# Lecture 10

Logistic regression: categorical variables

# Logistic regression

( ) /

- Logistic regression makes the transition from the basic (least-square-based) *general linear model* to the intermediate/advanced *generalised linear model*

- The generalised linear model extends linear regression models to variables not normally distributed, and to non-linear relationships

- For example, we may want to use regression techniques to predict *binary* responses:
  - we may want to predict the probability that someone is dead or alive, voted or did not vote in the last election etc. as a function of other variables (age, smoking, income etc.)

- In other words, we still want to use a regression:

Probability of binary outcome $= a + b_1X_1 + b_2X_2\ldots + b_nX_n = a + \Sigma b_iX_i$

with
a = intercept
$b_i$ = regression coefficients
$X_i$ = independent variables (continuous or categorical) ( )

# Applications

p ( yes/no)

- Logistic regression is a classification method used almost universally
  - it predicts whether an outcome happens or not (binary outcome yes/no) with a probability $p$

- It is frequently applied to predict binary outcomes (yes or no)
  - business: costumer choice (purchasing, being late for bills etc.)
  - medicine, pubic health (will develop a condition etc)
  - insurance (risk of event, credit decisions)
  - etc.

- Logistic regression is closely linked to neural networks and machine learning
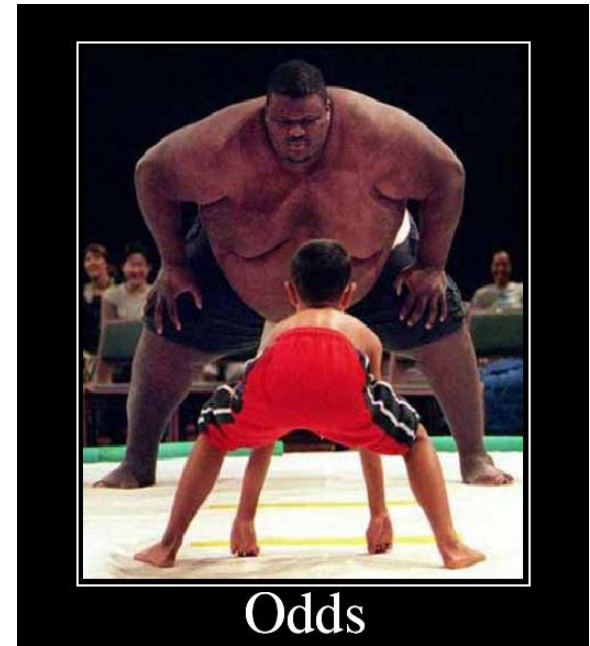
# Odds and log(odds)

- To understand logistic regressions, first we need to understand *odds* and *odds ratios*

  odd

- <mark>Important</mark>: odds are not the same as the *probability* of the event!

  $$\text{odds of event} = \frac{probability\ of\ event\ occurring}{probability\ of\ event\ not\ occurring}$$



Odds

# Odds and log(odds)

- Example: what is the *probability* of your birthday falling on a weekday this year?

  - probability of weekday=5/7=0.71         = p

$$\text{Odds of a weekday} = \frac{\text{probability of weekday}}{\text{probability of weekend day}}$$

  - odds of weekday = (5/7) / (2/7) = 5/2 = 2.5     = p/(1-p)

  - ln(odds of weekday) = log(2.5) = 0.91      = log(p/(1-p))

```
p        0.5      odd      1
log(odd)           1
```

```
p        0.5      odd      1
log(odd)
```

- And the probability of the non-event, i.e. weekend day?
  - probability of weekend day = 2/7=0.29          =1-p
  - odds of weekend day = 2/5 =0.4              =(1-p)/p
  - ln(odds of weekend day)= –0.91           = ln((1-p)/p)

Exercises

Calculate:

- Tossing a fair coin:
  - Probability of heads? `0.5`
  - Odds of heads? `1`
  - Odds of tails? `1`
    `0`
  - log(odds of heads)

- Now throwing a die:
  - Probability of 1? `1/6`
  - Odds of 1? `1/5`
    `5`
  - Odds of *not 1*? `log(0.2)`
  - log(odds of 1)?

# Odds ratio

- Now imagine you must choose between betting on coins (bet on 'heads') or dice (bet on '1'); what are the odds of winning in each?
  - odds of heads = 1/1 = 1
  - odds of a 1 = 1/5 =0.2


- So it is easier to win a coin toss; but how much easier?


- We can calculate the odds ratio of success in coins vs. dice


- Odds ratio $= \dfrac{odds\ of\ heads}{odds\ of\ 1} = \dfrac{1}{0.2} = 5$


- This means you are 5 times more likely to win by tossing a coin than throwing a die

# Summary

p    0   1

- probability p is always between 0 and 1

        0

- odds and odds ratio: 0 to $+\infty$

- log(odds) and log(odds ratio): $-\infty$ to $+\infty$

# Odd and probabilities

- If odds = p/(1-p), then:


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

Therefore

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

# Break
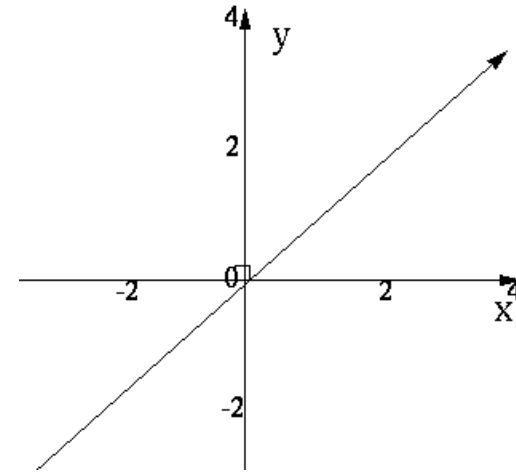
# Logistic function

- Back to logistic regression: we want to use a regression model to calculate the probability of binary events (dead/alive, head/tail etc.) from a set of predictors:

$$y = a + b_1X_1 + b_2X_2\ldots+ b_nX_n = a+\Sigma b_iX_i$$

- Problem:     $y$                                    $p$     $01$
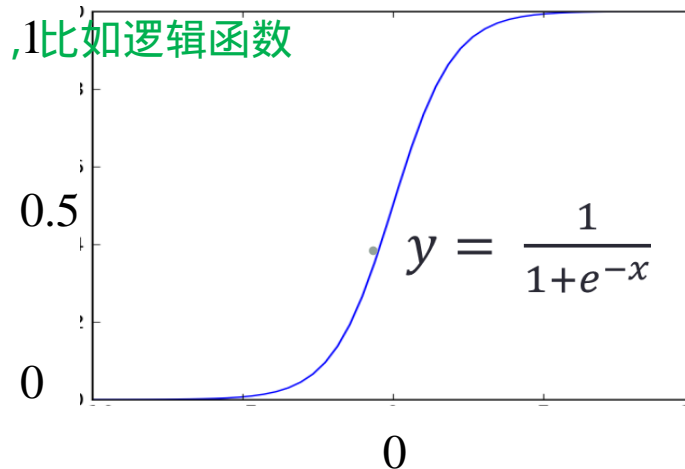  - linear regression predicts $y$ between $-\infty$ and $+\infty$
  - but probability is always between 0 and 1!

- Solution:
  - we want our probabilities to be estimated by a model such as the logistic function
  - why? Because whatever x, it will always return a value between 0 and 1

$x$
$01$
$$y = \frac{1}{1+e^{-x}} = \frac{1}{1+e^{1/x}}$$

$$y = \frac{1}{1+e^{-x}}$$

# Link function: Logit

- We need a *link function f* to be the *x* in the logistic function $y = \dfrac{1}{1+e^{-x}}$ and calculate *y* as probability p:

$$p = \frac{1}{1+e^{-f}} = \frac{1}{1+\frac{1}{e^f}}$$

- But $p = \dfrac{1}{1+\frac{1}{odds}}$

- Therefore $e^f = odds$; or $f = \log(odds)$

- The link function *f* is called **logit p**:

$$f = \text{logit } p = \log(odds) = \log\left(\frac{p}{1-p}\right)$$

p     f    x    y

---

another derivation:

- If we want $p = \dfrac{1}{1+e^{-f}}$ , then:

- $p = \dfrac{e^f}{e^f+1}$

- $p(e^f+1) = e^f$

- $pe^f + p = e^f$

- $p = e^f - pe^f$

- $p = e^f(1 - p)$

- $e^f = \dfrac{p}{1-p}$

- $\log(e^f) = \log\left(\dfrac{p}{1-p}\right)$

- $\boldsymbol{f = \log\left(\dfrac{p}{1-p}\right)}$

- note: logit is always natural log (i.e. log on base e=2.71)

# f = logit = log(odds of event)

- *f* = logit or log(odds) range from $-\infty$ to $+\infty$
  - therefore we can predict logits with a linear regression on our $X_1$, $X_2$ etc. variables
    X1X2                    logits
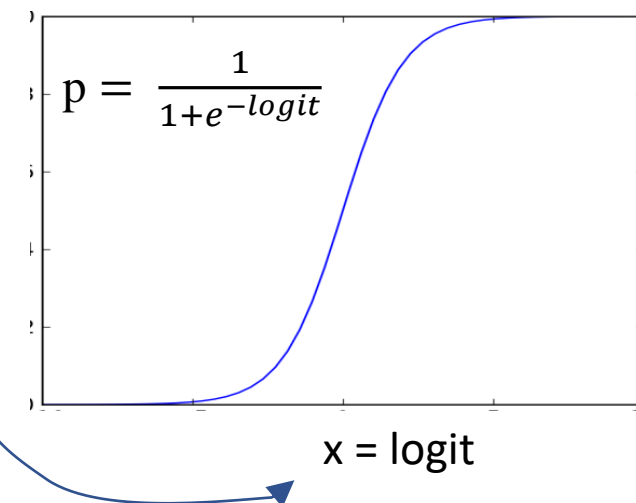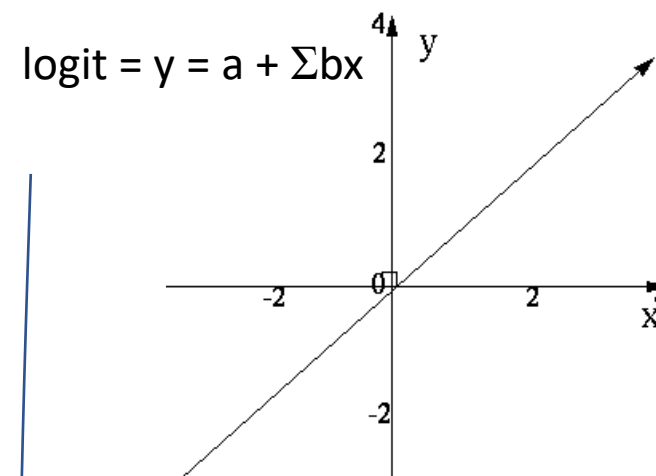
- The *logistic regression model* is thus

$$f = \text{logit} = \log\left(\frac{p}{1-p}\right) = \log(\text{odds}) = a + \Sigma b_i X_i$$

    logit        0   1

- Then we use the logits in the logistic regression to obtain probabilities between 0 and 1 !

$$p = \frac{1}{1+e^{-f}} = \frac{1}{1+e^{-\log(odds)}}$$

logit = y = a + Σbx



$$p = \frac{1}{1+e^{-logit}}$$

x = logit

# Fitting logistic regression

a    b        MML

- The parameters $a$ and $b_i$ are estimated by MML (method of maximum likelihood), not by least squares
  - (we can't expand on MML in this course)

                                                                    "        "

=

- For this reason, statistical significance and goodness of fit are based not on variance, but on measures of 'deviance' between observed and predicted values
  - = comparison between right and wrong predictions of individual cases

                              (           )                          p
t                    z                              t-test    F-test

- But as in linear regression, estimated parameters (coefficients, intercept) have a $P$-value that determines their significance
  - significance test based on a $z$-distribution related to $t$ and normal distributions
  - interpreted like $t$-tests or $F$-tests. i.e. parameter is significant if P<0.05; 95% confidence intervals are provided etc.

                                                              p<0. 05                      95%

# Logistic regression: categorical variable

Example: let's say we want to test the effect of smoking (X, yes or no) on hypertension (Y, also yes or no)

- Y=0: no hypertension; Y=1: hypertension    Y=0          Y=1          X=0          (          )
- X=0: non-smoker (baseline group); X=1: smoker (exposure group)   X=1          (          )

- Logistic regression model is then:

f = logit p = log(odds of hypertension) = a + b*X

In baseline group, X=0; Therefore

- **log(odds of having hypertension when X=0) = a + b*0 = a**

=the intercept = baseline = reference level (that is, the level of hypertension for non-smokers X = 0)          =          =

a

If we exponentiate a, we obtain odds at baseline

- $e^a$ = odds of hypertension for non-smokers
- $p = \dfrac{1}{1+e^{-a}} = \dfrac{odds\ of\ non-smokers}{1+odds\ of\ non-smokers}$ = probability of hypertension for non-smokers



- Those are the **baseline values**, i.e. the odds and probabilities for groups without exposure (when all $X_i=0$, i.e. even if nobody smoked)

# a + b = log (odds in the exposure group)

- Now the odds for smokers:

  - f = logit = a + bX = a + b.1 = a + b



**a + b = log(odds of hypertension for smokers)**

$e^{a+b} = e^a e^b$ = the odds of hypertension for smokers

$p = \dfrac{1}{1+e^{-(a+b)}} = \dfrac{odds\ of\ smokers}{1+odds\ of\ smokers}$ = probability of hypertension for smokers

Those are the results for the ***exposure group*** (smokers)

# Important: b=log(odds ratio)

So what is b then?

How likely is hypertension if you are a smoker compared a non-smoker?

- answer: it is the odds ratio (of hypertension in smokers to non-smokers)!

So:    b=log(        )

odds ratio = odds(hypert. in smokers)/odds(hypert. in non-smokers)= $e^a e^b / e^a = e^b$

And:

**log(odds ratio) = log($e^b$)=b**

b   log(                                )

- **The coefficient *b* in the logistic regression is the log(odds of hypertension in exposure group *relative to baseline*)**
  - in logistic regression, we test for significance of coefficient *b*
    - same as in linear regression!                    b                                         b        0   p<0.05
    - for a significant effect of variable, we need *b* different from 0 (i.e. P value $< 0.05$)
  - if b=0 (non-significant)       b=0                          /        =1
    - odds ratio for exposure vs. baseline = $e^b = e^0 = 1$       =    X        y
    - = the odds are the same for exposure and baseline (1 to1),
    - = the variable X has no effect on probability of event Y

# Odds ratio b

- Let's add some **_hypothetical_** numbers to the example:

  - odds of hypertension for smokers (=$e^{a+b}$)      = 0.3
  - odds of hypertension for non-smokers (=$e^a$)      = 0.1


- The odds of hypertension in smokers would be three times higher in smokers
  - **_odds ratio_** = odds smokers/odds non smokers = 0.3/0.1 = 3


- The **_odds ratio of the two groups (exposure/baseline)_** is a very useful representation of the effect of a factor on the occurrence of event

  b   log

- Logistic regression always reports b or **log of odds,** not odds of event in exposure group relative to baseline
  - more precisely, R reports _log(odds ratio of event in exposure vs. baseline)_
  - so in this example above, R would report b=log(3)=1.098612
    - We have to exponentiate b to obtain odds ratio = 3

  R           log

# Break

# Example in R: hypertension, smoking, obesity

- File *hypertension* presents data on people with or without hypertension as a function of two factors: smoking and obesity

- Cases coded as 'yes' or 'no' <span style="color:green">no</span>
  - 'no' comes first alphabetically and is read as baseline
  - alternatively: 'no'=0, 'yes'=1 (don't use 1 or 2!!!)   <span style="color:green">no=0  yes=1</span>

  - In this example, data are presented as a table
    - (we'll see a different way of presenting data with each case as a line)

>hypertension

|   | smoking | obesity | total | hyper | nonhyper |
|---|---------|---------|-------|-------|----------|
| 1 | no      | no      | 247   | 40    | 207      |
| 2 | yes     | no      | 102   | 15    | 87       |
| 3 | no      | yes     | 59    | 16    | 43       |
| 4 | yes     | yes     | 25    | 8     | 17       |

# Example in R: hypertension, smoking, obesity

- When data are presented as table
  - table has number of positives (hypertension, $Y = 1$) and negatives (no hypertension, $Y = 0$)
  - two predictors or X variables: $X_1$ = smoking, $X_2$ = obesity
    - For both, yes = 1, no = 0
  - this has been done already for you (file *hypnonhyp*)
    - i.e. the dependent variable will be the matrix *hypnonhyp*

| | hyper | nonhyper |
|---|---|---|
| **1** | 40 | 207 |
| **2** | 15 | 87 |
| **3** | 16 | 43 |
| **4** | 8 | 17 |

- Row 1: non-smoker, non-obese
- Row 2: smoker, non-obese
- Row 3: non-smoker, obese
- Row 4: smoker, obese

Note: don't worry about the table!
  - You will not be asked to create one!

glm

lm

# Running model

> model.hyper <- glm(hypnonhyp ~ smoking+obesity, binomial)

> summary(model.hyper)

Call:

glm(formula = hypnonhyp ~ smoking + obesity, family = binomial)

Deviance Residuals:

| 1 | 2 | 3 | 4 |
|---|---|---|---|
| 0.1593 | -0.2520 | -0.2653 | 0.4018 |

Coefficients:

|  | Estimate | Std. Error | z value | Pr(>\|z\|) |
|---|---|---|---|---|
| (Intercept) | -1.67143 | 0.16731 | -9.990 | < 2e-16 *** |
| smokingyes | -0.01654 | 0.27617 | -0.060 | 0.95224 |
| obesityyes | 0.76005 | 0.28270 | 2.689 | 0.00718 ** |

---

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

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 7.15022  on 3  degrees of freedom

Residual deviance: 0.32067  on 1  degrees of freedom

AIC: 23.935

Number of Fisher Scoring iterations: 3

- Logistic regression is an example of generalised linear model
  - function *glm*

- Syntax is simple and similar to linear regression

- Logistic model written like a multiple regression with *two* predictors:
  - *hypnonhyp ~ smoking+ obesity*
  - (ps. interactions later)
    binomial

- Argument *binomial* sets logistic regression
  - Never forget to add:

  family = binomial

  Otherwise it fits a linear rather than the logistic regression!!!

# Residuals

( )

> model.hyper <- glm(hypnonhyp ~ smoking+obesity, binomial)

> summary(model.hyper)

Call:

glm(formula = hypnonhyp ~ smoking + obesity, family = binomial)

Deviance Residuals:

|   1  |    2   |    3   |   4   |
|------|--------|--------|-------|
| 0.1593 | -0.2520 | -0.2653 | 0.4018 |

Coefficients:

|              | Estimate | Std. Error | z value | Pr(>|z|) |     |
|--------------|----------|------------|---------|----------|-----|
| (Intercept)  | -1.67143 | 0.16731    | -9.990  | < 2e-16  | *** |
| smokingyes   | -0.01654 | 0.27617    | -0.060  | 0.95224  |     |
| obesityyes   | 0.76005  | 0.28270    | 2.689   | 0.00718  | **  |

---

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

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 7.15022  on 3  degrees of freedom

Residual deviance: 0.32067  on 1  degrees of freedom

AIC: 23.935

Number of Fisher Scoring iterations: 3

- Residuals are given as deviance (not variance)

  logit ( yesno)

  - difference between observed and predicted logit values in each group (no/no, no/yes, yes/no, yes/yes)
  - residuals in logit scale (neither probability nor cell count)

# Intercept

> model.hyper <- glm(hypnonhyp ~ smoking+obesity, binomial)

> summary(model.hyper)

Call:

glm(formula = hypnonhyp ~ smoking + obesity, family = binomial)

Deviance Residuals:

   1      2      3      4

 0.1593  -0.2520  -0.2653   0.4018

Coefficients:

                     Estimate Std. Error z value Pr(>|z|)

(Intercept)            -1.67143   0.16731 -9.990  < 2e-16 ***

smokingyes          -0.01654   0.27617 -0.060  0.95224

obesityyes           0.76005   0.28270  2.689  0.00718 **

---

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

(Dispersion parameter for binomial family taken to be 1)

   Null deviance: 7.15022  on 3  degrees of freedom

Residual deviance: 0.32067  on 1  degrees of freedom

AIC: 23.935

Number of Fisher Scoring iterations: 3

$a=\log\left(\quad\right)$

a

z

0

1

0.5

- **Intercept a = -1.67**

- $a=\log$(odds of hypertension, baseline group)
  - =non-smokers (X1=0), non-obese (X2= 0)
  - $e^a$ =the odds of hypertension if you're non-smoker, non-obese
  - a=0.188

- z-test: intercept is significantly different from 0
  - odds of hypertension in baseline  $(e^a)$= not 1
  - probability of hypertension in baseline different from 0.5 in the sample

# Effect of smoking

> model.hyper <- glm(hypnonhyp ~ smoking+obesity, binomial)

> summary(model.hyper)

Call:

glm(formula = hypnonhyp ~ smoking + obesity, family = binomial)

Deviance Residuals:

   1     2     3     4

 0.1593 -0.2520 -0.2653  0.4018

Coefficients:

|  | Estimate | Std. Error | z value | Pr(>\|z\|) |
|---|---|---|---|---|
| (Intercept) | -1.67143 | 0.16731 | -9.990 | < 2e-16 *** |
| smokingyes | -0.01654 | 0.27617 | -0.060 | 0.95224 |
| obesityyes | 0.76005 | 0.28270 | 2.689 | 0.00718 ** |

---

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

(Dispersion parameter for binomial family taken to be 1)

   Null deviance: 7.15022  on 3  degrees of freedom

Residual deviance: 0.32067  on 1  degrees of freedom

AIC: 23.935

Number of Fisher Scoring iterations: 3

$b = \log($

- **Regression coefficient b for smoking:**
  - smokers (X=1) are shown as *smokingyes*,
    - variable name plus group ('yes')
  - b=log(odds ratio)=-0.0165
  - =log odds of hypertension for smokers relative to non-smokers
    $/$ $)$
- But P(z) = 0.95!  $b$  $0$
  - b is not significantly different from 0
  - odds ratio not different from $e^0 = 1$
    $1$
- So smokers are not more likely to have hypertension than non-smokers *in this hypothetical sample*
  - (don't start smoking because of me!)

# Effect of obesity

```
> model.hyper <- glm(hypnonhyp ~ smoking+obesity, binomial)

> summary(model.hyper)

Call:

glm(formula = hypnonhyp ~ smoking + obesity, family = binomial)

Deviance Residuals:
    1       2       3       4
 0.1593  -0.2520  -0.2653  0.4018

Coefficients:
                  Estimate Std. Error z value Pr(>|z|)
(Intercept)       -1.67143    0.16731  -9.990  < 2e-16 ***
smokingyes        -0.01654    0.27617  -0.060  0.95224
obesityyes         0.76005    0.28270   2.689  0.00718 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 7.15022  on 3  degrees of freedom
Residual deviance: 0.32067  on 1  degrees of freedom
AIC: 23.935

Number of Fisher Scoring iterations: 3
```

b=0.76

- **Regression coefficient b for obesity: b=0.76**
  - =log odds of hypertension for obese relative to non-obese

- $P(z) = 0.00718$   b   0
  - b is significantly different from 0
  - b = ln(odds of hypertension in obese relative to baseline) > 0
  - odds ratio= $e^{0.76}$ =2.14      2.14
    - odds ratio >1; obese at higher risk!

    >1

- Obesity more than doubles odds of hypertension *in this sample*

# Goodness of fit

> model.hyper <- glm(hypnonhyp ~ smoking+obesity, binomial)

> summary(model.hyper)

Call:

glm(formula = hypnonhyp ~ smoking + obesity, family = binomial)

Deviance Residuals:

```
    1       2       3       4
 0.1593  -0.2520  -0.2653   0.4018
```

Coefficients:

|  | Estimate | Std. Error | z value | Pr(>\|z\|) |  |
|---|---|---|---|---|---|
| (Intercept) | -1.67143 | 0.16731 | -9.990 | < 2e-16 | *** |
| smokingyes | -0.01654 | 0.27617 | -0.060 | 0.95224 |  |
| obesityyes | 0.76005 | 0.28270 | 2.689 | 0.00718 | ** |

---

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

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 7.15022  on 3  degrees of freedom

Residual deviance: 0.32067  on 1  degrees of freedom

AIC: 23.935

Number of Fisher Scoring iterations: 3

- MML does not use variance to measure goodness of fit
  - it includes no 'dispersion parameter', which has to be taken as 1

- In MML, deviance replaces variance          =                    (=
  - null deviance = deviance when model includes only intercept (=before predictors *smoking* and *obesity*)          )

  - residual deviance is unexplained deviance after predictors

  - difference between null and residual is the contribution of predictors to model

# Goodness of fit

> model.hyper <- glm(hypnonhyp ~ smoking+obesity, binomial)

> summary(model.hyper)

Call:

glm(formula = hypnonhyp ~ smoking + obesity, family = binomial)

Deviance Residuals:

```
    1       2       3       4
 0.1593 -0.2520 -0.2653  0.4018
```

Coefficients:

| | Estimate | Std. Error | z value | Pr(>|z|) | |
|---|---|---|---|---|---|
| (Intercept) | -1.67143 | 0.16731 | -9.990 | < 2e-16 | *** |
| smokingyes | -0.01654 | 0.27617 | -0.060 | 0.95224 | |
| obesityyes | 0.76005 | 0.28270 | 2.689 | 0.00718 | ** |

---

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

(Dispersion parameter for binomial family taken to be 1)

   Null deviance: 7.15022  on 3  degrees of freedom

Residual deviance: 0.32067  on 1  degrees of freedom

AIC: 23.935

Number of Fisher Scoring iterations: 3

<span style="color:green">AIC</span>

<span style="color:green">AIC(akaike ) R</span>

- Because there is no variance, goodness of fit is not measured by $R^2$
  - we use AIC (Akaike Information Criterion) instead


- Remember: adding additional predictors to regression may increase goodness of fit even when predictor is not significant


- AIC measures goodness of fit while punishing models for use of additional predictors
  - *the better and more parsimonious the model, the lower the AIC*

<span style="color:green">AIC</span>

- Models with lowest AIC are selected

<span style="color:green">AIC</span>

# Guide to interpretation and calculations:

- a = log(odds of event in baseline group)
- exp(a) = baseline odds of event
- Probability p of event in baseline: baseline odds/(1 + baseline odds)

    /(1+            )

Then
                        b                0
- b = log(odds ratio); if b is significant (different from 0):
- exp(b) = odds ratio
- exp(a+b) = exp(a)*exp(b) = odds(baseline)*odds ratio = odds in exposure group
- Probability p in exposure group = exposure odds/(1 + exposure odds)

```
model.hyper2 <- glm(hypnonhyp ~ obesity, binomial, data= )
```

## Exercises

- Since *smoking* is not significant, you must optimise the model by excluding *smoking*, and run model only with variable *obesity* (manually, or with *step* function)

1. Is a significant? What does that mean?

2. Is b significant? What does that mean?

- Calculate:

3. Baseline odds of hypertension

4. Odds ratio of hypertension (obese vs. non-obese)

5. Odds of hypertension in obese

6. Probability of hypertension in non-obese

7. Probability of hypertension in obese