# Lecture 11

Logistic regression: continuous variables

# Summary: calculation of odds and probabilities

- Model with obesity only (*model.hyper2*)

Coefficients:

> coef(model.hyper2)

(Intercept)     obesityyes

 -1.6762466   0.7599559

        b

Odds of hypertension in baseline $= e^a$, and odds ratio for obese $= e^b$

> exp(coef(model.hyper2))

(Intercept)  obesityyes

  0.1870748   2.1381818

Probability of hypertension in non-obese:

p = odds/(1+ odds) =
> 0.1870748/1.1870748
[1] 0.1576259

Odds of hypertension in obese group: $e^{a+b}$

=p(baseline)*odds ratio     p(     )*

> 0.1870748*2.1381818

[1] 0.4001163

Probability of hypertension in the obese:

p = odds/(1+ odds) =
> 0.4001163/1.4001163
> [1] 0.2857736

if odds = p/(1-p)
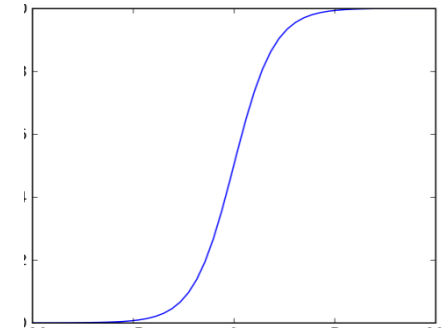
then

**p = odds/( 1+odds)**

# Continuous variables     (    )
=

- We can estimate the effects of a continuous variable (e.g. age) on the probability of an event (e.g. menarche) having occurred
  - = the *cumulative* probability of event having occurred

- How to interpret effect of continuous predictors on probability of event?

- Data: file *menar* (modified from *juul* in library *ISwR*)
  - girls aged 8 to 20 either had or haven't had menarche (i.e. they are either 'yes' or 'no' for menarche)
    - no: menarche=0      "    "     "    "
    - yes: menarche=1

  - logistic regression can estimate probability of menarche having occurred by age *x*

    x

```
model.menar <- glm(menarche~age,binomial, data=menar)
```

# Menarche and age

- Let's run a logistic regression of menarche against age

$a=$

$0$

- a=intercept not really meaningful in this case
  - log(odds) of menarche in people with 'no age' (age=0)

$b=\log( \quad / \quad )$

- b = log(odds ratio of menarche to no menarche)>1; P~0

- odds ratio= $e^b = e^{1.5176} = 4.56$
  - exposure to a 'unit of age' significantly increases odds of menarche to no menarche

" " /

# Menarche and age

```
> model.menar <- glm(menarche~age,binomial)
> summary(model.menar)
Call:
glm(formula = menarche ~ age, family = binomial)     b

Deviance Residuals:
    Min       1Q     Median       3Q       Max
-2.32759  -0.18998   0.01253   0.12132   2.45922
Coefficients:
              Estimate   Std. Error   z value   Pr(>|z|)
(Intercept)   -20.0132     2.0284      -9.867    <2e-16 ***
age             1.5173     0.1544       9.829    <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
(Dispersion parameter for binomial family taken to be 1)
    Null deviance: 719.39  on 518  degrees of freedom
Residual deviance: 200.66  on 517  degrees of freedom
AIC: 204.66
Number of Fisher Scoring iterations: 7                    4.56
```

( ) ( )

- In the case of continuous variable *age*, **odds ratio is change in odds of event (menarche) when age increases by 1 unit (i.e. per year)**

"                    " ( )

- *b* measures the effect of exposure to 'one unit of age' (i.e. one year)
  - it is the log of odds ratio of the event occurring at age N+1 to event occurring at age N

  n+1                    n

- =odds of menarche to no menarche increase by 4.56 per year

4.56

- 4.56 seems to be a large number, but don't forget that the odds start at ~ 0
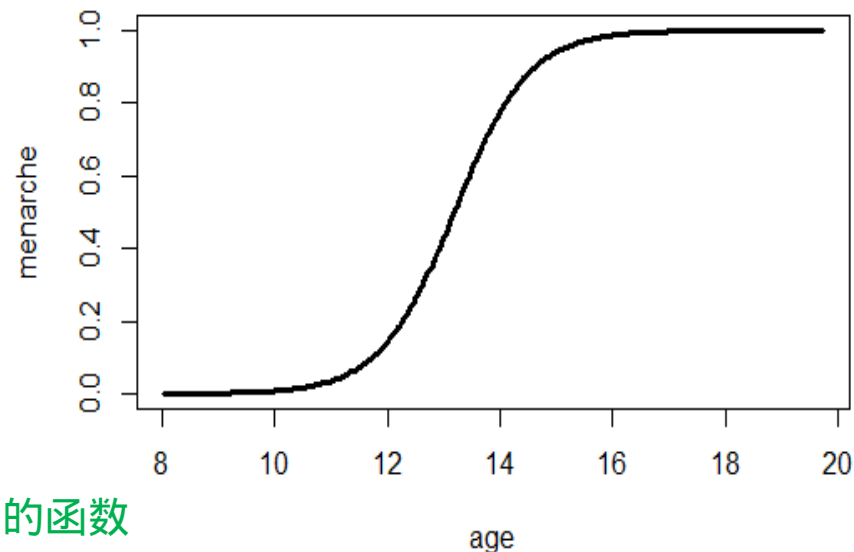
0

# Predicted probabilities

<span style="color:green">(     )</span>
<span style="color:green">odd</span>

- For categorical predictors like smoking (yes or no), we only need to calculate odds and probabilities for two groups: baseline and exposure

<span style="color:green">x    y</span>

- But for a continuous variable, there is a range of *x* values and *y* probabilities (e.g. probabilities of menarche as a function of ages from 8 to 20)

- To predict probabilities (= $\frac{1}{1+e^{-(a+bX)}}$ ) of menarche for all ages from 8 to 20:

  - function *predict*    <span style="color:green">predit</span>
  - add argument *type="response"*
    - otherwise *predict* returns logit values    <span style="color:green">predit    log</span>

- Saving probabilities in vector *prob*:    <span style="color:green">prob</span>
  <span style="color:#2e74b5">> probs <- predict(model.menar, type="response")</span>

- Plotting probability of menarche by age
  <span style="color:#2e74b5">> plot(probs~age, data=menar, pch=16, ylab="menarche")</span>

- To predict probabilities at a given point, use
  <span style="color:#2e74b5">> predict(your model, data.frame(X = value), type= "response")</span>

# Median age at event occurrence

- Median age at menarche = age where probability of menarche having occurred is 50%, or p=0.5

But when p=0.5:

- odds $= \dfrac{p}{1-p} = 1$

If odds = 1

then

- log(odds) = 0

# Median age at event occurrence

- Setting log(odds) = logit = 0, we can calculate
  age X at menarche from the equation:
  - logit p = a + bX

  - 0 = a + b(age)
  - 0 = -20.013 + 1.5173(age)
  - **median age = 13.19 years**

Exercise

- Still using file *menar*, run a logistic model using *igf1* (insulin-growth-like factor 1) as predictor of menarche

1. Interpret a and b

2. What is the probability of menarche for someone with igf1=500?

Tips:

- Estimate logit p ($=f$=a + bX) when *igf1* = 500
- Then use logistic function (p = 1/(1 +exp(-f)) to convert logit p into probability
- Confirm with this code:

> predict(model.menar, data.frame(igf1 = 500), type= "response")

# Break

# Categorical variables with > 2 levels

(1   0)

- Some categorical predictors have two levels (smoker=1, non-smoker=0), but others have more levels (month, location etc.)

- We can still run a logistic regression with those variables, but interpretation slightly changes again

- Example: month in *infant* dataset

2

- When predictor has >2 levels
  - first level is taken as baseline     `month=1,`          `as.factor(month)`
    - =January (coded as month=1 and then entered as *as.factor(month)*)
  - all other levels (=other months) are compared to the first *on an individual basis*, but not to each other
    - =February vs. January, March vs. January etc.
    - but not March vs. February etc.

# Categorical variables with > 2 levels

0        1

- Outcome: variable *healthy* (0=undernourished, 1=healthy)

- Predictor: *month* (birth month)
  - 9 levels (January=1 to September=9)                                        "    "
  - check file: month is a factor (although coded as numbers)
    - <mark>always make sure your categorical predictor is set as factor or character, NEVER numeric</mark>
    - <mark>otherwise you get a 'continuous' predictor (like age in the menarche analysis!)</mark>
  - "1" is baseline level "January", "2" is level "February"

> *model.infant <- glm(healthy ~ month,binomial, data=infant)*

# Categorical variables with > 2 discrete levels

```
> model.infant <- glm(healthy ~ as.factor(month),binomial, data=infant)
> summary(model.infant)
Coefficients:
```

|  | Estimate | Std. Error | z value | Pr(>\|z\|) |
|---|---|---|---|---|
| (Intercept) | -2.18407 | 0.25582 | -8.537 | <2e-16 *** |
| month2 | -0.06723 | 0.36678 | -0.183 | 0.8546 |
| month3 | -0.21383 | 0.37864 | -0.565 | 0.5723 |
| month4 | -0.81167 | 0.44350 | -1.830 | 0.0672 . |
| month5 | -0.81167 | 0.44350 | -1.830 | 0.0672 . |
| month6 | -1.28167 | 0.52124 | -2.459 | 0.0139 * |
| month7 | -1.25635 | 0.52140 | -2.410 | 0.0160 * |
| month8 | -0.97293 | 0.48909 | -1.989 | 0.0467 * |
| month9 | -1.14814 | 0.52210 | -2.199 | 0.0279 * |

```
   Null deviance: 642.54  on 1457  degrees of freedom
Residual deviance: 623.81  on 1449  degrees of freedom

AIC: 641.81
```

6-9      1   (      )                              b       <0

- Only months 6 to 9 (June to September) significantly differ from January (baseline)
  - coefficients b (log of odds ratio) are significantly <0

- Let's compare January and June:
- Exponentiating coefficients:

  >exp(coef(model.infant))

- Odds in January = exp(a)
  0.1125828

- Odds ratio June to January = odds ratio month 6 to month 1
  0.2775735

- Odds of being healthy in June:

=0.1125*0.277535

=0.03123487

# Categorical variables with > 2 discrete levels

- Remember:

    p

- When there are more than two levels in predictor, coefficients and P values reflect comparisons between each exposure level and baseline
  - Each month compared to January

- **But there is no comparison between exposure levels**
  - We know nothing directly about the difference between April and May

    4

- If we wanted to know about April vs. May:
  - Change baseline to April = 4 (with function *relevel*), then we would obtain a coefficient for May vs. April

```
levels(infant$month)                          levels(infant$month) <- c("1", "2", "3", "4", "5", "6", "7", "8",
infant$month <- relevel(infant$month, ref = "6")    "9")
        baseline
```

**Exercise**

- Re-run *model.infant,* but this time use month=9 (September) as baseline

(see code file)

```
infant$month <- relevel(infant$month, ref = "9")
        baseline
```

# Break

# Interactions

- Interactions occur when the effects of factors are not independent (they are not additive but multiplicative)

- Positive interaction:
  - Drug A causes small increase in odds of heart attack
  - Drug B causes small increase in odds of heart attack
  - But people taking drugs A and B have a large increase in odds of heart attack =positive interaction between A and B: their effects are stronger when combined

  =A   B

- Negative interaction:
  - Drug A causes increase in odds of heart attack
  - Drug B causes increase in odds of heart attack
  - But people taking drugs A and B show no increase in odds of heart attack =negative interaction between A and B: their effects are cancelled or reduced when combined

  =

# Interactions

- Interaction occurs if     1     2             1   2

  interaction

  - factors 1 and 2 are present in the same individual
  - their joint effect is different from the additive effects of 1 and 2 put together.

Example:

- Exposure to factor A doubles odds of an outcome
- Exposure to factor B also doubles odds of the outcome

What to expect from exposure to both A and B?

=(exposure to A) x (exposure to B) = 2 x 2 = 4 times the odds

=*additive* effect of A and B    A    B

      AB                4              AB

- But if exposure to both A and B results in odds different from 4 (the additive effect), A and B are *interacting*
  - joint effect > 4: *positive* interaction          >4
  - joint effect < 4: *negative* interaction          <4

# Interactions

- File *evans*: Evans county study of factors leading to coronary heart disease

- Let's focus son the effects of three factors:
    - *age*
    - *cat* (catecholamine levels) and
    - *chl* (cholesterol levels)

on the probability of coronary heart disease (*chd)*

<span style="color:green">R                  X1   X2                                     X1*X2</span>

- In R, to include all possible interactions between variables X1 and X2,
    - multiply them: Y ~ X1*X2
- X1*X2 is expanded into: <span style="color:green">X1*X2           X1+X2+X1X2            " "</span>
    - X1 + X2 + <span style="color:red">X1:X2</span>
    - interactions are represented by ":"

Our model with interactions is then:

```
> model.chd <- glm(chd ~ age*cat*chl, binomial, data=evans)
```

# Interactions: baseline

```
> model.chd <- glm(chd~age*cat*chl,binomial,data=evans)
> summary(model.chd)


Deviance Residuals:
     Min        1Q     Median        3Q        Max
 -2.3268   -1.1954    0.8112    1.1154     1.6543
 Coefficients:

                Estimate Std. Error z value Pr(>|z|)
(Intercept) -5.6564566  3.1060236   -1.821  0.06859 .
age          0.0929091  0.0589832    1.575  0.11522
cat         28.3920812 10.9032473    2.604  0.00921 **
chl          0.0223684  0.0140188    1.596  0.11058
age:cat     -0.5281193  0.1861421   -2.837  0.00455 **
age:chl     -0.0003483  0.0002650   -1.314  0.18873
cat:chl     -0.1302252  0.0546123   -2.385  0.01710 *
age:cat:chl  0.0024763  0.0009319    2.657  0.00788 **


Null deviance: 840.31  on 608  degrees of freedom
Residual deviance: 809.76  on 601  degrees of freedom
AIC: 825.76
```

- First let's understand how to read interactions:

cat                                b=28.39

- *cat* increases odds of coronary disease
  - b = 28.39
- *age* does not have a significant effect
- however, *age* and *cat* show a significant and b=-0.52 negative interaction
  - b= -0.52

# Example: interaction *age:cat*

```
> model.chd <- glm(chd~age*cat*chl,binomial,data=evans)
> summary(model.chd)
```

Deviance Residuals:
```
     Min        1Q    Median        3Q       Max
 -2.3268   -1.1954    0.8112    1.1154    1.6543
```
 Coefficients:
```
               Estimate  Std. Error  z value  Pr(>|z|)
(Intercept)  -5.6564566   3.1060236   -1.821   0.06859 .
age           0.0929091   0.0589832    1.575   0.11522
cat          28.3920812  10.9032473    2.604   0.00921 **
chl           0.0223684   0.0140188    1.596   0.11058
age:cat      -0.5281193   0.1861421   -2.837   0.00455 **
age:chl      -0.0003483   0.0002650   -1.314   0.18873
cat:chl      -0.1302252   0.0546123   -2.385   0.01710 *
age:cat:chl   0.0024763   0.0009319    2.657   0.00788 **
```

```
Null deviance: 840.31  on 608  degrees of freedom
Residual deviance: 809.76  on 601  degrees of freedom
AIC: 825.76
```

cat
b

a cat

a        a+b(age)+b(cat)

age:cat

age  cat

- **Important:**

- Effect of *cat* is measured relative to baseline
  - to calculate it, we use intercept *a* and *cat* coefficient *b*

- By contrast, the interaction effect *age:cat* now has a different baseline: a person where *age* and *cat* had their additive effects calculated

- In other words, baseline for interaction term is not intercept a, but
  - a + b(age) + b(cat)

# Example: interaction *age:cat*

```
> model.chd <- glm(chd~age*cat*chl,binomial,data=evans)
> summary(model.chd)


Deviance Residuals:
    Min        1Q     Median        3Q        Max
-2.3268   -1.1954     0.8112    1.1154     1.6543
 Coefficients:
                Estimate Std. Error z value Pr(>|z|)
(Intercept) -5.6564566  3.1060236   -1.821  0.06859 .
age          0.0929091  0.0589832    1.575  0.11522
cat         28.3920812 10.9032473    2.604  0.00921 **
chl          0.0223684  0.0140188    1.596  0.11058
age:cat     -0.5281193  0.1861421   -2.837  0.00455 **
age:chl     -0.0003483  0.0002650   -1.314  0.18873
cat:chl     -0.1302252  0.0546123   -2.385  0.01710 *
age:cat:chl  0.0024763  0.0009319    2.657  0.00788 **

Null deviance: 840.31  on 608  degrees of freedom
Residual deviance: 809.76  on 601  degrees of freedom
AIC: 825.76
```

A:B   3                    4

4              3        =7

- **Important:**

- Back to our hypothetical example: if
  - effect of A = 2
  - effect of B = 2
  addictive effects: 2 x 2 = 4

- then if interaction A:B = 3
  - total effect = 2 x 2 x 3 = 12

- A:B is 3; less than the additive effects of 4
  - this doesn't mean a negative interaction! Not a reduction!
  - interaction term is an effect of 3 on top of the additive effects of 4 = total effect of 7

# Example: interaction *age:cat*

```
> model.chd <- glm(chd~age*cat*chl,binomial,data=evans)
> summary(model.chd)
```

```
Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.3268  -1.1954   0.8112   1.1154   1.6543
 Coefficients:

              Estimate Std. Error z value Pr(>|z|)
(Intercept) -5.6564566  3.1060236  -1.821  0.06859 .
age          0.0929091  0.0589832   1.575  0.11522
cat         28.3920812 10.9032473   2.604  0.00921 **
chl          0.0223684  0.0140188   1.596  0.11058
age:cat     -0.5281193  0.1861421  -2.837  0.00455 **
age:chl     -0.0003483  0.0002650  -1.314  0.18873
cat:chl     -0.1302252  0.0546123  -2.385  0.01710 *
age:cat:chl  0.0024763  0.0009319   2.657  0.00788 **

Null deviance: 840.31  on 608  degrees of freedom
Residual deviance: 809.76  on 601  degrees of freedom
AIC: 825.76
```

- **Important**:

- In log form (where multiplications become additions):
  - effect of A = 0.5
  - effect of B = 0.5
  - additive effects: A + B = 1

- A:B= 0.1: positive interaction
  - Total effect = 0.5 + 0.5 + 0.1 = 1.1

# Example: interaction *age:cat*

```
> model.chd <- glm(chd~age*cat*chl,binomial,data=evans)
> summary(model.chd)
```

```
Deviance Residuals:
    Min      1Q    Median      3Q      Max
-2.3268  -1.1954   0.8112   1.1154   1.6543
 Coefficients:

                Estimate Std. Error z value Pr(>|z|)
(Intercept) -5.6564566  3.1060236  -1.821  0.06859 .
age          0.0929091  0.0589832   1.575  0.11522
cat         28.3920812 10.9032473   2.604  0.00921 **
chl          0.0223684  0.0140188   1.596  0.11058
age:cat     -0.5281193  0.1861421  -2.837  0.00455 **
age:chl     -0.0003483  0.0002650  -1.314  0.18873
cat:chl     -0.1302252  0.0546123  -2.385  0.01710 *
age:cat:chl  0.0024763  0.0009319   2.657  0.00788 **
```

```
Null deviance: 840.31  on 608  degrees of freedom
Residual deviance: 809.76  on 601  degrees of freedom
AIC: 825.76
```

chd                cat

- Now back to *chd* model

- *age:cat* interaction: b=-0.528 = negative interaction;

- But total effect on odds of chd is still positive:

$age + cat + age:cat =$

$0.09 + 28.39 - 0.528 = 27.952$

- In summary, the effects of *cat* and *age* are reduced when they act together
  - but they still increase odds of *chd* because of the strong additive effect of *cat*

# Break

# Model optimisation and the Hierarchy Principle

- When regressions return non-significant terms, we must optimise models to obtain a *minimal adequate model*

- Optimisation must follow a hierarchical procedure: higher-order interactions are tested first, individual factors last
  - if an interaction is significant, all lower level interactions and single terms must be kept *even if they are not significant*

  XA: X2: X3                                      X1: X2  X1: X3  X2: X3        X1  X2  X3

- EXAMPLE: if an interaction *X1:X2:X3*, is significant, final model must also include
  - all their lower interactions *X1:X2, X1:X3, X2:X3*
  - its single terms *X1, X2, X3*

                                          =

- Why? Because we need lower levels as baselines to estimate total effects
  = same reason we need odds in baseline to estimate odds in exposure group

# The Hierarchy Principle

age:cat:chl

```
> model.chd <- glm(chd~age*cat*chl,binomial,data=evans)
> summary(model.chd)


Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.3268  -1.1954   0.8112   1.1154   1.6543
 Coefficients:
               Estimate Std. Error z value Pr(>|z|)
(Intercept) -5.6564566  3.1060236  -1.821  0.06859 .
age          0.0929091  0.0589832   1.575  0.11522
cat         28.3920812 10.9032473   2.604  0.00921 **
chl          0.0223684  0.0140188   1.596  0.11058
age:cat     -0.5281193  0.1861421  -2.837  0.00455 **
age:chl     -0.0003483  0.0002650  -1.314  0.18873
cat:chl     -0.1302252  0.0546123  -2.385  0.01710 *
age:cat:chl  0.0024763  0.0009319   2.657  0.00788 **

Null deviance: 840.31  on 608  degrees of freedom
Residual deviance: 809.76  on 601  degrees of freedom
AIC: 825.76
```

- Back to our model: triple interaction *age:cat:chl* is significant
  - (unfortunately!)

- Therefore, if we optimised this model, we would not be able to discard any non-significant terms

# The Hierarchy Principle

```
> model.chd <- glm(chd~age*cat*chl,binomial,data=evans)
> summary(model.chd)
```

```
Deviance Residuals:
    Min       1Q     Median      3Q       Max
-2.3268   -1.1954    0.8112    1.1154    1.6543
 Coefficients:
```

| | Estimate | Std. Error | z value | Pr(>\|z\|) | |
|---|---|---|---|---|---|
| (Intercept) | -5.6564566 | 3.1060236 | -1.821 | 0.06859 | . |
| age | 0.0929091 | 0.0589832 | 1.575 | 0.11522 | |
| cat | 28.3920812 | 10.9032473 | 2.604 | 0.00921 | ** |
| chl | 0.0223684 | 0.0140188 | 1.596 | 0.11058 | |
| age:cat | -0.5281193 | 0.1861421 | -2.837 | 0.00455 | ** |
| age:chl | -0.0003483 | 0.0002650 | -1.314 | 0.18873 | |
| cat:chl | -0.1302252 | 0.0546123 | -2.385 | 0.01710 | * |
| age:cat:chl | 0.0024763 | 0.0009319 | 2.657 | 0.00788 | ** |

```
Null deviance: 840.31  on 608  degrees of freedom
Residual deviance: 809.76  on 601  degrees of freedom
AIC: 825.76
```

- When calculating the total effect of *age* and *cat and chl* on chd, we need:
  - intercept a
  - $b_1$ for *age*
  - $b_2$ for cat
  - $b_3$ for chl
  - $b_{12}$ for age:cat
  - $b_{13}$ for age:chl
  - $b_{23}$ for cat:chl
  - $b_{123}$ for age:cat:chl

logit $chd$ = a + $b_1(age)$ + $b_2(cat)$ + $b_{12}(age:cat)$
= a + $b_1$*X1 + $b_2$*X2 + $b_{12}$*X1*X2…etc

Example: for subject id=51, *age*=56, *cat*=1, and chl=201, logit chd
= -5.66 + 0.09*(56) + 28.4*(1) + 0.02*(201)…etc

step()

# Model optimisation: function *step( )*

=

- We optimise models (=discarding unnecessary variables) using the function <mark>*step*</mark>, which follows the hierarchical principle

  ANOVA
  - *step* runs ANOVAs comparing models with and without a given term; if absence of term does not significantly change the model, it should be eliminated

  log-likehood  aic
- Optimisation is based on the log-likelihood and AIC (Akaike information criterion, a function both of significance and number of variables in a model)
  - AIC comparisons only work for models that are hierarchically organised, i.e. when variables in model 1 are a subset of variables in model 2

AIC                                     1              2

- In practical terms:              log-likelihood      AIC
  - we eliminate a variable if this increases log-likelihood and reduces AIC
  - we test variables according to the hierarchical principle (higher-interactions first, single terms last)

# Example: menarche

Call:
glm(formula = menarche ~ age * igf1, family = binomial, data = menar)

Deviance Residuals:
```
    Min       1Q   Median       3Q      Max
-2.41072  -0.03565   0.01761   0.09315   2.60345
```

Coefficients:

|             | Estimate   | Std. Error | z value | Pr(>\|z\|) |    |
|-------------|------------|------------|---------|-----------|----|
| (Intercept) | -3.162e+01 | 1.021e+01  | -3.096  | 0.00196   | ** |
| age         | 2.100e+00  | 7.633e-01  | 2.752   | 0.00593   | ** |
| igf1        | 1.794e-02  | 1.996e-02  | 0.899   | 0.36886   |    |
| age:igf1    | 7.769e-04  | 1.522e-03  | -0.511  | 0.60962   |    |

---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
    Null deviance: 564.83  on 410  degrees of freedom
Residual deviance: 111.44  on 407  degrees of freedom
  (108 observations deleted due to missingness)
AIC: 119.44
Number of Fisher Scoring iterations: 8

- Now let us run a model of menarche on *age* and *igf1* with interactions,

- Neither *igf1* nor *age:igf1* is significant

- They should be both eliminated from the optimised model

# Example

> summary(step(model.menar2)) **

Call:
glm(formula = menarche ~ age + igf1, family = binomial, data = menar)

Deviance Residuals:
| Min | 1Q | Median | 3Q | Max |
|---|---|---|---|---|
| -2.43884 | -0.04581 | 0.01931 | 0.09146 | 2.58392 |

Coefficients:

| | Estimate | Std. Error | z value | Pr(>|z|) | |
|---|---|---|---|---|---|
| (Intercept) | -26.887594 | 3.650184 | -7.366 | 1.76e-13 | *** |
| age | 1.739611 | 0.238325 | 7.299 | 2.89e-13 | *** |
| igf1 | 0.007814 | 0.001880 | 4.157 | 3.23e-05 | *** |

p

igf1                              age:igf1

        age    igf1

step

- We optimise model with function *step*

** we obtain more output than this; see code

IMPORTANT
- When you remove a term, P values of all remaining terms may change!
- Therefore you can only tell which term is significant after optimisation
  - in the full model (previous slide), *igf1* was not significant, but once *age:igf1* was removed, it became significant in the final model (this slide)
- Only use coefficients from the final, optimised model!
- In this example, only age and igf1 are significant

```
model.chd2 <- glm(chd~age*cat,binomial, data=evans)
summary(model.chd2)
summary(step(model.chd2))
```

**Exercise**

- Run a logistic regression of *chd* on *age*, *cat* and their interaction.

- What is the optimal model and its AIC?