 + b1 * woman + b2 * age45 + b3 * lnBMI25
```

That is, the estimated equation and the formula.

18

---

**Stata special – plotting "response" curves**

Furthermore the estimated coefficients are stored as "
global macros":

```
. macro list
b0:            4.856592269392944
b3:            .2583398993331004
b2:            .0065383788673611
b1:            .0036328605876014

S_E_depv:      lnSBP
S_E_cmd:       regress
.....
```

The global macros b0 to b3 contains the coefficients
and can be used in calculations.
If you want to use the estimated coefficient to age45,
then you just write $b2.

19

---

**Stata special – plotting "response" curves**

The expected log(SBP) for a 30 year old man with BMI=27
remember: Y= b0 + b1 * woman + b2 * age45 + b3 * lnBMI25

```
display $b0+$b1*0+$b2*(30-45) +$b3*ln(27/25)
4.7783987
```

You could also get this (with CI) using the lincom command:

```
 display ln(27/25)
.07696104

. lincom -15*age45 + .07696104*lnBMI25+_cons

 ( 1) - 15 age45 + .076961 lnBMI25 + _cons = 0

------------------------------------------------------------------------------
  lnSBP |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
--------+---------------------------------------------------------------------
    (1) |  4.778399    .0266891   179.04   0.000     4.725764     4.831033
------------------------------------------------------------------------------
```

20

---

**Remember:** `Y = b0+ b1 * woman + b2 * age45 + b3 * lnBMI25`

The expected log(SBP) for a 30 year old man as a function of the **BMI** is given as:

```
Y = b0 + b1 *0 +b2 * (30-45) + b3 * ln(BMI/25)
```

We can plot this by using the plot function in Stata:

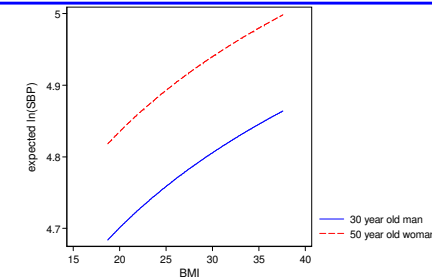```
twoway                                                            ///
( function Y=$b0 + $b1 * 0 +$b2 * (30-45) + $b3 * ln(x/25), range(bmi) ) ///
, legend(off) ytit("expected ln(SBP)") xtit("BMI") xlab( 15(5)40)
```



21

---

**Stata special – plotting response curves**

The expected log(SBP) for a 30 year old **man** and a 50 year old **woman** as a function of the BMI is given as:

```
twoway                                                            ///
( function Y=$b0 + $b1 * 0 + $b2 * (30-45) + $b3 * ln(x/25)        ///
                              , range(bmi) lco(blue) ) ///
( function Y=$b0 + $b1 * 1 + $b2 * (50-45) + $b3 * ln(x/25)        ///
                              , range(bmi) lco(red)  ) ///
, ytit("expected ln(SBP)") xtit("BMI") xlab( 15(5)40)             ///
 legend(label(1 "30 year old man") label(2 "50 year old woman"))
```



22

---

**Confidence intervals**

Just like in the simple regression we get :
  (except we have $n$-$k$-1 degrees of freedom).

**Exact** 95% confidence intervals , CI's, for $\beta_p$ is found from the estimates and standard errors

$$95\% \text{ CI for } \beta_p : \hat{\beta}_p \pm t_{n-k-1}^{0.975} \cdot \text{se}\left(\hat{\beta}_p\right)$$

Where $t_{n-k-1}^{0.975}$ is the upper 97.5 percentile in the t-distribution $n$-$k$-1 degrees of freedom.

These confidence intervals are found in the output.

A confidence interval for $\sigma$ can be found by `cisd`

Note that if $n$-$k$-1 is large then this percentile is close to 1.96 and one can use the **approximate confidence intervals:**
  Approx. 95% CI for $\beta_p : \hat{\beta}_p \pm 1.96 \cdot \text{se}\left(\hat{\beta}_p\right)$

23

---

**The ANOVA table and the F-test**

The first part of the output:

An **an**alysis **o**f **va**riance table dividing the variation in $y$ in two components: explained by the **model** (i.e. the **3** variables) and the **residual** (the rest)

```
  Source |       SS       df       MS              Number of obs =     200
---------+------------------------------           F( 3,   196) =   16.46
   Model | 1.05572698        3  .351908994         Prob > F      =  0.0000
Residual | 4.18969066      196  .021375973         R-squared     =  0.2013
---------+------------------------------           Adj R-squared =  0.1890
   Total | 5.24541764      199  .026358883         Root MSE      =  .14621
```

A $F$-test testing the hypothesis: "all $\beta$s (except $\beta_0$) is zero."

Here the test is highly significant: The model explains a statistically significant part of the variation in $y$!

24

---

**The F-test and R-squared**

The F- test calculated as:  $F = \dfrac{0.35519}{0.02138} = 16.46$

```
   Source |       SS       df       MS              Number of obs =      200
----------+------------------------------          F(  3,   196) =    16.46
    Model | 1.05572698        3 .351908994          Prob > F      =   0.0000
 Residual | 4.18969066      196 .021375973          R-squared     =   0.2013
----------+------------------------------          Adj R-squared =   0.1890
    Total | 5.24541764      199 .026358883          Root MSE      =   .14621
```

And under the hypothesis it follows an F-distribution with 3 and 196 degrees of freedom.

The **R-squared** is the amount of the total variation explained by the model(=1.0557/5.2454).

As this will **increase**, if we include more variables in the model, one can look at the **adjusted R-squared =**(0.02636-0.02138 )/0.02636

25

---

**Predicted values, residuals and leverages**

$$Y = \beta_0 + \sum_{p=1}^{k} \beta_p \cdot x_p + E \quad E \sim N\left(0, \sigma^2\right)$$

As in the simple linear regression one can find **predicted** values, **residuals**, **leverages** and **standardized residuals**:

Predicted value :  $\hat{y}_i = \hat{\beta}_0 + \sum_{p=1}^{k} \hat{\beta}_p \cdot x_{pi}$

Residual :  $r_i = y_i - \hat{y}_i = y_i - \left( \hat{\beta}_0 + \sum_{p=1}^{k} \hat{\beta}_p \cdot x_{pi} \right)$

Leverage :  $h_i = $ a complicated formula

Standardized-Residual :  $z_i = \dfrac{r_i}{\hat{\sigma}\sqrt{1 - h_i}}$

26

---

**Leverage**

Although the formula for the leverage is complicated, the **interpretation** of leverage is the same:

A **high leverage** indicates that the data point has **extreme** values of the explanatory variables and hence a **high influence** on the estimates.

27
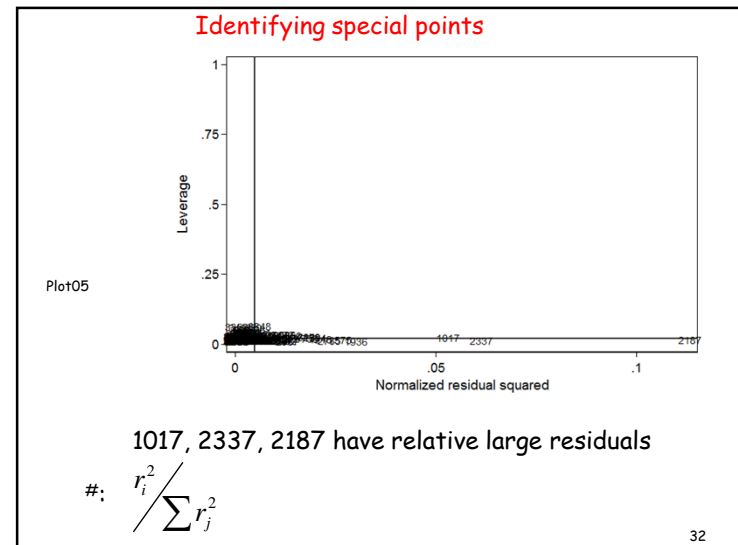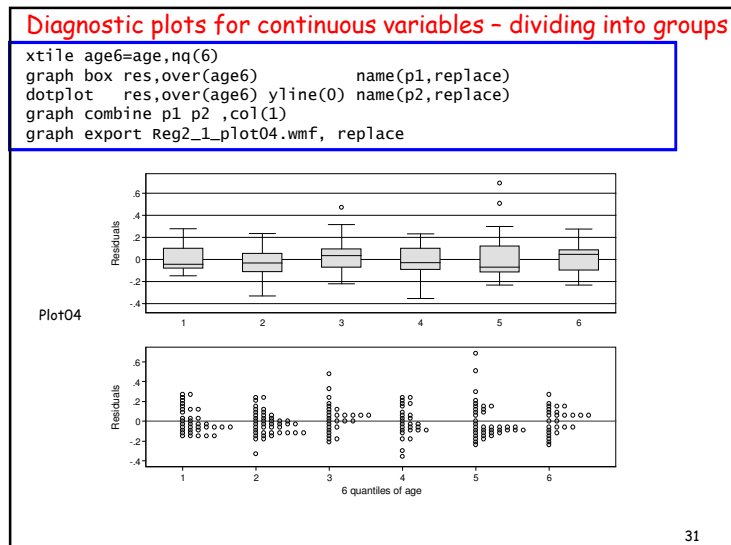
---

**Checking the model 1:**

As the model is much more complicated than the simple linear regression checking the model is also complicated

Again **assumption no. 2**: *the errors should be independent*, is mainly checked by considering how the data was collected.

The **distribution of the error** is checked by the same type of plot as for the simple linear regression.

•Plots of residuals versus **fitted**

•Plots of residuals versus **each of** the **explanatory** variables.

•Histogram and QQ-plot of the residuals.

28

---

### Slide 29

```
rvfplot        ,name(p1,replace)
rvpplot age45  ,name(p2,replace)
rvpplot lnBMI25 ,name(p3,replace)
rvpplot woman  ,name(p4,replace)
graph combine p1 p2 p3 p4
```

**r**esidual **v**ersus **f**itted
**r**esidual **v**ersus **p**redictor

Plot02

Not informative se next page

29

### Slide 30

Diagnostic plots for categorical variables – here woman (sex)

```
predict res if e(sample),res
qplot res,            over(sex) ///
    trscale(invnorm(@)) mco(blue red) msy(x Oh)    name(p3,replace)
graph box res      , over(sex)                     name(p4,replace)
graph combine p3 p4,col(1)
by woman: sum res
```

???

Plot03

sd=0.131          sd=0.157

30

### Slide 31

Diagnostic plots for continuous variables – dividing into groups

```
xtile age6=age,nq(6)
graph box res,over(age6)       name(p1,replace)
dotplot  res,over(age6) yline(0) name(p2,replace)
graph combine p1 p2 ,col(1)
graph export Reg2_1_plot04.wmf, replace
```

Plot04

31

### Slide 32

Identifying special points

Plot05

1017, 2337, 2187 have relative large residuals

$$\#: \quad r_i^2 \Big/ \sum r_j^2$$

32

**Checking the model 2: Independent errors ?**

**Assumption no. 2**: *the errors should be **independent**,* is mainly checked by considering **how the data was collected**.

The assumption is **violated** if

•some of the persons are **relatives** (and some are not) and the dependent variable have some **genetic** component.

•some of the persons were **measured** using one instrument and others with another.

•in general if the persons were sampled in clusters.

33

**Checking the model 3: Extending the model**

One should **also** try to check the validity of the linearity assumption that is the assumption of **additivity**, **proportionality** and **no effect modification** (no interaction).

It can be done by:

1. Introducing the explanatory variable in a **different scale**, e.g. adding $age^2$ or $\log(age)$ ….

2. Introducing the explanatory variable as a **categorical** variable instead e.g. use $age$ divided into **agegroups** instead as age in years.

3. Introducing **interactions** between some of the eplanatory variables.

4. ….

34