



**Zahra Abdulla**

---

Department: Biostatistics and Health  
Informatics

Institute of Psychiatry, Psychology and Neuroscience  
08/2020

**Module Title:** Introduction to Statistics

**Session Title:** Risks and Odds

---

**Topic title:** Binary Logistic Regression



---

After working through this session, you should be able to understand:

- Describe odds and odds ratios
- Understand why and when odds and odds ratios can be used
- How to interpret an odds ratio for a binary outcome
- Why risk ratios cannot always be used as measures of relative risk
- How to compare categorical variables across groups using chi-square, fisher's exact test and odds ratios



# Some Scenarios

---

- Are clients with high scores on a personality test more likely to respond to psychotherapy than are clients with low scores?
- Do children have a better chance of surviving a severe illness than do adults?
- Does gender defined at birth increase the risk of being depressed?

**What kind of methods would give us an answer to these questions?**



# Why look at odds?

---

- We use parametric analysis techniques when we are comparing between groups on a **continuous outcome variable** (e.g. weight differences between gender defined at birth, or cholesterol levels before and after an intervention).
- We use linear regression to look at the change in a **continuous outcome variable** (e.g., mental health dimension)

**But...**

- What if we don't have a continuous outcome variable?
- What about disease outcomes?
  - Develop disease vs. not develop disease?
  - Depressed vs. not depressed?
  - Death vs. survival?
- These variables would typically be measured as binary and so we need to use a different techniques (**odds ratios and logistic regression**) to examine them.



# Odds

---

In health care, the odds describes the ratio of the number of people with the event to the number without.

Odds of 10-1 at the bookmakers ...

- ... means: the probability that the outcome will not happen is 10 times the probability that it will

Or odds of developing a disorder...

- odds of disorder A = the probability that disorder A **does** happen versus the probability that disorder A **does not** happen
- Other examples:
  - an odds of 0.01 is often written as 1:100,
  - odds of 0.33 as 1:3, and
  - odds of 3 as 3:1



# Risk

---

Risk describes the probability with which a health outcome (usually an adverse event) will occur.  
Risk is commonly expressed as a decimal number between 0 and 1

A new drug reduced cancer incidence by 50%

- In absolute terms, the new drug reduced cancer incidence from 2 in 1000 to 1 in 1000.

Relative risk ...

- is the probability of an adverse outcome in an exposure group versus its likelihood in an unexposed group. This statistic indicates whether exposure corresponds to increases, decreases, or no change in the probability of the adverse outcome.
- The exposed group has 0.6 times the risk of the outcome (or 40% less risk of the outcome) compared to the unexposed group.
- More examples
  - when the risk is 0.1, about 10 people out of every 100 will have the event;
  - when the risk is 0.5, about 50 people out of every 100 will have the event.
  - In a sample of 1000 people, these numbers are 100 and 500 respectively.



# How would we calculate the odds and risk?

---

- Could use a **contingency table**

| abuse        | psychosis  | no psychosis | total        |
|--------------|------------|--------------|--------------|
| exposed +    | 127        | 275          | 402          |
| exposed -    | 187        | 1,081        | 1,268        |
| <b>total</b> | <b>314</b> | <b>1,356</b> | <b>1,670</b> |

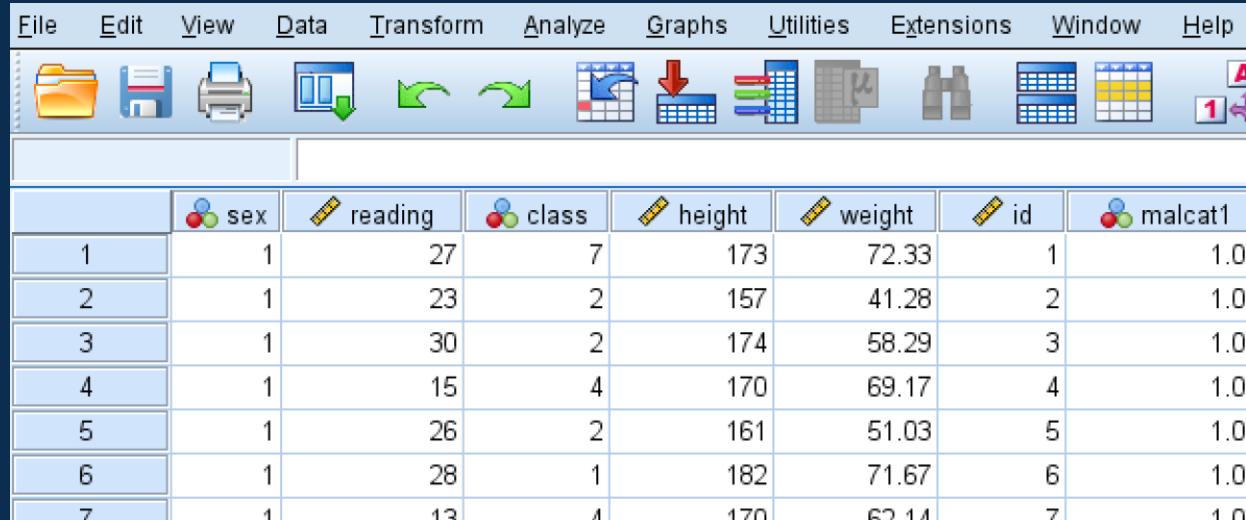
- A contingency table summarizes the frequency distribution of each of two categorical variables as well as the **association between two categorical variables**
- Each cell contains the frequency at which combination of its row and column categories occurred
- A contingency table allows us to check
  - How each of the potential explanatory variables are related to the dependent variable, one by one
  - That categories for explanatory variables are large enough (suggest at least 5 cases per category)
  - How many missing cases there are for each variable



# SPSS Slide

---

Download the data that we are going to use during the lecture. The dataset is the [lecture\\_10\\_data.sav](#).



|   | sex | reading | class | height | weight | id | malcat1 |
|---|-----|---------|-------|--------|--------|----|---------|
| 1 | 1   | 27      | 7     | 173    | 72.33  | 1  | 1.00    |
| 2 | 1   | 23      | 2     | 157    | 41.28  | 2  | 1.00    |
| 3 | 1   | 30      | 2     | 174    | 58.29  | 3  | 1.00    |
| 4 | 1   | 15      | 4     | 170    | 69.17  | 4  | 1.00    |
| 5 | 1   | 26      | 2     | 161    | 51.03  | 5  | 1.00    |
| 6 | 1   | 28      | 1     | 182    | 71.67  | 6  | 1.00    |
| 7 | 1   | 13      | 4     | 170    | 62.14  | 7  | 1.00    |

The dataset contains data from 42 babies, with respect to their  
**Specific body measurements at birth** : headcircumf, length, weight (lbs)

**Gestation:** Gestational age at birth

**Information about the baby's mother:** smoker, motherage, mnocig, mheight, mppwgt

**Information about the baby's father:** fage, fedyrs, fnocig, fheight

**lowbwt:** Low birthweight Baby 0 = No, 1 = Yes

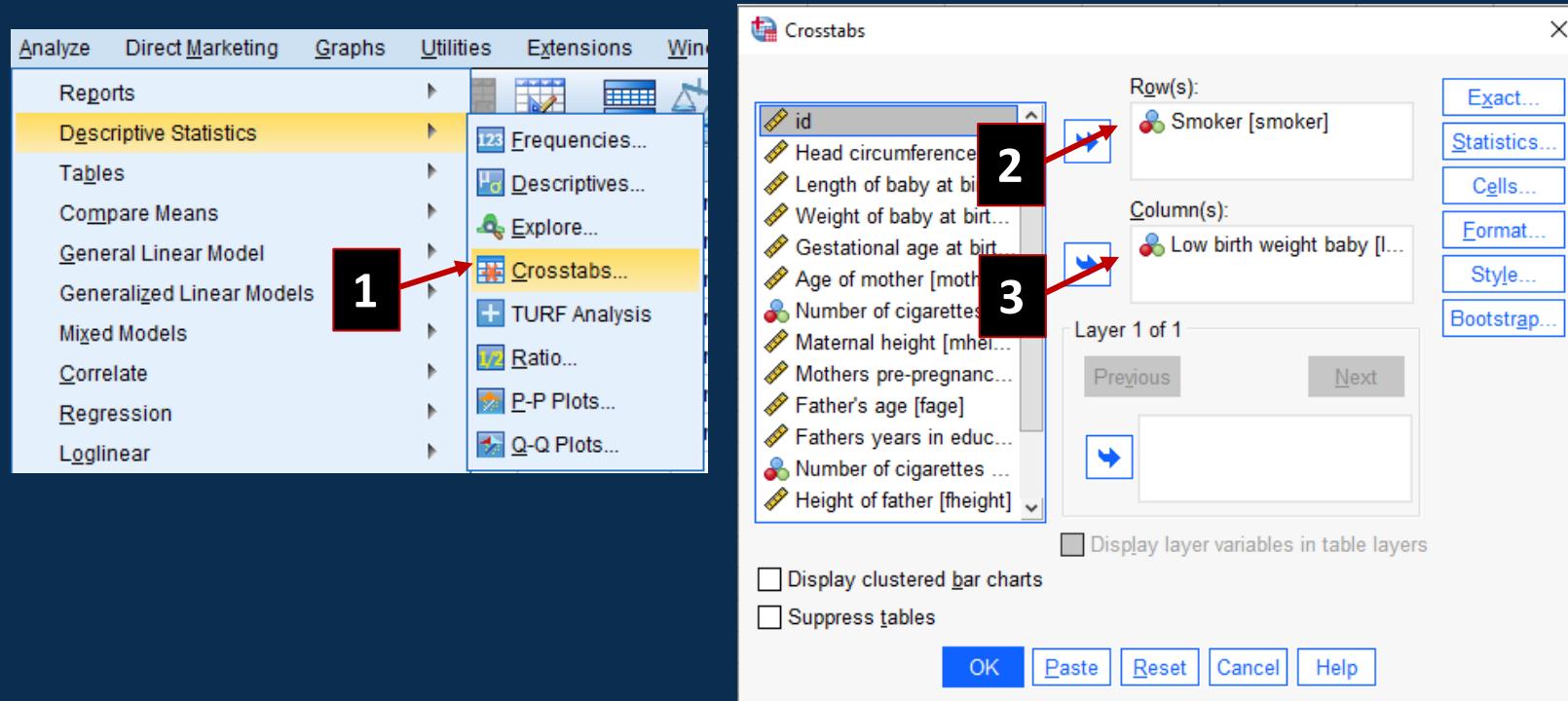
**Mage35:** 0=under 35, 1=Over 35

# SPSS Slide: 'how to'

The next question is: Are the proportions of low weight babies different from mothers who smoked through pregnancy compared to those who did not smoke through pregnancy?

Step 1: Create a contingency table

Analyse -> Descriptive Statistics-> Crosstabs



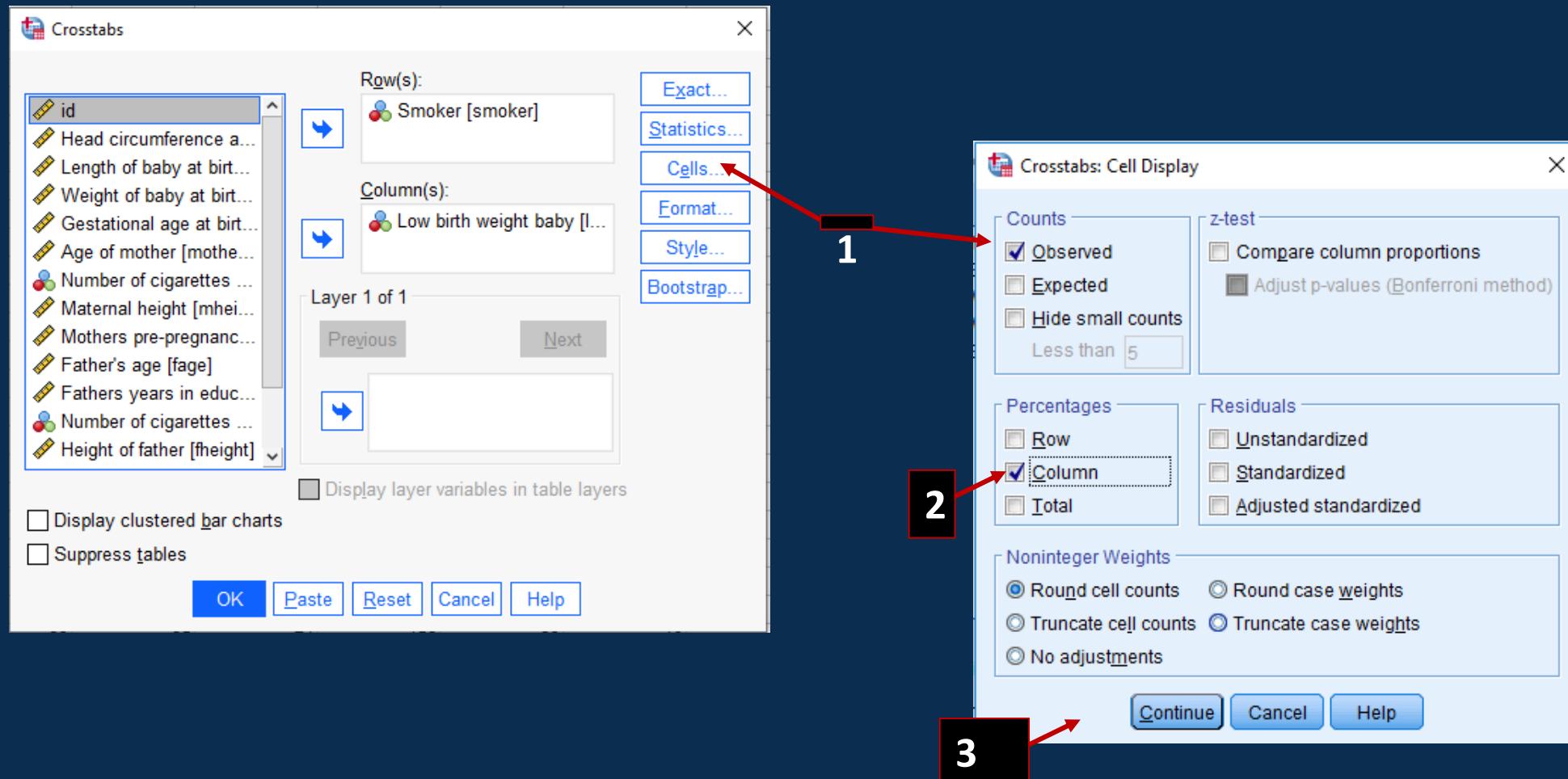
Add the variable of interest (our outcome) (Low Birth Wgt) in to the 'columns box'.

Add the second variable interest (Smoker) in the 'rows box'



# SPSS Slide: 'how to'

## Step 1: Choose the most appropriate 'Percentages'



# Output and Interpretations

Percentages  
 Row  
 Column  
 Total

|        |            | Low birth weight baby |       |       |
|--------|------------|-----------------------|-------|-------|
|        |            | No                    | Yes   | Total |
| Smoker | Non-smoker | Count                 | 15    | 5     |
|        |            | % within Smoker       | 75.0% | 25.0% |
| Total  | Smoker     | Count                 | 9     | 13    |
|        |            | % within Smoker       | 40.9% | 59.1% |
|        |            | Count                 | 24    | 18    |
|        |            | % within Smoker       | 57.1% | 42.9% |
|        |            | Total                 | 20    | 22    |

Among those babies who were low birth weight, the proportion of those whose mothers smoked during pregnancy was higher than the proportion of whose mothers did not smoke during pregnancy (59.1% versus 25.0%, respectively).

Percentages  
 Row  
 Column  
 Total

|        |            | Low birth weight baby          |        |        |
|--------|------------|--------------------------------|--------|--------|
|        |            | No                             | Yes    | Total  |
| Smoker | Non-smoker | Count                          | 15     | 5      |
|        |            | % within Low birth weight baby | 62.5%  | 27.8%  |
| Total  | Smoker     | Count                          | 9      | 13     |
|        |            | % within Low birth weight baby | 37.5%  | 72.2%  |
|        |            | Count                          | 24     | 18     |
|        |            | % within Low birth weight baby | 100.0% | 100.0% |
|        |            | Total                          | 20     | 22     |

Among those mothers who smoked during pregnancy there was a higher proportion who had a baby of low birthweight compared to a baby of normal birthweight (72.2% versus 37.5%, respectively).

NEVER compare percentages which add up to 100%!



# Pearson's chi-square test

---

## When to use

To test if, according to the current data, the proportions in the population of babies being born of low-birth-weight changes based on mothers smoking status during pregnancy

## Hypotheses:

$H_0$ : there is no association between the mother's smoking status and baby's birth weight

$H_a$ : there is an association between the mother's smoking status and baby's birth weight

## Assumptions:

- The observations are randomly and independently drawn
- The number of cells with expected frequencies less than 5, are less than 20%
- The minimum expected frequency is at the very least 1.
- The observations are not paired

# Output and Interpretations

Computations: 'Pearson's chi-square test'.

| Smoker * Low birth weight baby<br>Crosstabulation |            |                       |     |       |
|---|------------|-----------------------|-----|-------|
|   |            | Count                 |     |       |
|   |            | Low birth weight baby |     | Total |
|   |            | No                    | Yes |       |
| Smoker  | Non-smoker | 15                    | 5   | 20    |
|   | Smoker     | 9                     | 13  | 22    |
| Total   |            | 24                    | 18  | 42    |

Row(s):  
Smoker [smoker]

Column(s):  
Low birth weight baby [l...

Counts  
 Observed  
 Expected  
 Hide small counts  
Less than 5

Counts  
 Observed  
 Expected  
 Hide small counts  
Less than 5

| Smoker * Low birth weight baby<br>Crosstabulation |            |                       |      |       |
|---|------------|-----------------------|------|-------|
|   |            | Expected Count        |      |       |
|   |            | Low birth weight baby |      | Total |
|   |            | No                    | Yes  |       |
| Smoker  | Non-smoker | 11.4                  | 8.6  | 20.0  |
|   | Smoker     | 12.6                  | 9.4  | 22.0  |
| Total   |            | 24.0                  | 18.0 | 42.0  |

$$\sum \frac{(O-E)^2}{E} = \frac{(15-11.4)^2}{11.4} + \frac{(5-8.6)^2}{8.6} + \frac{(9-12.6)^2}{12.6} + \frac{(13-9.4)^2}{9.4} = 5.05$$



# Output and Interpretation Slide

| Smoker * Low birth weight baby Crosstabulation |            |                                |        |        |
|--|------------|--------------------------------|--------|--------|
|  |            | Low birth weight baby          | Total  |        |
|  |            | No                             | Yes    |        |
| Smoker   | Non-smoker | Count                          | 15     | 5      |
|  |            | % within Low birth weight baby | 62.5%  | 27.8%  |
|  | Smoker     | Count                          | 9      | 13     |
|  |            | % within Low birth weight baby | 37.5%  | 72.2%  |
| Total  |            | Count                          | 24     | 18     |
|  |            | % within Low birth weight baby | 100.0% | 100.0% |

| Chi-Square Tests                   |                    |    |                                   |                      |                      |
|------------------------------------|--------------------|----|-----------------------------------|----------------------|----------------------|
|                                    | Value              | df | Asymptotic Significance (2-sided) | Exact Sig. (2-sided) | Exact Sig. (1-sided) |
| Pearson Chi-Square                 | 4.972 <sup>a</sup> | 1  | .026                              | .033                 | .027                 |
| Continuity Correction <sup>b</sup> | 3.677              | 1  | .055                              |                      |                      |
| Likelihood Ratio                   | 5.104              | 1  | .024                              | .033                 | .027                 |
| Fisher's Exact Test                |                    |    |                                   | .033                 | .027                 |
| Linear-by-Linear Association       | 4.853 <sup>c</sup> | 1  | .028                              | .033                 | .027                 |
| N of Valid Cases                   | 42                 |    |                                   |                      | .022                 |

<sup>a</sup>. 0 cells (0.0%) have expected count less than 5. The minimum expected count is 8.57.  
<sup>b</sup>. Computed only for a 2x2 table  
<sup>c</sup>. The standardized statistic is 2.203.

Among those mothers who smoked during pregnancy there was a higher proportion who had a baby of low birthweight compared to a baby of normal birthweight (72.2% versus 37.5%, respectively). This difference was statistically significant (Pearson  $\chi^2=4.972$ , df=1, p=0.026).

Therefore, we conclude that mothers who smoked during pregnancy tend to have babies born with low birthweight than women who did not smoke during pregnancy. The variables 'Smoker' and 'Lowbirthwgt' are associated.

# Quantifying the risk

Compare risk of smoking between babies classed as having a low birthweight and those who are not

|        |            | Smoker * Low birth weight baby<br>Crosstabulation |     |       |
|--------|------------|---|-----|-------|
|        |            | Count   |     |       |
|        |            | Low birth weight baby                             |     | Total |
|        |            | No  | Yes |       |
| Smoker | Non-smoker | 15  | 5   | 20    |
|        | Smoker     | 9   | 13  | 22    |
|        | Total      | 24  | 18  | 42    |

The risk of an outcome is the number of times the outcome of interest occurs divided by the total number of possible outcomes.

For example: In the above study out of 22 mothers who smoked during pregnancy, there were 13 babies who were born with low birthweight. So, we get the risk of smoking during pregnancy and having a low birthweight baby by the following calculation:

$$\text{Risk} = 13 \div 22 = 0.59$$



# Calculating the Risk Ratio (Relative Risk)

If we want to compare the effects of smoking during pregnancy and not smoking during pregnancy we could calculate the risk of having a baby of low birth weight for each group:

**Risk** of having baby of low birth weight in smokers =  $13 \div 22 = 0.59$

**Risk** of having baby of low birth weight in non-smokers =  $5 \div 20 = 0.25$

We can compare the risk for each of the groups using the risk ratio.

$$\begin{aligned} &(\text{Risk when smoker}) \div (\text{Risk when non-smoker}) = \\ &0.59 \div 0.25 = 2.36 \end{aligned}$$

So, the risk of having a low birthweight baby when the mother smoked through pregnancy is 2.36 times that of when the mother did not smoke during pregnancy.

|       |            | Smoker * Low birth weight baby Crosstabulation |     |       |
|-------|------------|--|-----|-------|
|       |            | Low birth weight baby                          |     | Total |
| Count | Smoker     | No   | Yes |       |
|       | Non-smoker | 15   | 5   | 20    |
|       | Smoker     | 9  | 13  | 22    |
| Total |            | 24   | 18  | 42    |

# Interpreting the Risk

---

**Relative Risk = 1:** The risk ratio equals one when the numerator and denominator are equal.

- This equivalence occurs when the probability of the event occurring in the exposure group equals the likelihood of it happening in the unexposed group.
- E.g: There is no association between mothers smoking status and a baby being born with a low birth weight.

**Relative Risk > 1:** The numerator is greater than the denominator in the risk ratio.

- Therefore, the event's probability is greater in the exposed group than in the unexposed group.
- E.g. If the RR = 1.4, the smoking status corresponds to a 40% greater probability of a mother having a child with low birthweight.

**Relative Risk < 1:** The numerator is less than the denominator in the risk ratio.

- Consequently, the probability of the event is lower for the exposed group than for the unexposed group.
- E.g. If the RR = 0.4, the smoking status corresponds to a 60% lower probability of a mother having a child with low birthweight.



# Calculating the Odds

The odds in favour of a particular outcome is the number of times the outcome occurs divided by the number of times it doesn't occur.

If we want to compare the effects of smoking and Not smoking during pregnancy we could calculate the odds for each group:

Odds for having a baby of low birthweight when mother is a smoker =  $13 \div 9 = 1.44$

Odds for having a baby of low birthweight when mother is a non-smoker =  $5 \div 15 = 0.33$

We can compare the odds using the odds ratio. The odds ratio for having a baby of low birthweight when mother smokes during pregnancy compare to a mother who did not smoke during pregnancy

$$(\text{Odds when smoker}) \div (\text{Odds when non-smoker}) = 1.44 \div 0.33 = 4.33$$

- So the odds of having a low birthweight baby when the mother smoked during pregnancy is about 4.36 times larger than the odds for mothers who did not smoke during pregnancy.

|        |            | Smoker * Low birth weight baby Crosstabulation |     |       |
|--------|------------|--|-----|-------|
|        |            | Count  |     |       |
|        |            | Low birth weight baby                          |     |       |
|        |            | No   | Yes | Total |
| Smoker | Non-smoker | 15   | 5   | 20    |
| Smoker | Smoker     | 9  | 13  | 22    |
| Total  | Total      | 24   | 18  | 42    |

# Interpreting the odds

---

## OR = 1

Odds of 1 mean the outcome occurs at the same rate in both groups

- Exposure does not affect the odds of outcome
- E.g There is no difference in the odds of low birth rate between smokers and non-smokers.

## OR < 1

Odds of less than 1 mean the outcome occurs less often in the first group than the second group

- Indicates that the exposure is associated with a decreased risk of developing the disease
- E.g if the odds ratio = 0.339 then the odds of a non-smoker having a low birth weight baby is a third (33.9%) of smokers.

## OR > 1

Odds of less than 1 mean the outcome occurs more often in the first group than the second group

- Indicates that the exposure is associated with an increased risk of developing the disease
- E.g. if the odds ratio = 1.5 then the odds of smokers having a low birth weight baby is 1.5 times that of the odds of non-smokers.

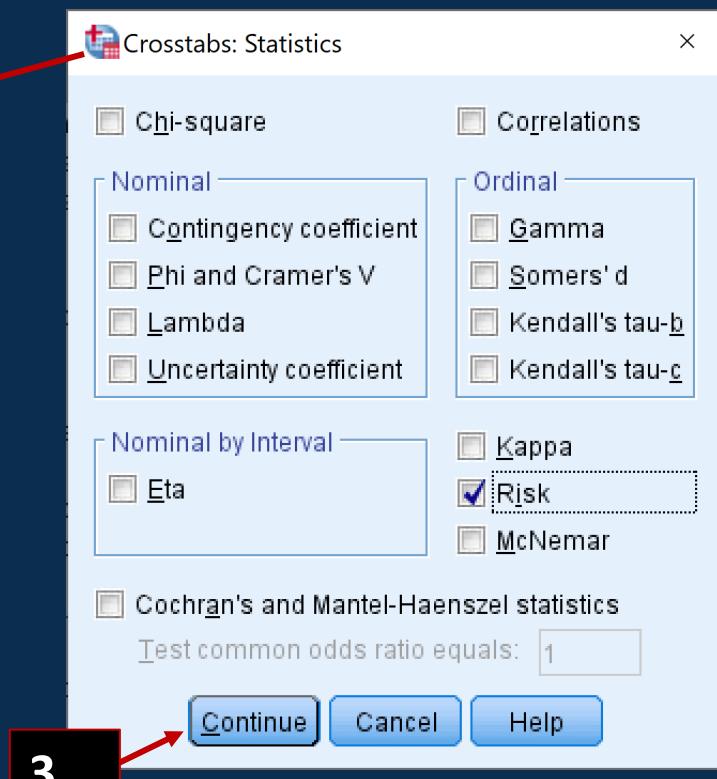
# SPSS Slide: 'how to'

Calculate the risk of having a baby of low birthweight if mothers smoked during pregnancy versus if they did not.

Step 1: Create a contingency table



Step 2: Click on 'Statistics' and Choose "Risk"



3

2

# Output and Interpretation

Odds Ratio

Risk Ratio  
(Relative Risk)

| Risk Estimate                                  |       |                         |        |
|--|-------|-------------------------|--------|
|  | Value | 95% Confidence Interval |        |
|  |       | Lower                   | Upper  |
| Odds Ratio for Smoker<br>(Non-smoker / Smoker) | 4.333 | 1.156                   | 16.248 |
| For cohort Low birth weight<br>baby = No       | 1.833 | 1.045                   | 3.217  |
| For cohort Low birth weight<br>baby = Yes      | .423  | .184                    | .975   |
| N of Valid Cases                               | 42    |                         |        |

The risk ratio given here is 0.423, it is the risk of having a low birthweight baby when the mother did not smoke through pregnancy. To understand the risk of a mother who smoked, we take the reciprocal  $1/0.423 = 2.36$ . So, the risk of a mother who smokes to have a low-birth-weight baby is 2.36 times that of a non-smoker.

# Risk ratios and Odds ratios

---

## Case-control studies

- The risk ratio cannot be used in a case-control study, the odds ratio can be used. Risk ratios cannot be used in studies where selection of subjects is based on the outcome.

## Rare outcomes

- When an outcome is rare the risk ratio and odds ratio will be approximately equal.

Clinical practitioners often prefer the risk ratio due to its more direct interpretation. Statisticians tend to prefer the odds ratio as it applies to a wide range of study designs, allowing comparison between different studies and meta-analysis based on many studies. It also forms the basis of **logistic regression**.

# Knowledge Check

Q1: In the paper Caries prevalence in northern Scotland before and 5 years after, water defluoridation (Stephen et al., 1987, BDJ 163: 324-326) the researchers studied two groups of children in Wick; one group whilst the water was fluoridated and one group after defluoridation.

| Water Type      | Caries |    |  | Total |
|-----------------|--------|----|--|-------|
|                 | Yes    | No |  |       |
| Fluoridated     | 77     | 29 |  | 106   |
| Non-fluoridated | 95     | 31 |  | 126   |
| Total           | 172    | 60 |  | 232   |

Calculate the Risk Ratio and the Odds Ratio for the risk of Caries considering the water is Fluoridated.

# Knowledge Check

Q2: Suppose we conducted a randomised trial to investigate the effect of citalopram on depression. A group of patients who are at risk for depression are randomly assigned to either a placebo or citalopram. At the end of one year, the number of patients suffering with depression is recorded.

|            |  | Depression |      |       |
|------------|--|------------|------|-------|
| Group      |  | Yes +      | No - | Total |
| Placebo    |  | 20         | 80   | 100   |
| Citalopram |  | 15         | 135  | 150   |
| Total      |  | 35         | 215  | 250   |

Calculate the Risk Ratio and the Odds Ratio for the risk of depression for placebo versus citalopram.

# Q1: Knowledge Check Solutions

---

## Risk Ratio

If we want to compare the effects of fluoridated and non-fluoridated water we could calculate the risk of having caries for each group:

Risk of having caries when water is fluoridated =  $77 \div 106 = 0.73$

Risk of having caries when water is not fluoridated =  $95 \div 126 = 0.75$

We can compare the risk for each of the groups using the risk ratio. The risk ratio for being caries free when water is fluoridated compared to when it is not fluoridated is:

$(\text{Risk when fluoridated}) \div (\text{Risk when not fluoridated}) = 0.73 \div 0.75 = 0.96$

So, the risk of having caries when the water is fluoridated is only 0.96 that of when the water is not fluoridated.

# Q1: Knowledge Check Solutions

---

## Odds Ratio

If we want to compare the effects of fluoridated and non-fluoridated water we could calculate the odds for each group:

Odds for having caries when water is fluoridated =  $77 \div 29 = 2.66$

Odds for having caries when water is not fluoridated =  $95 \div 31 = 3.06$

We can compare the odds using the odds ratio. The odds ratio for having caries when water is fluoridated compared to when it is not fluoridated is:

$(\text{Odds when fluoridated}) \div (\text{Odds when not fluoridated}) = 2.66 \div 3.06 = 0.87$

So, the odds of having caries when the water is fluoridated are about 90% those of when the water is not fluoridated.

## Q2: Knowledge Check Solutions

---

### Risk Ratio

If we want to compare the effects of Citalopram and placebo we could calculate the risk of having depression for each group:

Risk of having depression when on Citalopram =  $15 \div 150 = 0.1$

Risk of having depression when on placebo =  $20 \div 100 = 0.2$

We can compare the risk for each of the groups using the risk ratio. The risk ratio for being depression free when on Citalopram compared to placebo is:

$(\text{Risk when on Citalopram}) \div (\text{Risk when on placebo}) = 0.1 \div 0.2 = 0.5$

So the risk of having depression halves when the subject is on Citalopram compared to placebo.

## Q2: Knowledge Check Solutions

---

### Odds Ratio

If we want to compare the effects of Citalopram and placebo we could calculate the odds of having depression for each group:

Odds for having depression when on Citalopram =  $15 \div 135 = 0.11$

Odds for having depression when on placebo =  $20 \div 80 = 0.25$

We can compare the odds using the odds ratio. The odds ratio for having depression when on Citalopram compared to placebois:

$(\text{Odds when citalopram}) \div (\text{Odds when on placebo}) = 0.11 \div 0.25 = 0.44$

So the odds of having depression when on Citalopram are about 45% of those of on placebo. Or another interpretation (1/odds ratio). The odds of having depression when on placebo is 2.3 times the odds of those on citalopram

# References

---

- Altman, D.G., 1991. Practical statistics for medical research, pp. 49-50, 250-253, 259-271. Chapman and Hall, London.
- Bowers, D., 1997. Statistics further from scratch : for health care professionals, pp. 156-157, 162-163. John Wiley & Sons, Chichester.
- Field, Andy. Discovering statistics using IBM SPSS statistics. Sage, 2013. (Chapter 19)
- Agresti, Alan. Categorical data analysis. John Wiley & Sons, 2014.



# Thank you

Contact details/for more information:

Zahra Abdulla

[Zahra.abdulla@kcl.ac.uk](mailto:Zahra.abdulla@kcl.ac.uk)

Department of Biostatistics and Health Informatics (BHI)

IoPPN



**Zahra Abdulla**

---

Department: Biostatistics and Health  
Informatics

Institute of Psychiatry, Psychology and Neuroscience

**Module Title:** Introduction to Statistics

**Session Title:** Binary Logistic Regression

---

**Topic title:** Binary Logistic Regression



---

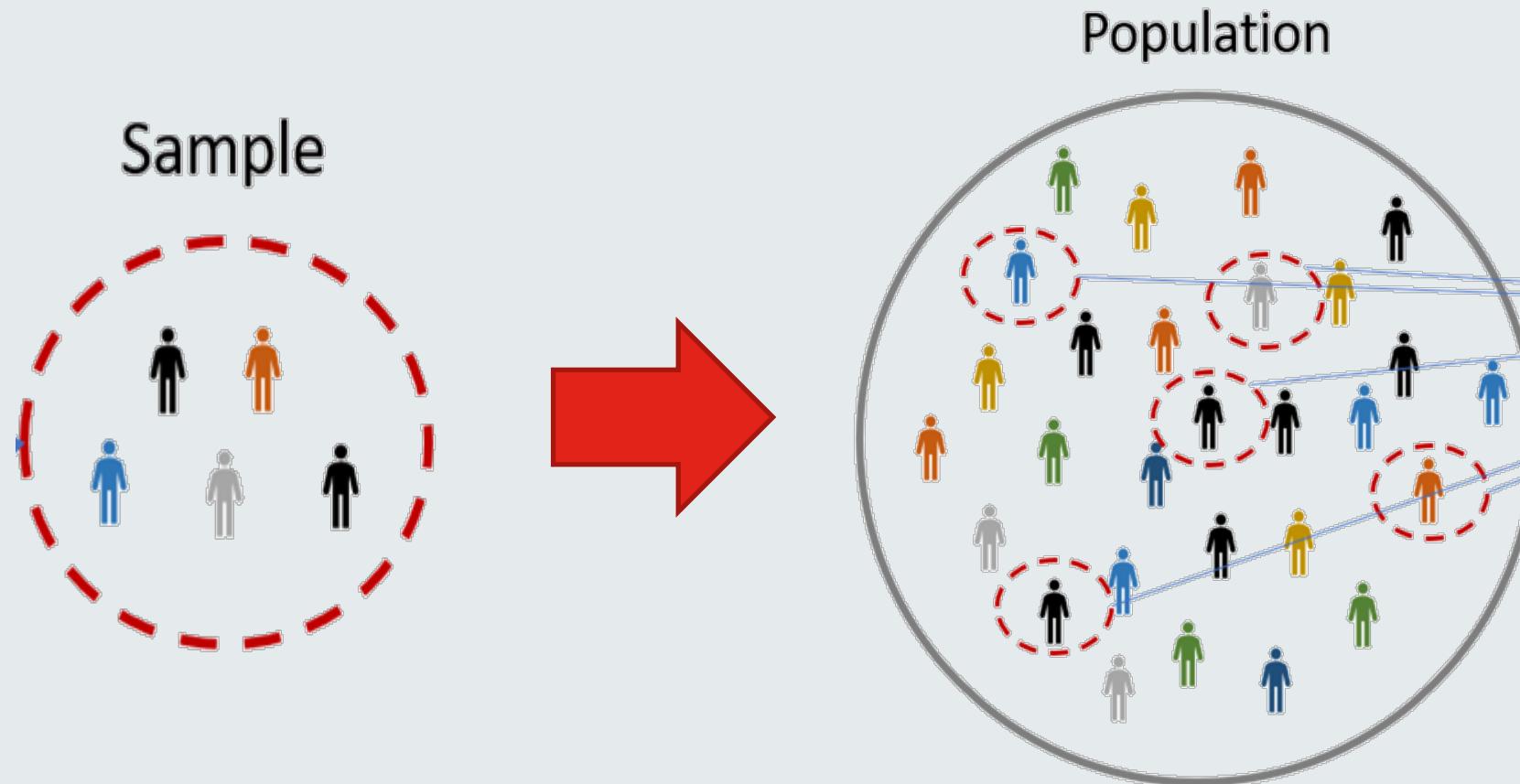
After working through this session, you should be able to:

- Understand what modelling is and why we do it
- Recognise when a binary logistic regression analysis is suitable
- Run a binary logistic regression analysis in a software package



# What is modelling? Recap

---



# Why is statistical modelling important

---

- Example: Investigating the effect of a new app for treating depression
- In a randomised trial we observe reduced depression scores in the group that used the app

**Could the difference have occurred by chance?**

- Modelling allows us to calculate the **association between variables** (e.g., **odds ratios**) as well as **uncertainty about the association** (e.g., **confidence intervals and p values**)



# General linear model (linear regression)

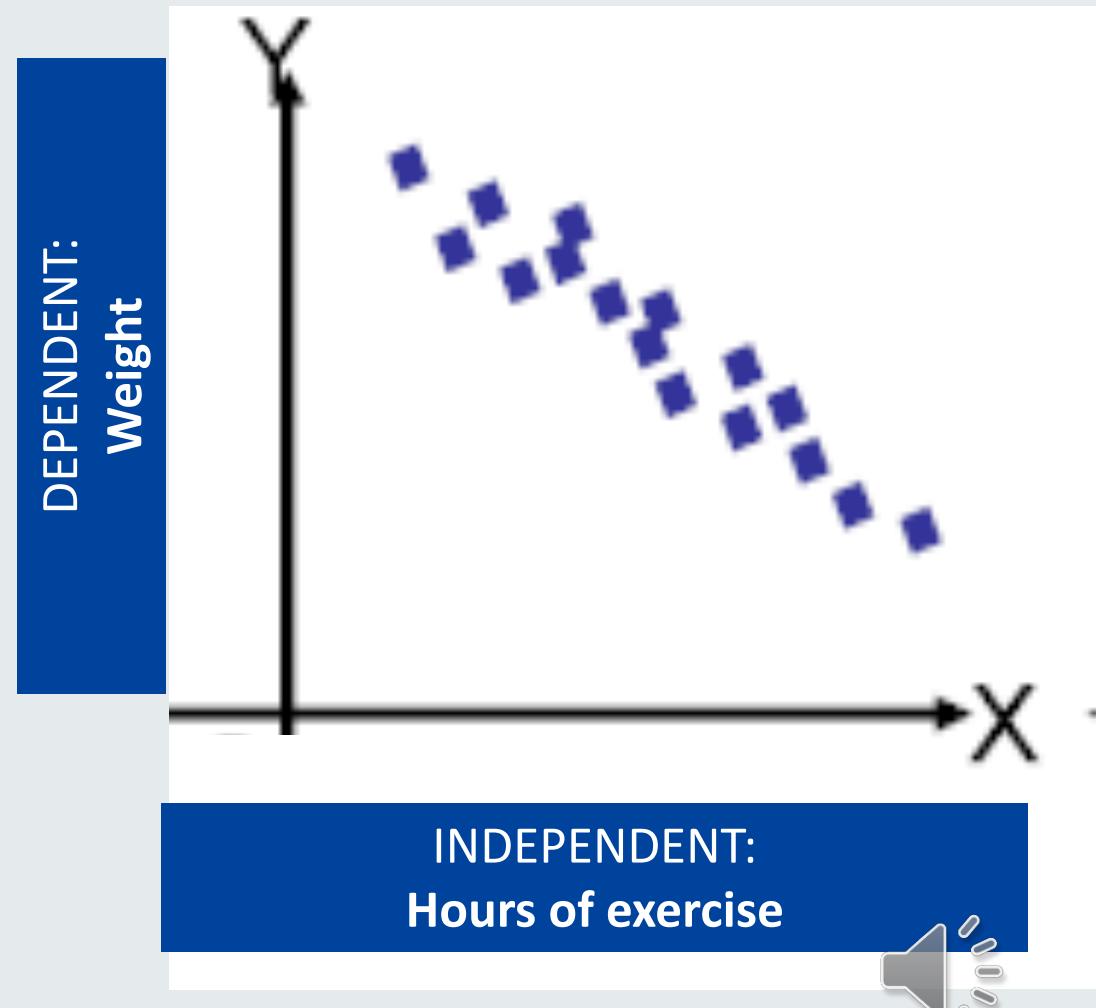
16 people were observed to see if the weight of a person is related to exercise:

**Hypothesis** 'The greater the number of hours of exercise, the lower the weight'.

The plot of data points  $(x,y)$  with  $x = \text{hours of exercise}$  and  $y = \text{weight}$  of a person where the data is continuous is called a **scatterplot**.

Correlation Coefficient (Pearson)  $r=-0.85$

**There is a strong, negative, linear association between hours of exercise and weight loss ( $r=-0.85$ )**



# General linear model (linear regression)

## Interpretation

The relationship is expressed as a linear equation

$$y = \beta_0 + \beta_1 x$$

where  $\beta_0$  is the y intercept = 70

where  $\beta_1$  is the slope of the line = -5

- $\beta_0 = 70$ , When hours of exercise = 0, weight is 70kg.
- $\beta_1 = -5$ , Each additional hour of exercise decreases weight by 5kg.

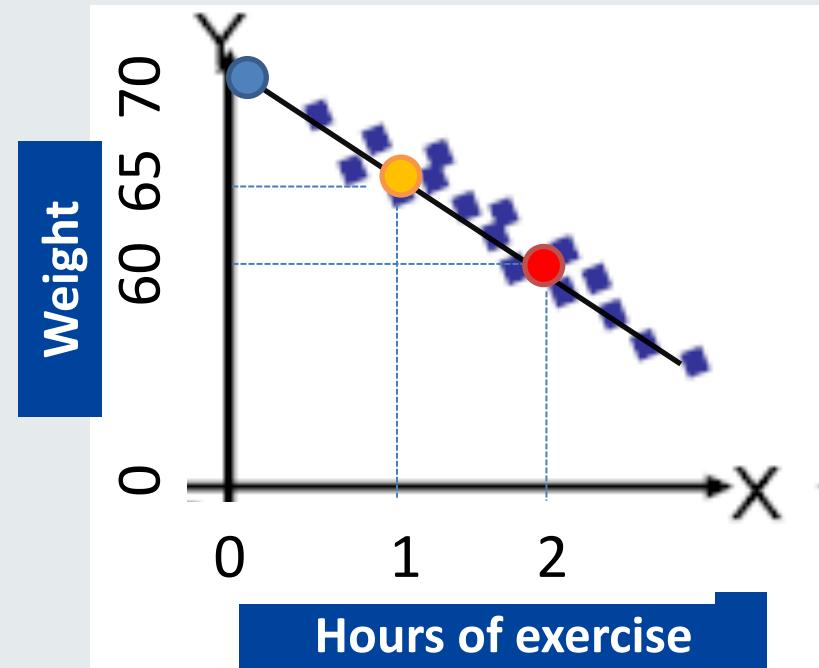
## Linear regression model:

- To measure to what extent there is a linear relationship between two continuous variables, where the outcome variable (dependent variable) is continuous.
- A rule that predicts the dependent variable given the independent variable

$$\beta_0=70; \beta_1=-5;$$

$$y = 70 - 5x$$

| X | Y  |
|---|----|
| 0 | 70 |
| 1 | 65 |
| 2 | 60 |



# Some Scenarios

---

- Are clients with high scores on a personality test more likely to respond to psychotherapy than are clients with low scores?
- Do children have a better chance of surviving a severe illness than do adults?
- Do income, socio economic status and education distinguish persons who are depressed from persons who are not depressed.

**Can we use linear regression to answer these questions?**



# Generalised linear model (logistic regression)

Not all data are suitable for general linear models (linear regression)



What happens when we have other types of data e.g., binary data?

An example: Imagine we wanted to predict whether a person starts smoking or not based on the price of cigarettes at the time they were born

- Here, we have a **binary dependent variable: starts smoking (yes, no)**
- And a **numerical continuous independent variable: price of cigarettes**
  - As the *independent variable is continuous, we can't use cross-tabs.*

We want to know the **probability** that any given person will start smoking or not, at each price



And hence the **proportion** of people that will start smoking at each price on average

# Examples of binary Outcomes

---

Outcomes in Psychology and Psychiatry are often binary:

- Illness (Schizophrenia, Autism,...)
- Passing some threshold (Depression, Anxiety, Obese, ....)
- Recurrence of psychosis
- Hospitalization
- Survival
- Hospital discharge
- Relapse to alcohol use

Often you need to define a timeframe:

- Depressive symptoms within the last year

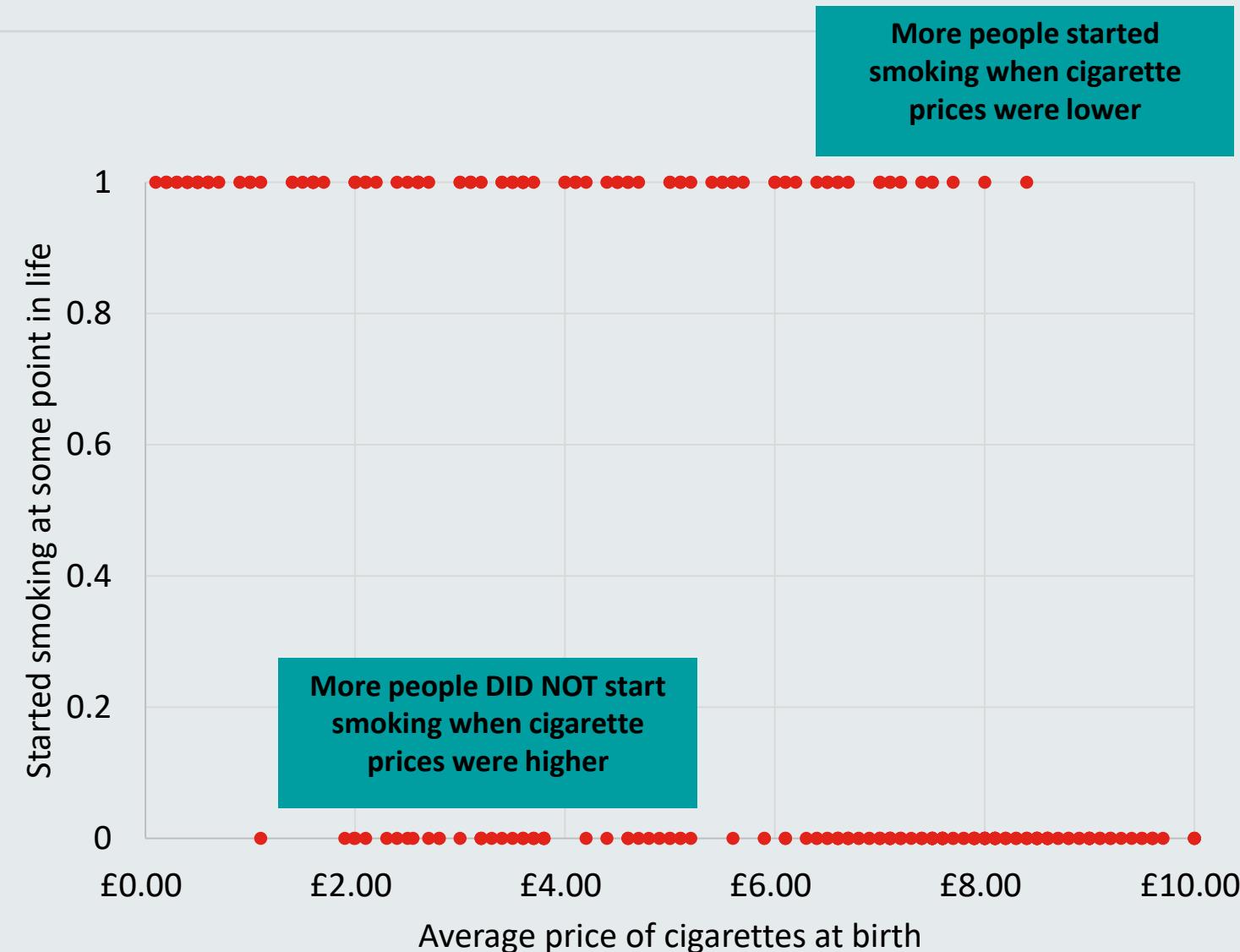
**Do not dichotomize if not necessary (Loss of information)**

# What is wrong with Simple Linear Regression?

We want to predict a probability; this can only vary between zero and 1

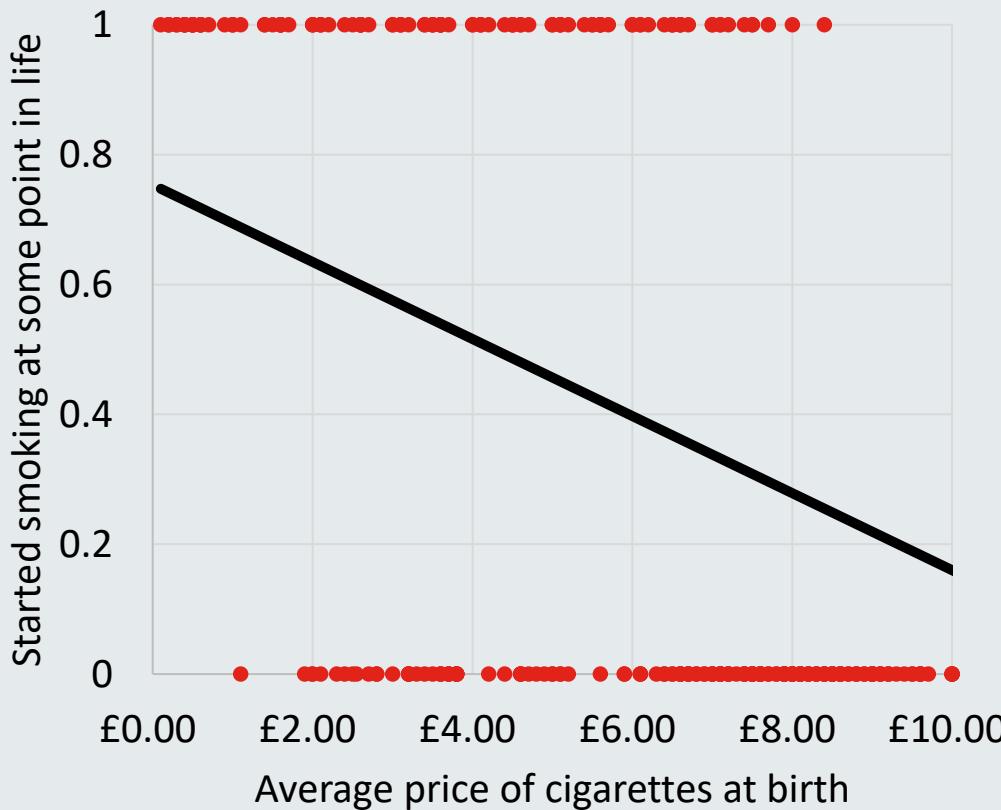
But our simple linear regression may predict values that are below zero or above 1

This is a scatter plot of 400 people who answered a survey about their smoking behaviours, plotted against the average price of cigarettes at the time they were born



# More problems with Simple Linear Regression

Could add a linear regression line, but prediction would not make much sense (not below 0 or 1)



For linear regression we assumed that the population distribution was normally distributed around the mean, for each value of the X variable.

That's not going to be the case if we've got a binary response. The distribution around the mean is going to be quite different.



# Non-linear relationship

---

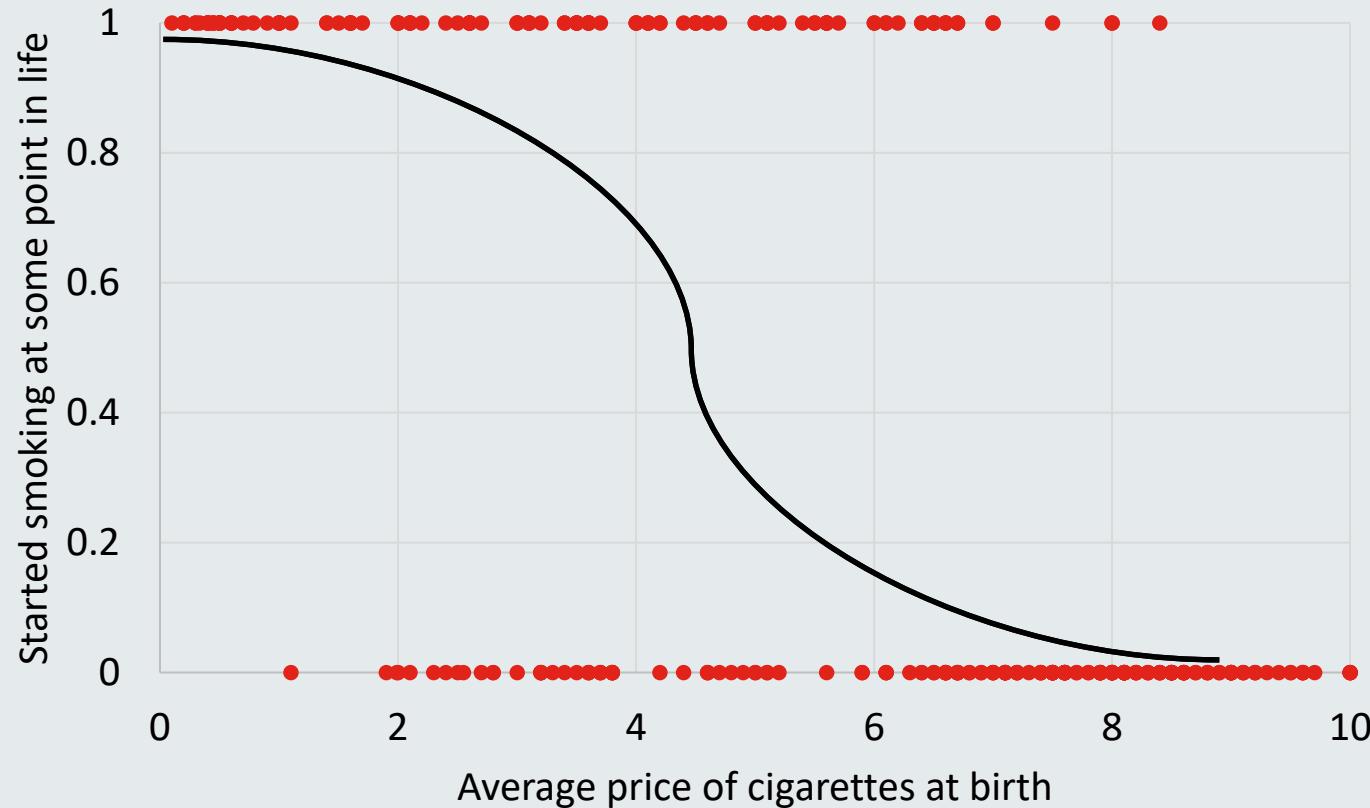
Linear relationship does not make sense for binary outcomes

We rather assume a **nonlinear** relationship as

- Output variable is limited to **[0,1]**, some of our observations are outside this range
- Our goal is to separate best the two groups not to **minimize Mean Squared Error**.
- Linear regression would be **highly sensitive to influential cases**
- Assumptions of linear regression are violated (esp. homogeneity of variances) and hence inference is not valid

# Non-linear relationship

- We assume there is a non-linear S-shaped (or sigmoid) relationship between cigarette price and starting smoking.
- Here is a more realistic representation of the relationship between the probability of cigarette price and starting smoking:



- Lower the price = more likely to start smoking
- There is never a probability of 0 or 100%



# The Link Function

---

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

Linear regression describes the **linear** relationship between the outcome  $y$  and the predictor variable(s)  $x_i$  in a general linear model, where  $\epsilon$  describes the random component (error) which is assumed to be normal distributed.

**Generalised Linear Models (GLM)** extend the ordinary regression model and allow the response variable (dependent, outcome)  $y$  to have an error distribution other than the normal distribution.

In a **logistic** regression, we relate  $x_i$  and the mean outcome at  $x_i$  ( $\mu$ ) by way of a **function**, known as **link function  $g(\mu)$** :

$$g(\mu) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

a link function will connect a model's outcome to its predictors in a linear way, so that we can model a linear relationship between the left- and right-hand side of the equation.



# Logistic Regression: The Link Function

The link function in **logistic regression** is called the **Logit** link (used when data are binary):

$$g(\mu) = \ln\left(\frac{\pi}{1-\pi}\right)$$

When we have a binary outcome, the errors will follow a binomial distribution, where the mean of outcome  $y$  is represented by the probability (proportion)  $\pi$  of an event bounded by 0 and 1, as a function of the predictor variables. The logit link function will transform the data into a logit scale so that we can model a linear relationship between the left and right hand side of the equation.

$$\ln\left(\frac{\pi}{1-\pi}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_i x_i$$

Natural log.

This is just the **odds**,  
the probability that expected outcome **does**  
happen divided by the probability that expected  
outcome **does not** happen

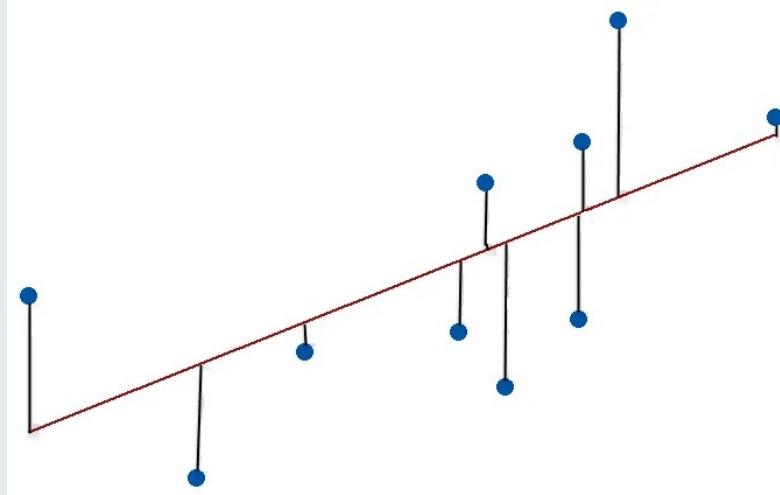
The (adjusted) odds ratio is the estimated change in odds  
for a unit change in  $x_1$  (holding  $x_2 x_3, \dots x_i$  constant)

For variables coded as binary or dummy variables 'one  
unit' usually means a comparison between the group  
of interest and a reference group.

# Fitting this model (1)

---

- With SLR we tried to **minimize the squares of the residuals**, to get the best fitting line.



- This doesn't really make sense here (remember the errors won't be normally distributed as there's only two values).
- We use something called **maximum likelihood** to estimate the coefficients of the linear predictors

## Fitting this model (2)

---

- **Maximum likelihood** is an **iterative process** that estimates the best fitted equation.
- The coefficients maximise the probability (likelihood) of obtaining the actual group membership for cases in the sample (e.g. depressed)
- Coefficients are known as **Maximum Likelihood parameters**



## An example...

| Variable        | Coefficient value | Standard error | p-value |
|-----------------|-------------------|----------------|---------|
| Cigarette price | -0.07             | 0.01           | 0.00    |
| Intercept       | 3.69              | 0.72           | 0.00    |

- In OLS linear regression, a change of one unit on the X variable meant that the Y variable would increase by the coefficient for X.
- That's not what the coefficient associated with X in our logistic regression means.
  - It's clear that cigarette price has a negative (and statistically significant) effect on starting smoking – i.e., as cigarette price increases the probability of starting smoking decreases.
  - But what does the -0.07 actually mean?

In logistic regression an increase in X of 1 unit will decrease our log (odds) by 0.07.

The anti-log ( $e^x$ ) of -0.07 gives us the odds ratio for price



# Logistic Regression: The Logistic Model

---

$\pi$

This is the **Probability** of an event

$\frac{\pi}{1-\pi}$

This is the **Odds** of an event

**Model:**  $\ln\left(\frac{\pi}{1-\pi}\right) = \alpha + \beta x$

This is called the **Logit**

$L = \alpha + \beta x$

This is called the **Linear Predictor**

The model is **linear** in the **logit** but **non-linear** in the probability  $\pi$ . The data are fitted using **logits** and then one **transforms** the **fitted** parameters to **probabilities** afterwards.

$\exp(L) = e^L$

This is the **Odds** of an event

$\hat{\pi} = \frac{\text{odds}}{1+\text{odds}}$

This is the **Estimated Probability** of an event

$\hat{\pi} = \frac{\exp(L)}{1+\exp(L)} = \frac{1}{1+\exp(-L)}$



# Binary Logistic Regression

---

## When to use

To test, according to the current data, if in the population there is an association between babies being born of low birth weight and mothers' smoking status during pregnancy

## Hypotheses:

- $$H_0: \text{there is no association between the mother's smoking status and baby's birth weight}$$
- $$H_a: \text{there is an association between the mother's smoking status and baby's birth weight}$$

## Assumptions:

- Binary dependent variable which has a **Bernoulli (binomial)** Distribution
- Continuous variables have a linear effect on the log-odds scale (Is linearly related to the predictor variables only after transforming into the **logit** scale )
- Observations are independent
- Adequate Sample size
- Absence of multicollinearity
- No outliers

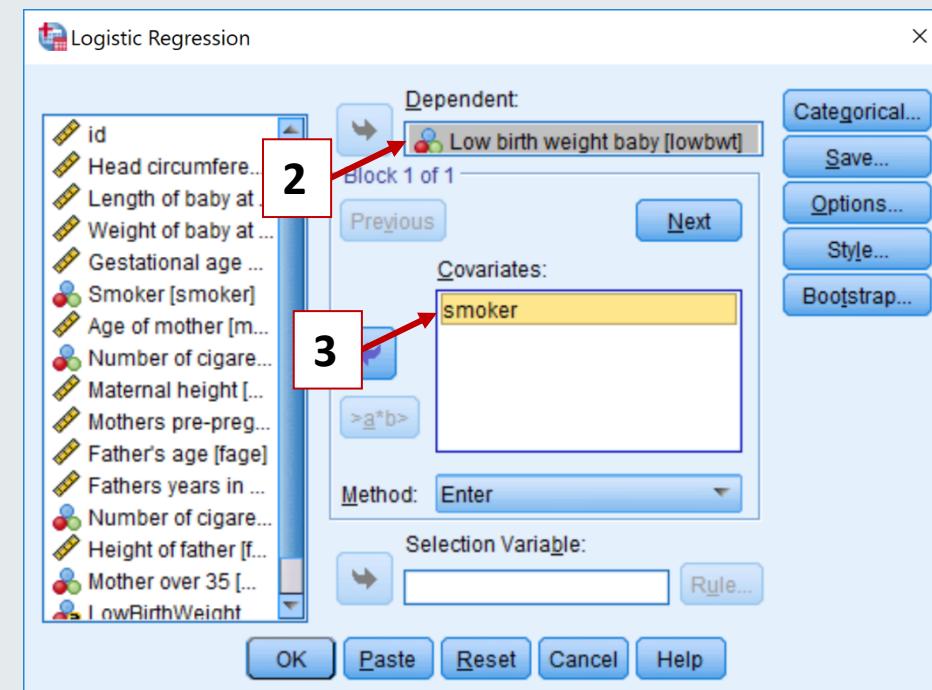
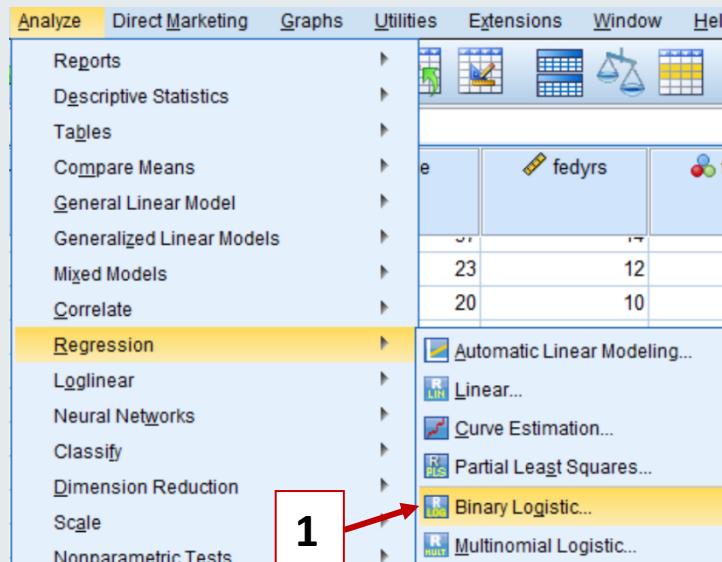


# SPSS slide: 'how to'

Is there an association between having a baby of low birth weight with mothers who smoked through pregnancy? Use the Lecture\_10\_data.sav

**Step 1:** Use the appropriate test, here: 'Binary Logistic Regression'.

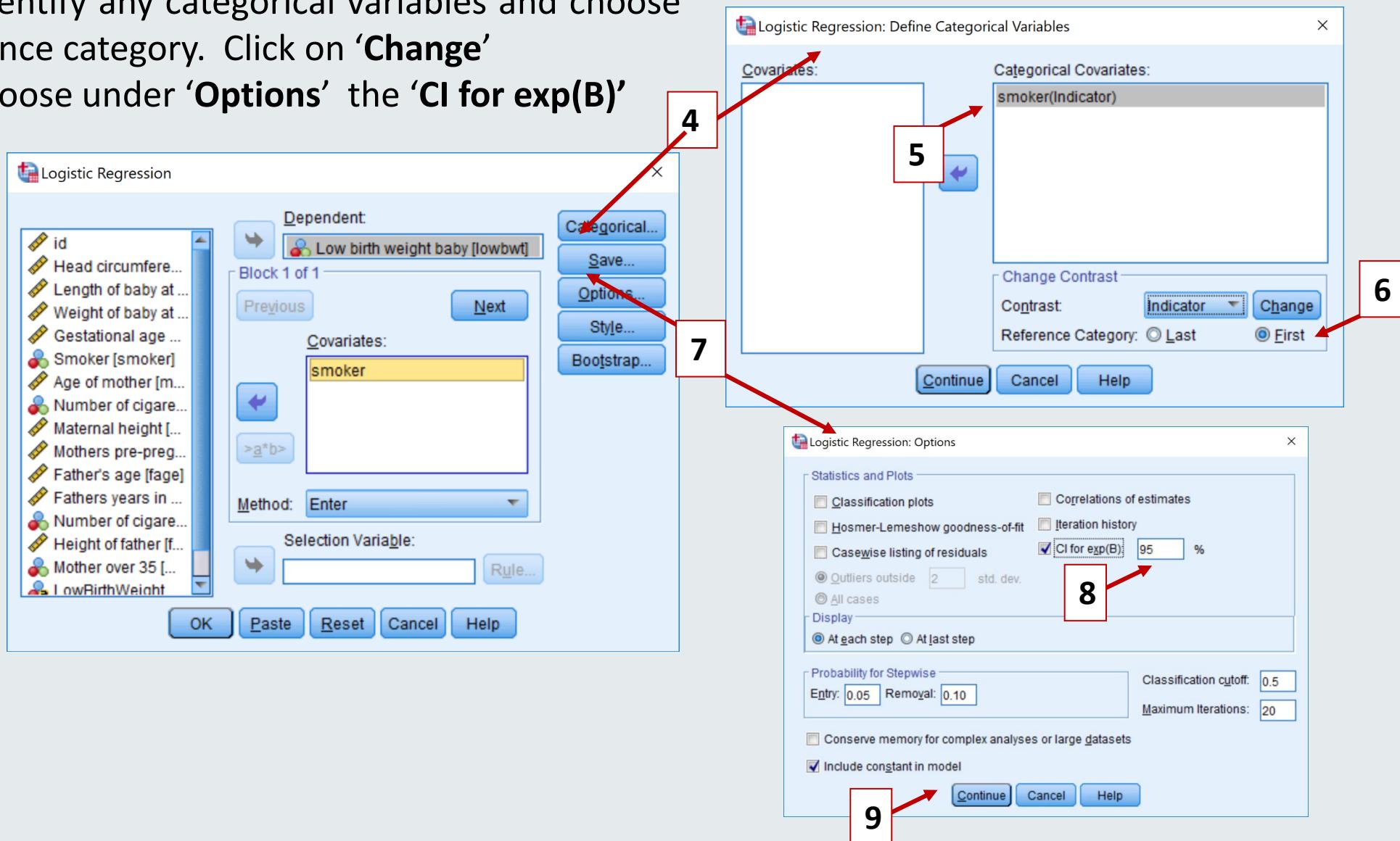
Analyse -> Regression> Binary Logistic



# SPSS slide: 'how to'

**Step 2:** Identify any categorical variables and choose the reference category. Click on '**Change**'

**Step 3:** choose under '**Options**' the '**CI for exp(B)**'



# Output and Interpretation

| Variables in the Equation |           |        |      |       |      |        |                     |              |
|---------------------------|-----------|--------|------|-------|------|--------|---------------------|--------------|
|                           | B         | S.E.   | Wald | df    | Sig. | Exp(B) | 95% C.I. for EXP(B) |              |
|                           |           |        |      |       |      |        | Lower               | Upper        |
| Step 1 <sup>a</sup>       | Smoker(1) | 1.466  | .674 | 4.729 | 1    | .030   | 4.333               | 1.156 16.248 |
|                           | Constant  | -1.099 | .516 | 4.526 | 1    | .033   | .333                |              |

a. Variable(s) entered on step 1: Smoker.

## Regression Equation

$$\ln \frac{p}{1-p} = -1.099 + 1.466 \text{smoker}$$

Odds ratio for the effect of mothers who smoked during pregnancy on low-birth-weight  $\text{Exp}(\beta) = 4.333$ .

There is significant evidence ( $p=.030$ ) of an association between mothers smoking status during pregnancy and a baby being born at a low birth weight. Mothers who smoke during pregnancy have a **4.33 times larger** odds of having a baby born with low birth weight compared to a mother who did not smoke during pregnancy **95%CI 1.156 to 16.248, p=0.030**.



## Thinking back to why this is important... - rerecord

---

- *Odds ratio for the effect of mothers who smoked during pregnancy on low birth weight  $\text{Exp}(\beta) = 4.333$ . There is significant evidence ( $p=.030$ ) of an association between mothers smoking status during pregnancy and a baby being born at a low birth weight. Mothers who smoke during pregnancy have a **4.33 times larger** odds of having a baby born with low birth weight compared to a mother who did not smoke during pregnancy **95%CI 1.156 to 16.248, p=0.030**.*
- We now have a measure of association (odds ratio):  $\text{Exp}(\beta) = 4.333$
- As well as measures of uncertainty (confidence intervals and p values):  
 $95\% \text{CI} = 1.156 \text{ to } 16.248, p=0.030$

# Knowledge Check

---

Q1:

Consider research by Wuensch & Poteat, published in the *Journal of Social Behavior and Personality* in 1998

College students ( $N = 315$ ) were asked to pretend that they were serving on a university research committee regarding a complaint against animal research being conducted by a faculty member. The complaint included a description of the research in simple but emotional language.

In his defence, the researcher made a case about the benefits of his research and steps taken to ensure the animals did not experience pain.

After reading the case materials, each participant was asked to decide whether or not to withdraw the faculty members' authorization to conduct the research.

# Knowledge Check

Q1 Cont. Let us first consider a simple (bivariate) logistic regression, using subjects' decisions as the dichotomous criterion variable and their gender as a dichotomous predictor variable. gender coded with 0 = Female, 1 = Male, and decision with 0 = "Stop the Research" and 1 = "Continue the Research".

Write the regression equation and interpret the output below

Model Summary

| Step | -2 Log likelihood    | Cox & Snell R Square | Nagelkerke R Square |
|------|----------------------|----------------------|---------------------|
| 1    | 399.913 <sup>a</sup> | .078                 | .106                |

a. Estimation terminated at iteration number 3 because parameter estimates changed by less than .001.

Variables in the Equation

|                        | B        | S.E.  | Wald | df     | Sig. | Exp(B) |
|------------------------|----------|-------|------|--------|------|--------|
| Step <sup>a</sup><br>1 | gender   | 1.217 | .245 | 24.757 | 1    | .000   |
|                        | Constant | -.847 | .154 | 30.152 | 1    | .000   |

a. Variable(s) entered on step 1: gender.

# Knowledge Check Solution

- We can interpret **Nagelkerke R<sup>2</sup>** 10.6% of the variation in ‘Stop the research’ can be explained by the model including Gender.
- The exp( $\beta$ ) tells us that the model predicts that the odds of deciding to continue the research are 3.376 times higher for men than they are for women.

| Model Summary |                      |                      |                     |
|---------------|----------------------|----------------------|---------------------|
| Step          | -2 Log likelihood    | Cox & Snell R Square | Nagelkerke R Square |
| 1             | 399.913 <sup>a</sup> | .078                 | .106                |

a. Estimation terminated at iteration number 3 because parameter estimates changed by less than .001.

| Variables in the Equation |          |       |      |        |      |        |
|---------------------------|----------|-------|------|--------|------|--------|
|                           | B        | S.E.  | Wald | df     | Sig. | Exp(B) |
| Step <sup>a</sup><br>1    | gender   | 1.217 | .245 | 24.757 | 1    | .000   |
|                           | Constant | -.847 | .154 | 30.152 | 1    | .000   |

a. Variable(s) entered on step 1: gender.

# References

---

**J Scott Long, Sage, 1997** Regression Models for Categorical and Limited Dependent Variables

**A Agresti, Wiley, 2002** An Introduction to Categorical Data Analysis 2nd ed by

Wuensch and Poteat example: <http://core.ecu.edu/psyc/wuenschk/MV/Multreg/Logistic-SPSS.PDF>



# Thank you

Contact details/for more information:

Zahra Abdulla

[Zahra.abdulla@kcl.ac.uk](mailto:Zahra.abdulla@kcl.ac.uk)

Department of Biostatistics and Health Informatics (BHI)

IoPPN



**Zahra Abdulla**

---

Department: Biostatistics and Health  
Informatics

Institute of Psychiatry, Psychology and Neuroscience  
08/2020

**Module Title:** Introduction to Statistics

**Session Title:** Multiple Independent Variables

---

**Topic title:** Binary Logistic Regression



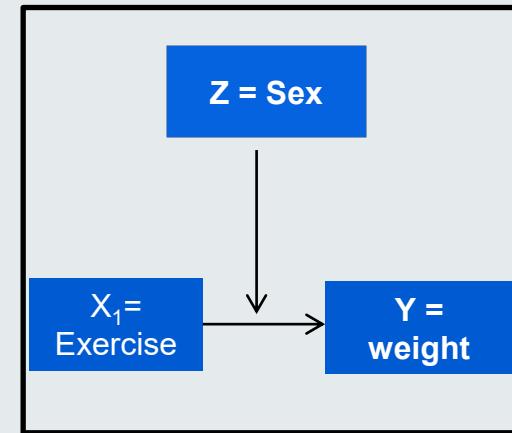
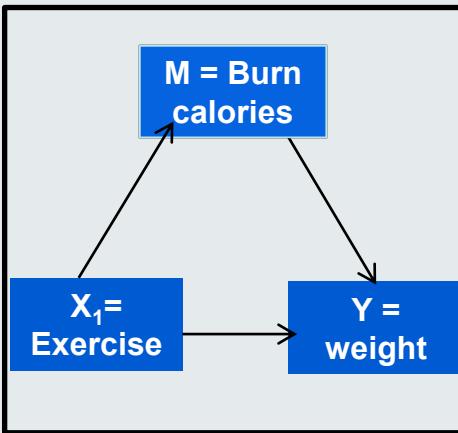
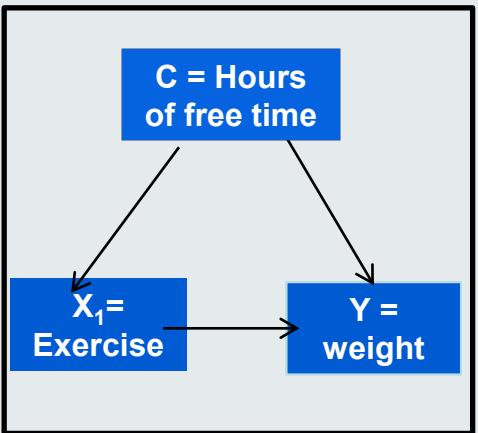
---

After working through this session you should be able to:

- Interpret a binary logistic regression model with multiple independent variables.
- Run a binary logistic regression analysis with multiple predictors in a software package.

# Dealing with third variables

Both confounder, mediator and moderator, are third variables that explain a part (or most) of the association between an independent and dependent variable.



A **confounder (C)** has a common effect on the independent and dependent variables. A confounder is **extrinsic to the causal pathway**.

A **mediator (M)** is caused by the independent variable which in turn causes the dependent variable. A mediator is **in the causal pathway**.

A **moderator (Z)** modifies the effect of an independent variable on a dependent variable. The association varies depending on the values/levels of Z



# The logistic transformation: Multiple predictors

Just as we would be able to develop a Multiple Linear Regression model we are able to build a Binary logistic regression with multiple independent variables. This includes investigating

- Confounding Variables
- Mediators
- Effect Modifiers or Interaction Terms.

Independent or Predictor variables can be numerical or categorical

$$\ln \left( \frac{\pi}{1 - \pi} \right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_i x_i$$

This is just the *odds*.

The (adjusted) odds ratio is the estimated change in odds for a unit change in  $x_1$  (holding  $x_2, x_3, \dots, x_i$  constant)

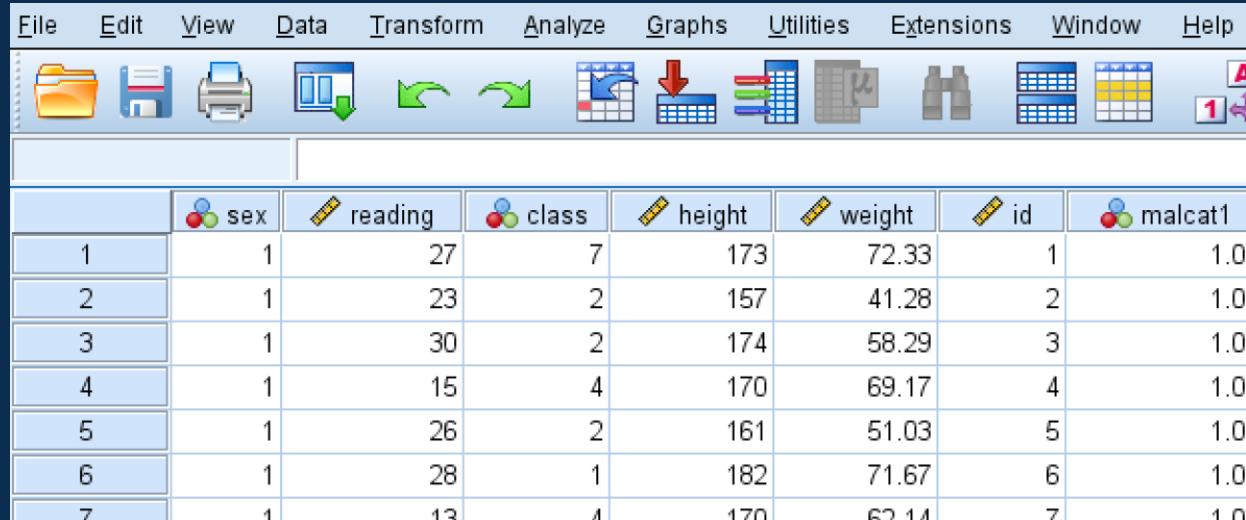
For variables coded as binary or dummy variables 'one unit' usually means a comparison between the group of interest and a reference group.



# SPSS Slide

---

Download the data that we are going to use during the lecture. The dataset is the [lecture\\_10\\_data.sav](#).



|   | sex | reading | class | height | weight | id | malcat1 |
|---|-----|---------|-------|--------|--------|----|---------|
| 1 | 1   | 27      | 7     | 173    | 72.33  | 1  | 1.00    |
| 2 | 1   | 23      | 2     | 157    | 41.28  | 2  | 1.00    |
| 3 | 1   | 30      | 2     | 174    | 58.29  | 3  | 1.00    |
| 4 | 1   | 15      | 4     | 170    | 69.17  | 4  | 1.00    |
| 5 | 1   | 26      | 2     | 161    | 51.03  | 5  | 1.00    |
| 6 | 1   | 28      | 1     | 182    | 71.67  | 6  | 1.00    |
| 7 | 1   | 13      | 4     | 170    | 62.14  | 7  | 1.00    |

The dataset contains data from 42 babies, with respect to their  
**Specific body measurements at birth** : headcircumf, length, weight (lbs)

**Gestation:** Gestational age at birth

**Information about the baby's mother:** smoker, motherage, mnocig, mheight, mppwgt

**Information about the baby's father:** fage, fedyrs, fnocig, fheight

**lowbwt:** Low birthweight Baby 0 = No, 1 = Yes

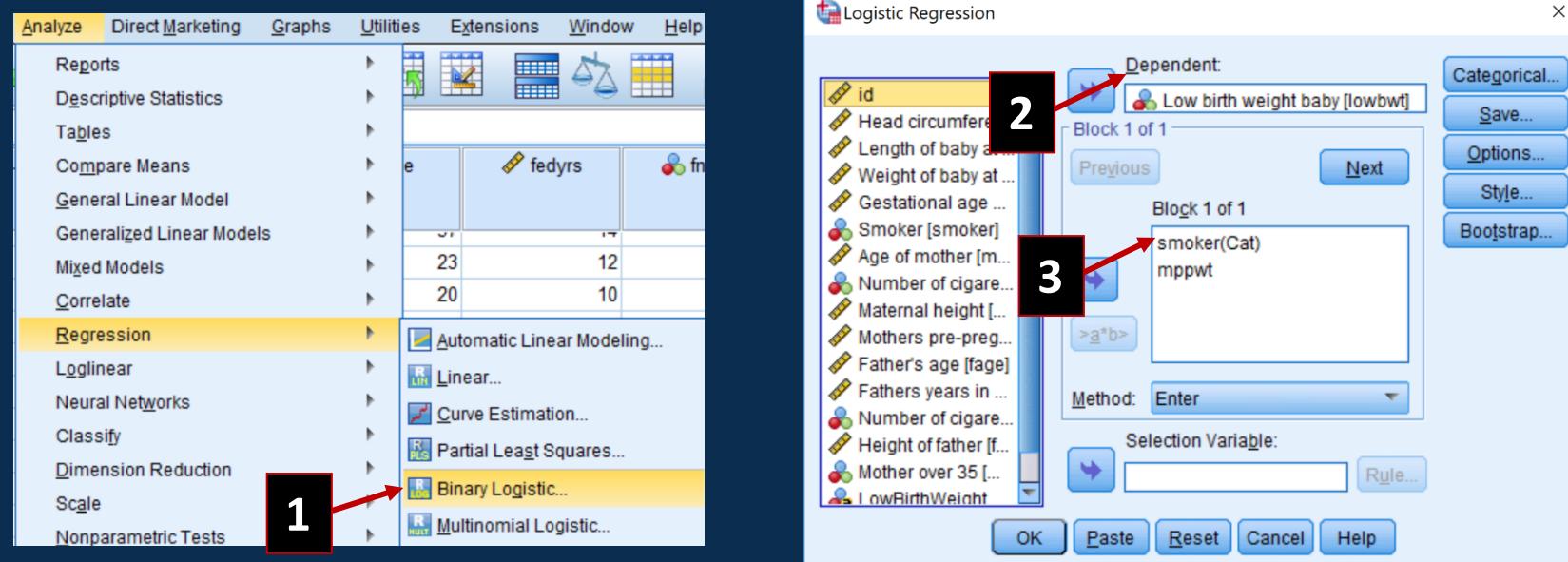
**Mage35:** 0=under 35, 1=Over 35

# SPSS slide: ‘how to’

Is there an association between having a baby of low birth weight with mothers who smoked through pregnancy adjusting for mother's weight pre-pregnancy?

**Step 1:** Use the appropriate test, here: ‘Binary Logistic Regression’.

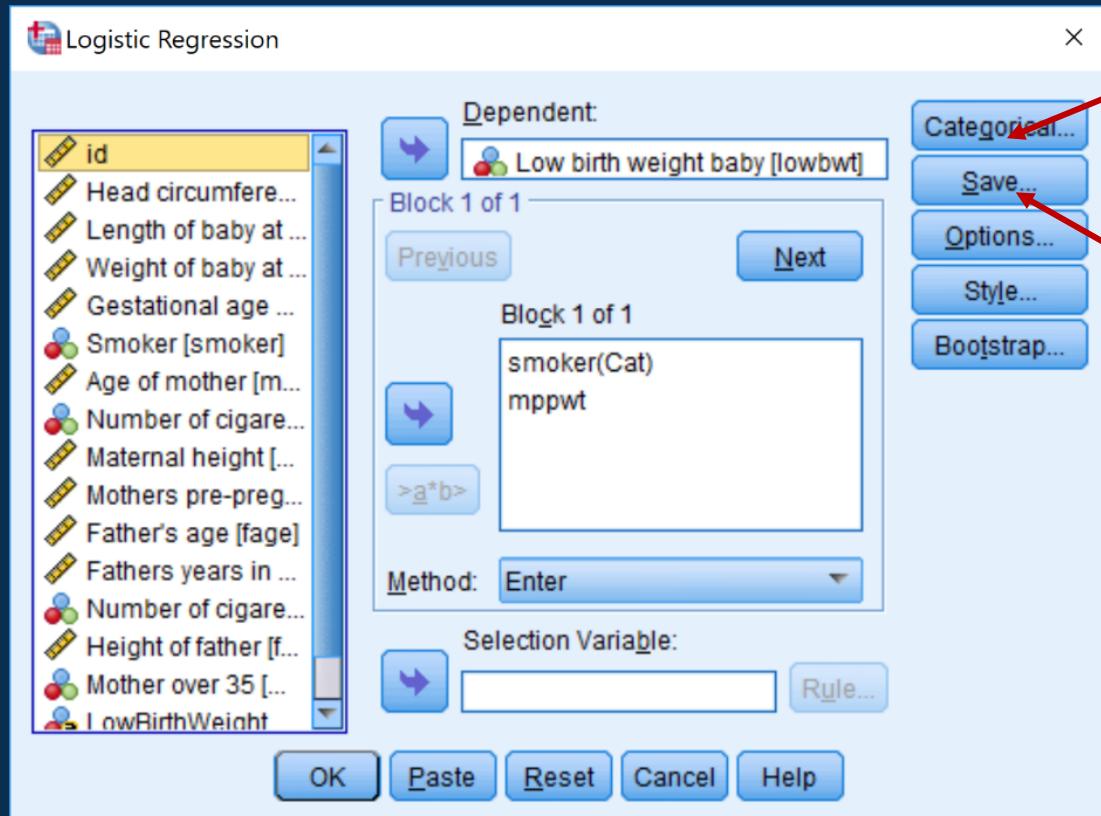
Analyse -> Regression> Binary Logistic



# SPSS slide: 'how to'

**Step 2:** Define any categorical variables and choose the Reference category

**Step 3:** In Options choose the CI for exp ( $\beta$ )



Three dialog boxes are shown, each with numbered callouts:

- Logistic Regression: Define Categorical Variables**: Shows 'smoker(Indicator)' in the 'Categorical Covariates:' list. Callout 5 points to this list.
- Logistic Regression: Options**: Shows the 'Statistics and Plots' section with 'CI for exp(B)' checked at 95%. Callout 8 points to this checkbox.
- Logistic Regression: Options**: Shows the 'Display' section with 'At each step' selected. Callout 9 points to the 'Continue' button at the bottom.



# Output and Interpretation

| Omnibus Tests of Model Coefficients |            |    |      |
|-------------------------------------|------------|----|------|
|                                     | Chi-square | df | Sig. |
| Step 1 Step                         | 8.573      | 2  | .014 |
| Block                               | 8.573      | 2  | .014 |
| Model                               | 8.573      | 2  | .014 |

A p-value (sig) of less than 0.05 for block means that the final model is a significant improvement to the constant only model. (**chi-square=8.573, df=2, p=.014**)

**Nagelkerke R<sup>2</sup> = 24.8%** of the variation in lowbwt can be explained by the final model.

| Model Summary |                     |                      |                     |
|---------------|---------------------|----------------------|---------------------|
| Step          | -2 Log likelihood   | Cox & Snell R Square | Nagelkerke R Square |
| 1             | 48.791 <sup>a</sup> | .185                 | .248                |

a. Estimation terminated at iteration number 4 because parameter estimates changed by less than .001.

| Observed           |                       | Predicted |     |                       | Percentage Correct |
|--------------------|-----------------------|-----------|-----|-----------------------|--------------------|
|                    |                       | No        | Yes | Low birth weight baby |                    |
| Step 1             | Low birth weight baby | No        | 19  | 5                     | 79.2               |
|                    |                       | Yes       | 7   | 11                    | 61.1               |
| Overall Percentage |                       |           |     |                       | 71.4               |

a. The cut value is .500

The correct classification rate has increased by **14.3% to 71.4%**



# Output and Interpretation

| Variables in the Equation |                                    |       |       |       |      |        |                     |              |
|---------------------------|------------------------------------|-------|-------|-------|------|--------|---------------------|--------------|
|                           | B                                  | S.E.  | Wald  | df    | Sig. | Exp(B) | 95% C.I. for EXP(B) |              |
|                           |                                    |       |       |       |      |        | Lower               | Upper        |
| Step 1 <sup>a</sup>       | Smoker(1)                          | 1.575 | .709  | 4.936 | 1    | .026   | 4.831               | 1.204 19.386 |
|                           | Mothers pre-pregnancy weight (lbs) | -.040 | .023  | 3.130 | 1    | .077   | .961                | .919 1.004   |
|                           | Constant                           | 3.898 | 2.840 | 1.884 | 1    | .170   | 49.306              |              |

a. Variable(s) entered on step 1: Smoker, Mothers pre-pregnancy weight (lbs).

## Regression Equation

$$\ln \frac{p}{1-p} = 3.898 + 1.575smoker + -0.040mppwt$$

Odds ratio for the effect of mothers who smoked during pregnancy on low birth weight  $\text{Exp}(\beta) = 4.831$  once adjusted for mothers pre-pregnancy wgt (lbs). Mothers who smoke during pregnancy have **a 4.831 times larger** odds of having a baby born with low birth weight compared to a mother who did not smoke during pregnancy adjusting for mother's pre-pregnancy weight. This was a significant association **95%CI 1.204 to 19.386, p=0.026**.

**One lbs increase** in mothers pre-pregnancy weight would lead to **a 4% reduction ( $\exp(\beta) = 0.961$ )** in the odds of having a baby of low birth weight, if the mother is a non-smoker. **This is not a significant association 95% CI (0.919 to 1.004), p=0.077**

# Reference Categories and Dummy Variables

---

- Categorical Independent **dichotomous** variables:
  - E.g. Gender defined at birth
  - One category is treated as a baseline, or reference category.
  - Reference Category is arbitrarily coded 0, comparison group coded 1
- Categorical independent variables with **more than two levels** need to be recoded into **dummy** variables
  - A “**dummy variable**” is a numerical variable used in regression analysis to represent subgroups of the sample in your study.
  - E.g. Variable X has three levels, create two new variables, each comparing one level to the baseline or reference category
  - Coding represents a contrast between categories.



# Building Models

---

Which predictor variables should I include?

- Literature
- Researcher theory
- Iterative Multivariable Logistic Regression
  - Often have too many variables to legitimately include in the logistic regression model.
    - At least 50 times as many subjects as predictors
  - Used to find a good subset of variables
    - A subset that includes only **statistically** significant predictors and that results in good negative and positive predictive values (more about this in the next section).
- Forward, backward Stepwise regression



# Model Building Strategies

The log likelihood (LL), the deviance (-2LL), or the likelihood ratio(LR) give an overall goodness of fit measurement for the model.

## Forward Selection

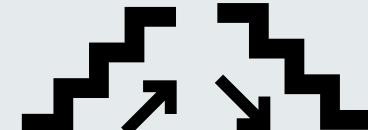
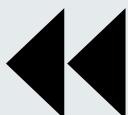
- Variables are tested one at a time.
- First variable added has the smallest LR (and is statistically significant).
- Other variables added if their LR is also significant when adjusted for other variables in the model.
- Model Building stops.
  - All variables have been entered.
  - LR is non-significant for all variables not entered.

## Backward Selection

- Start with all the predictors (significant and not significant).
- Variables are tested one at a time.
- First variable removed has a LR with the largest probability that is greater than alpha.
- Continue until only statistically significant variables remain.

## Stepwise Selection

- Combination of forward and backward.
- Each variable is tested for entry to the model.
- When a predictor is entered, other variables are tested for removal.
- Continue until no more variables can be entered or removed.



# Knowledge Check

Q1. The researcher was also interested to see if the length of gestation had a impact on low birth weight of babies alongside other factors already tested. Interpret these results.

| Omnibus Tests of Model Coefficients |            |       |      |      |
|-------------------------------------|------------|-------|------|------|
|                                     | Chi-square | df    | Sig. |      |
| Step 1                              | Step       | 9.078 | 3    | .028 |
|                                     | Block      | 9.078 | 3    | .028 |
|                                     | Model      | 9.078 | 3    | .028 |

| Model Summary |                     |                      |                     |
|---------------|---------------------|----------------------|---------------------|
| Step          | -2 Log likelihood   | Cox & Snell R Square | Nagelkerke R Square |
| 1             | 48.286 <sup>a</sup> | .194                 | .261                |

a. Estimation terminated at iteration number 4 because parameter estimates changed by less than .001.

| Variables in the Equation |                                    |       |       |       |      |        |                     |              |
|---------------------------|------------------------------------|-------|-------|-------|------|--------|---------------------|--------------|
|                           | B                                  | S.E.  | Wald  | df    | Sig. | Exp(B) | 95% C.I. for EXP(B) |              |
| Step 1 <sup>a</sup>       | Smoker(1)                          | 1.557 | .715  | 4.746 | 1    | .029   | 4.746               | 1.169 19.271 |
|                           | Mothers pre-pregnancy weight (lbs) | -.037 | .023  | 2.496 | 1    | .114   | .964                | .921 1.009   |
|                           | Gestational age at birth (weeks)   | -.100 | .141  | .497  | 1    | .481   | .905                | .686 1.194   |
|                           | Constant                           | 7.326 | 5.701 | 1.651 | 1    | .199   | 1519.126            |              |

a. Variable(s) entered on step 1: Smoker, Mothers pre-pregnancy weight (lbs), Gestational age at birth (weeks).



# Knowledge Check Solutions

Q1. The chi-square is significant (chi-square=9.078, df=3, p=0.028) so our new model is significantly better. Nagelkerke's R<sup>2</sup> suggests that the model explains roughly 26.1% of the variation in the outcome.

For every unit increase in the length of gestation, the odds of a mother having a lowbwt baby is decreased by 9.5%, 95% CI of the odds ( 0.686, 1.194) adjusting for Mothers smoking status and and mothers pre-pregnancy weight, this result was statistically non significant (Wald = 0.497, df=1, p=0.481)

| Omnibus Tests of Model Coefficients |       |            |    |      |
|-------------------------------------|-------|------------|----|------|
|                                     |       | Chi-square | df | Sig. |
| Step 1                              | Step  | 9.078      | 3  | .028 |
|                                     | Block | 9.078      | 3  | .028 |
|                                     | Model | 9.078      | 3  | .028 |

| Model Summary   |                     |                      |                     |
|---|---------------------|----------------------|---------------------|
| Step  | -2 Log likelihood   | Cox & Snell R Square | Nagelkerke R Square |
| 1   | 48.286 <sup>a</sup> | .194                 | .261                |
| a. Estimation terminated at iteration number 4 because parameter estimates changed by less than .001. |                     |                      |                     |

| Variables in the Equation |                                    |       |       |       |      |        |                     |              |
|---------------------------|------------------------------------|-------|-------|-------|------|--------|---------------------|--------------|
|                           | B                                  | S.E.  | Wald  | df    | Sig. | Exp(B) | 95% C.I. for EXP(B) |              |
| Step 1 <sup>a</sup>       | Smoker(1)                          | 1.557 | .715  | 4.746 | 1    | .029   | 4.746               | 1.169 19.271 |
|                           | Mothers pre-pregnancy weight (lbs) | -.037 | .023  | 2.496 | 1    | .114   | .964                | .921 1.009   |
|                           | Gestational age at birth (weeks)   | -.100 | .141  | .497  | 1    | .481   | .905                | .686 1.194   |
|                           | Constant                           | 7.326 | 5.701 | 1.651 | 1    | .199   | 1519.126            |              |

a. Variable(s) entered on step 1: Smoker, Mothers pre-pregnancy weight (lbs), Gestational age at birth (weeks).



# References

---

Field, Andy. Discovering statistics using IBM SPSS statistics. Sage, 2013. (Chapter 19)

Agresti, Alan. Categorical data analysis. John Wiley & Sons, 2014.



# Thank you

Contact details/for more information:

Zahra Abdulla

[Zahra.abdulla@kcl.ac.uk](mailto:Zahra.abdulla@kcl.ac.uk)

Department of Biostatistics and Health Informatics (BHI)

IoPPN



**Zahra Abdulla**

---

Department: Biostatistics and Health  
Informatics

Institute of Psychiatry, Psychology and Neuroscience  
08/2020

**Module Title:** Introduction to Statistics

**Session Title:** Prediction, Goodness of Fit and Classification

---

**Topic title:** Binary Logistic Regression



---

After working through this session, you should be able to:

- Make predictions and describe these as probabilities
- Assess the goodness of fit of the model.

# Prediction

---

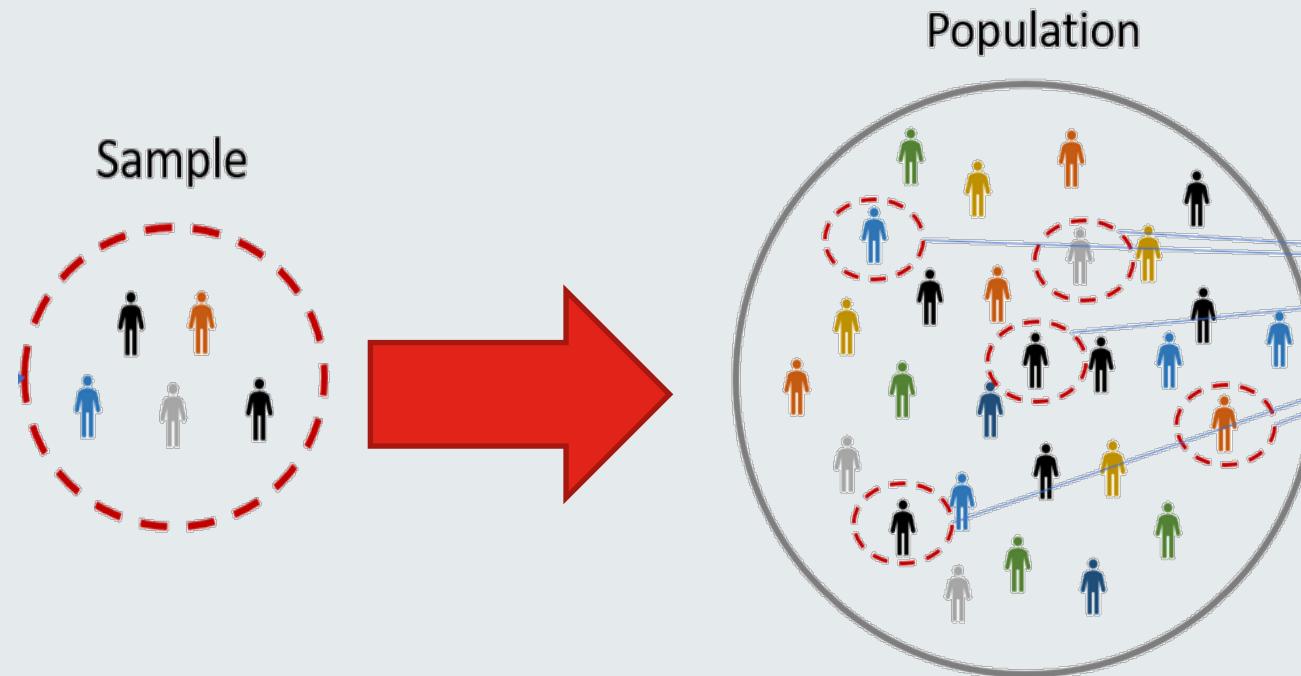
- A logistic regression model can be used to make predictions
- The prediction is the value of the linear predictor
- We need to obtain the odds of the person experiencing an event - exponentiate the linear predictor.
- To get the probability you rearrange the odds equation.



# Why is prediction important?

---

- Because we're modelling!
- We want to make predictions about what would happen in the general population at a given point



# The logistic transformation: a recap

$$\ln\left(\frac{\pi}{1-\pi}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_i x_i$$

This is just the *odds*.

The (adjusted) odds ratio is the estimated change in odds for a unit change in  $x_1$  (holding  $x_2, x_3, \dots, x_i$  constant)

$$L = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_i x_i$$

This is called the **Linear Predictor**

$$\exp(L) = e^L$$

This is the **Odds of an event**

$$\hat{\pi} = \frac{odds}{1+odds} \quad \hat{\pi} = \frac{\exp(L)}{1+\exp(L)} = \frac{1}{1+\exp(-L)}$$

This is the **Estimated Probability of an event**

# Estimating the probability of an event

What is the probability of a person starting smoking, if when they were born cigarettes cost £2?

We know

$$\log\left(\frac{\pi}{1-\pi}\right) = L, \text{ where } L = 3.69 - 0.07x$$

To calculate the probability of starting smoking, as per the conditions above

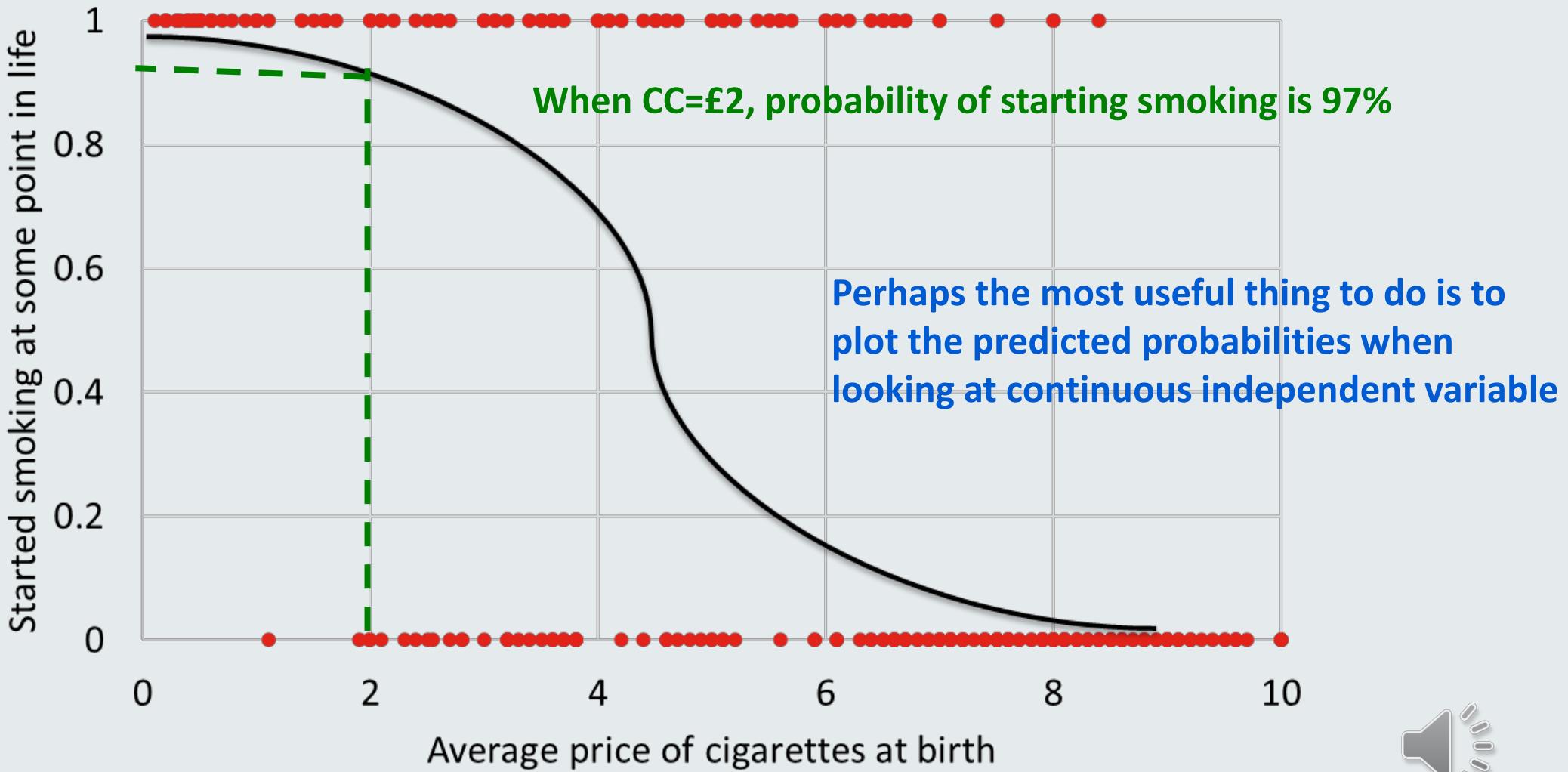
$$\hat{\pi} = \frac{\exp(L)}{1+\exp(L)}$$

$$\hat{\pi} = \frac{e^{3.69-0.07x}}{1+e^{3.69-0.07x}}$$

$$\hat{\pi} = \frac{e^{3.69-0.07\times 2}}{1+e^{3.69-0.07\times 2}} = 0.97$$



# Thinking about prediction



# Estimating probabilities

---

What is the probability of a mother whose pre-pregnancy weight is 110 LLbs and a smoker of having a baby of low birth weight?

The **Linear Predictor (L)** is given by

$$L = 3.898 + 1.575 \times Smoker - 0.040 \times Mppwgt$$

$$L = 3.898 + (1.575 \times 1) - (0.040 \times 110)$$

$$L = 1.073$$

The **Probability (P)** is given by

$$P = \frac{\exp(L)}{1 + \exp(L)}$$

$$P = \frac{\exp(1.073)}{1 + \exp(1.073)}$$

$$P = \frac{2.924}{3.924}$$

$$P = 0.745$$

## Interpretation

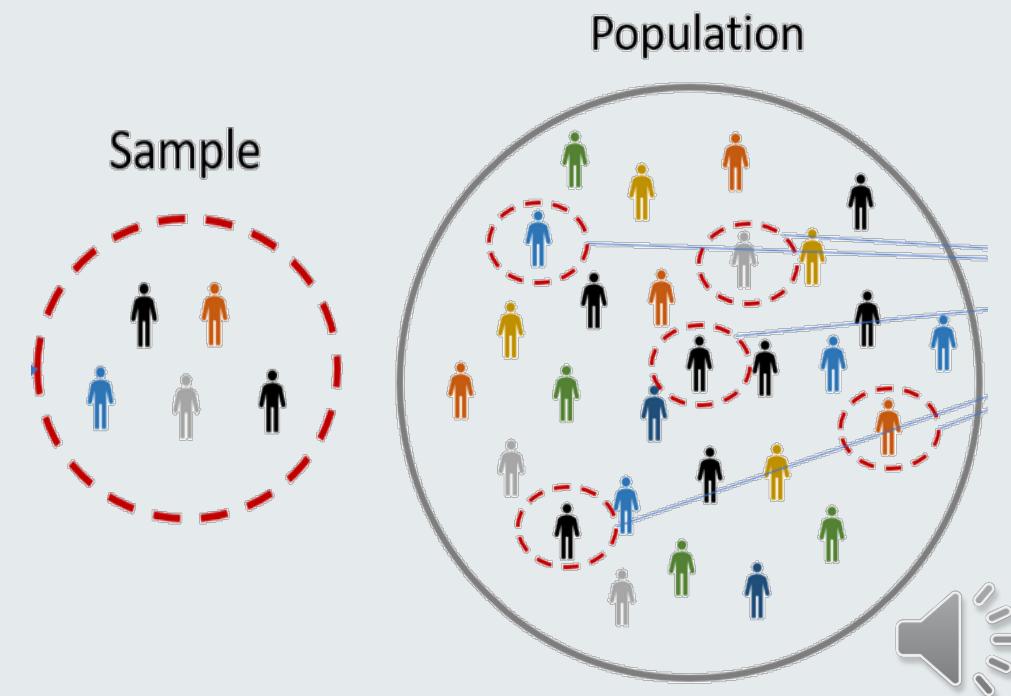
**The probability of a baby born with a low birth weight is 74.5%**



# Goodness of fit

---

- Goodness-of-Fit tests help determine if observed data aligns with what is expected in the actual population.
- More specifically, it is used to test if sample data fits a distribution from a certain population (e.g., a population with a normal distribution)
- Remember, we're still modelling...



# Goodness of fit

---

Here we will discuss two ways of assessing goodness of fit:

1. Classification analysis
2. Hosmer and Lemeshow test



# Classification Analysis

---

One way of assessing goodness of fit is to use a **classification table**.

This allows us to evaluate **predictive accuracy** of the logistic regression model.

Classification tables are useful because they provide information that allow us to consider goodness of fit in different ways e.g., specificity and sensitivity (we will come back to these).

They are built on regression models used to predict **probability** of an outcome. When we use classification tables we identify a **threshold probability**, beyond which, an outcome is expected.

For example, if we want to identify a threshold probability, beyond which, a healthcare worker is encouraged to remove a breathing tube from an intensive care patient – we could do this based on a regression model in which we predict the probability of success, when removing a breathing tube, under different conditions.



# An example with birth weight

|                    |                       | Classification Table <sup>a</sup> |     | Percentage<br>Correct |  |
|--------------------|-----------------------|-----------------------------------|-----|-----------------------|--|
|                    |                       | Predicted                         |     |                       |  |
|                    |                       | Low birth weight baby             |     |                       |  |
| Observed           |                       | No                                | Yes |                       |  |
| Step 1             | Low birth weight baby | No                                | 15  | 9                     |  |
|                    |                       | Yes                               | 5   | 13                    |  |
| Overall Percentage |                       |                                   |     | 66.7                  |  |

a. The cut value is .500

So, following our regression model, the observed values for the DV and the predicted values are cross-classified.

We can then **classify individuals** by saying that all individuals with a predicted value higher than a certain threshold probability are positive i.e. will have babies with a low birth weight.

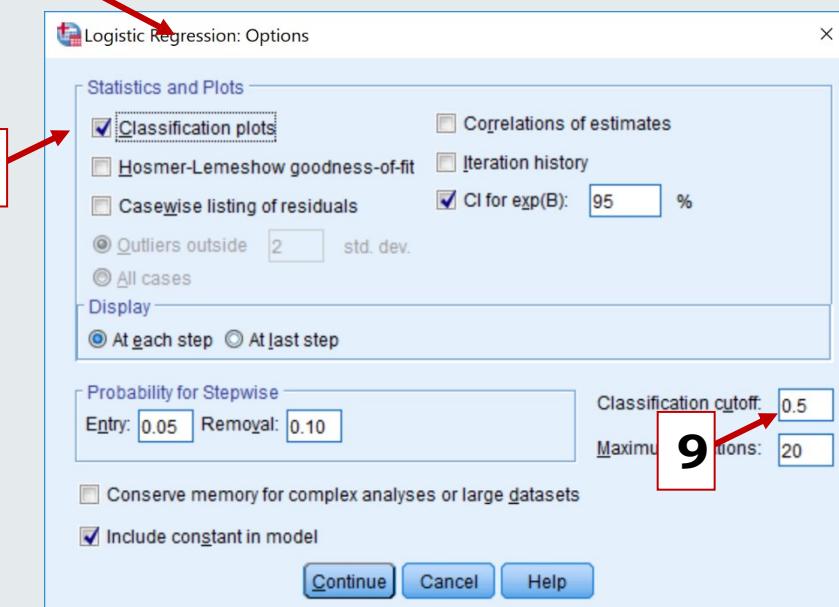
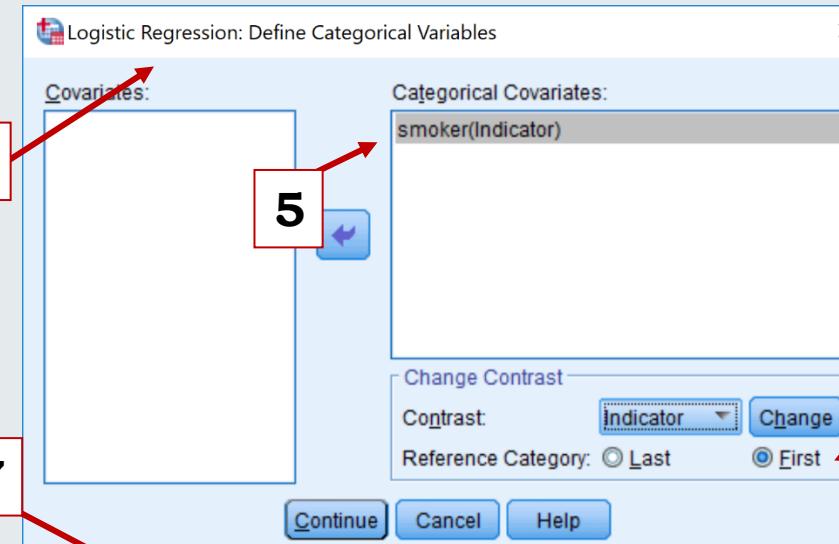
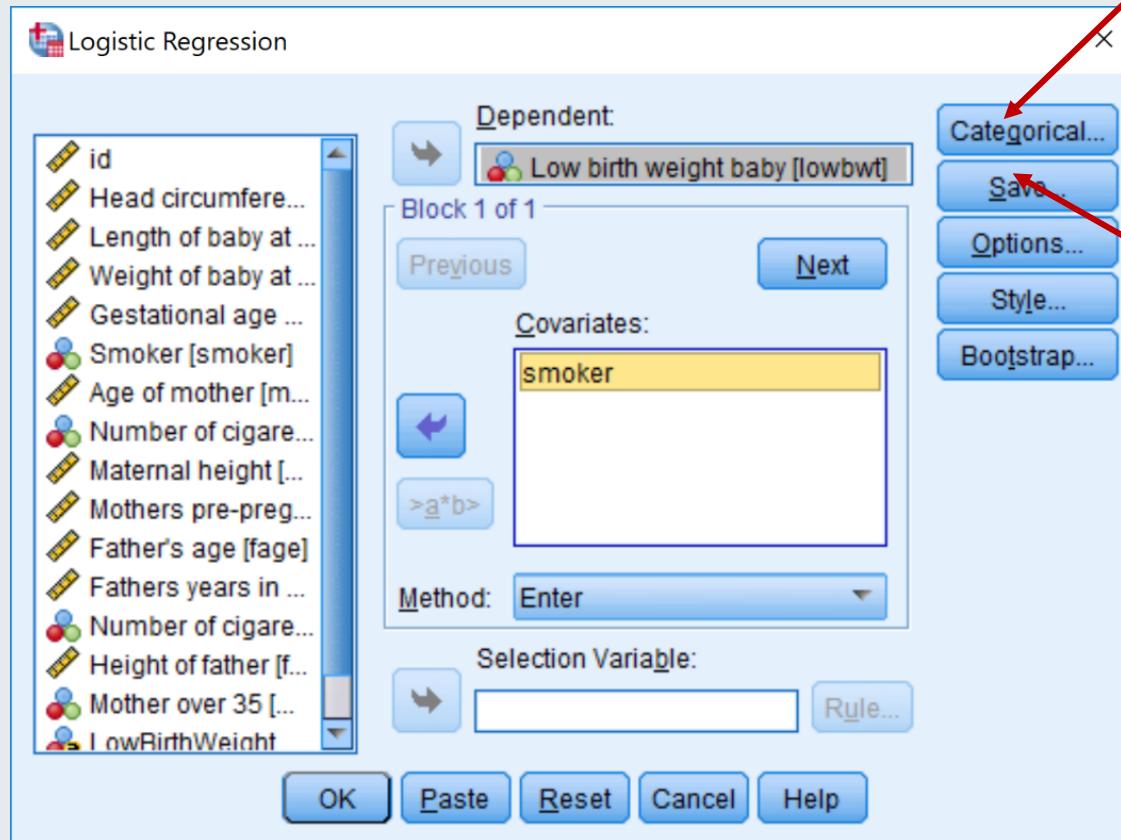
- For every individual we use the linear predictor to estimate their probability of having a binary outcome (e.g., babies of low birth weight)
- Based on some cut-off probability we classify them as positive or negative
- Cross tabulate the predicted values versus the true values



# SPSS slide: 'how to'

Step 1: Use the appropriate test, here: 'Binary Logistic regression'.

Step 2: Under Options choose 'Classification Plots'



# Classification Table

Based on a cut-off of 0.5, 62.5% of those without low birth weight are correctly predicted to be negative and 72.2% of those with babies with low birth weight is correctly predicted to be positive.

| Observed           |                       | Predicted             |        | Percentage<br>Correct |
|--------------------|-----------------------|-----------------------|--------|-----------------------|
|                    |                       | Low birth weight baby | No Yes |                       |
| Step 1             | Low birth weight baby | No                    | 15 9   | 62.5                  |
|                    |                       | Yes                   | 5 13   | 72.2                  |
| Overall Percentage |                       |                       |        | 66.7                  |

a. The cut value is .500

"The cut value is .500". This means that if the probability of a case being classified into the "yes" category is greater than .500, then that particular case is classified into the "yes" category. Otherwise, the case is classified as in the "no" category (as mentioned previously).



# Sensitivity and specificity

---

In order to choose a threshold probability to turn a probability model into a classification model we usually consider the quantities **sensitivity** and **specificity**

**Sensitivity**, which is the percentage of cases that had the observed characteristic (e.g., "yes" for baby with low birth weight) which were correctly predicted by the model (i.e., true positives).

**Specificity**, which is the percentage of cases that did not have the observed characteristic (e.g., "no" for baby with low birth weight) and were also correctly predicted as not having the observed characteristic (i.e., true negatives).

In an ideal world we would like to maximise both sensitivity and specificity, but there is often a trade-off

We select an optimal threshold by considering what degree of sensitivity and specificity are acceptable

# Positive and negative predicted values

---

We can also use the classification table to look at **positive and negative predictive values**

Remember again we're still modelling...

**The positive predictive value** is the percentage of correctly predicted cases "with" the observed characteristic compared to the total number of cases predicted as having the characteristic.

**The negative predictive value** is the percentage of correctly predicted cases "without" the observed characteristic compared to the total number of cases predicted as not having the characteristic.



# How can I calculate these!?

|                             | Outcome successful   | Outcome unsuccessful |
|-----------------------------|----------------------|----------------------|
| Classification successful   | True positives (TP)  | False positives (FP) |
| Classification unsuccessful | False negatives (FN) | True negatives (TN)  |

The formulae for the various quantities are as follows:

$$\text{Sensitivity} = \frac{TP}{(TP+FN)}$$

$$\text{Specificity} = \frac{TN}{(FP+TN)}$$

$$\text{Positive Predictive Value (PPV)} = \frac{TP}{(TP+FP)}$$

$$\text{Negative Predictive Value (NPV)} = \frac{TN}{(FN+TN)}$$



# Calculation

| Observed           |                       | Predicted             |    | Percentage<br>Correct |
|--------------------|-----------------------|-----------------------|----|-----------------------|
|                    |                       | Low birth weight baby | No |                       |
| Step 1             | Low birth weight baby | No                    | 15 | 9                     |
|                    | Yes                   | 5                     | 13 | 72.2                  |
| Overall Percentage |                       |                       |    | 66.7                  |

a. The cut value is .500

**Percentage Accuracy in Classification (PAC)** is the overall percentage of cases correctly classified by the model =  $(15 + 13) / (15+9+5+13) = 66.7$

**Sensitivity** =  $13 / (13+5) = 72.2 \%$

**Specificity** =  $15 / (9+15) = 62.5 \%$

**Positive Predictive Value (PPV)** =  $13 / (13+9) = 29.1\%$

**Negative Predictive Value (NPV)** =  $15 / (5+15) = 75\%$



# Interpretation

| Observed           |                       | Predicted             |    | Percentage<br>Correct |
|--------------------|-----------------------|-----------------------|----|-----------------------|
|                    |                       | Low birth weight baby | No |                       |
| Step 1             | Low birth weight baby | No                    | 15 | 9                     |
|                    | Yes                   | Yes                   | 5  | 13                    |
| Overall Percentage |                       |                       |    | 66.7                  |

a. The cut value is .500

Overall, the model correctly classified 66.7% of the cases. Sensitivity, 72.2% is high compared to specificity, which is 62.5%. The positive predictive value, computed for low-birth-weight baby, is 29.1%; the negative predictive value, computed for no low-birth-weight baby, is 75%. The low PPV may be indicative that the model is not a good predictor of low birth weight, as only 29.1% of cases predicted to have a baby of low birthweight had babies of low birthweight.



# Hosmer and Lemeshow Goodness of fit

---

Another way of assessing goodness of fit is (i.e., is our model any good?) is to use a [Hosmer and Lemeshow test](#).

This is a [statistical test for goodness of fit](#) for the logistic regression model.

The data are divided into approximately ten groups defined by increasing order of estimated risk.

The observed and expected number of cases in each group is calculated and a [Chi-squared statistic](#) is produced.

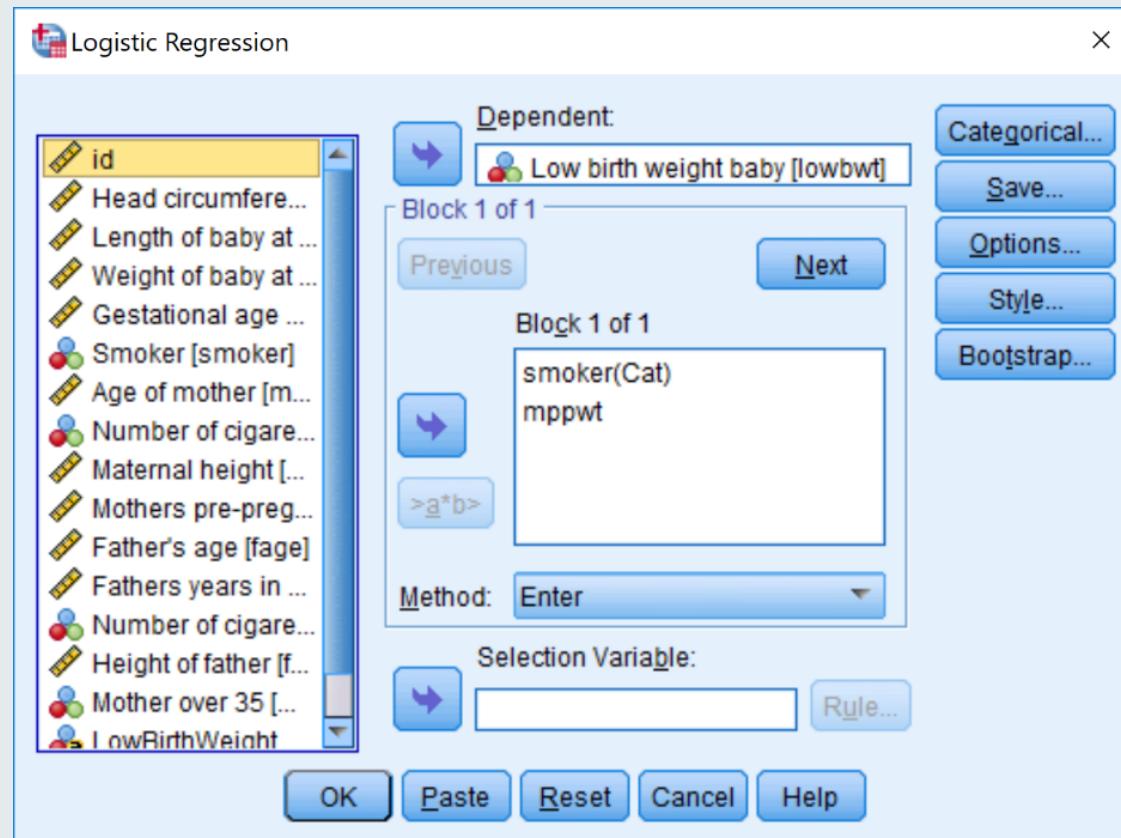
**You can only do this test with multiple predictors**



# SPSS slide: 'how to'

**Step 1:** Use the appropriate test, here: 'Binary Logistic regression'.

**Step 2:** Under Options choose "Hosmer-Lemeshow.."



**Logistic Regression: Define Categorical Variables**

- Covariates:** (empty)
- Categorical Covariates:** smoker(Indicator)
- Change Contrast:**
  - Contrast: Indicator
  - Reference Category: Last (radio button)
  - First (radio button) **6**

**Logistic Regression: Options**

**Statistics and Plots:**

- Classification plots
- Hosmer-Lemeshow goodness-of-fit **8**
- Casewise listing of residuals
- Correlations of estimates
- Iteration history
- CI for exp(B): 95 %

**Display:**

- Outliers outside 2 std. dev.
- All cases
- At each step
- At last step

**Probability for Stepwise:**

- Entry: 0.05
- Removal: 0.10

**Classification cutoff:** 0.7

**Maximum Iterations:** 20

Conserve memory for complex analyses or large datasets

Include constant in model

Buttons: Continue, Cancel, Help.



# Hosmer and Lemeshow Goodness of fit

| Hosmer and Lemeshow Test |            |    |      |
|--------------------------|------------|----|------|
| Step                     | Chi-square | df | Sig. |
| 1                        | 7.199      | 8  | .515 |

Null hypothesis: The model is consistent with the data. i.e. **a non-significant p-value indicates good fit.**

A large value of Chi-squared (with small p-value < 0.05) indicates poor fit and small Chi-squared values (with larger p-value closer to 1) indicate a good logistic regression model fit.

The Contingency Table for Hosmer and Lemeshow Test table shows the details of the test with observed and expected number of cases in each group



## To conclude...

---

You should now be able to analyse data using binary logistic regressions

You should be able to run binary logistic regressions adjusting for covariates

You should understand goodness of fit

You should be able to make predictions based on data with dichotomous outcomes and continuous predictors

# Knowledge Check

---

What is the probability of a mother whose pre pregnancy weight is 210lbs and a non-smoker of having a baby of low birth weight?

If we were to raise the cutoff to 0.70, how well is the model predicting babies with low birth weight?



# Knowledge Check Solutions

---

What is the probability of a mother whose pre pregnancy weight is 210lbs and a non-smoker of having a baby of low birth weight?

$$\begin{aligned}L &= 3.898 + 1.575 \times Smoker - 0.040 \times Mppwgt \\&= 3.898 + (1.575 \times 0) - (0.040 \times 210) \\&= -4.502\end{aligned}$$

$$\begin{aligned}P &= (\exp(L))/(1+\exp(L)) \\&= (\exp(-4.502))/(1+\exp(-4.502)) \\&= (0.0112)/(1.0112) \\&= 0.0112 = 1.1\%\end{aligned}$$



# Knowledge Check Solutions

If we were to raise the cutoff to 0.70, how well is the model predicting babies with low birth weight?

|                    |                                   | Classification Table <sup>a</sup> |    | Predicted             |   | Percentage Correct |
|--------------------|-----------------------------------|-----------------------------------|----|-----------------------|---|--------------------|
|                    |                                   |                                   |    | Low birth weight baby |   |                    |
| Step 1             | Observed<br>Low birth weight baby | No                                | 22 | Yes                   | 2 | 91.7               |
|                    |                                   | Yes                               | 14 |                       | 4 | 22.2               |
| Overall Percentage |                                   |                                   |    |                       |   | 61.9               |

a. The cut value is .700

Based on a cut-off of 0.7, 91.7% of those without low birth weight are correctly predicted to be negative but only 22.2% of those with babies with low birth weight is correctly predicted to be positive.

# References

---

Field, Andy. Discovering statistics using IBM SPSS statistics. Sage, 2013. (Chapter 19)

Agresti, Alan. Categorical data analysis. John Wiley & Sons, 2014.

Binomial Logistic Regression using SPSS Statistics, Laerd Statistics.

<https://statistics.laerd.com/spss-tutorials/binomial-logistic-regression-using-spss-statistics.php>: Accessed 05/01/22



# Thank you

Contact details/for more information:

Zahra Abdulla

[Zahra.abdulla@kcl.ac.uk](mailto:Zahra.abdulla@kcl.ac.uk)

Dr. Silia Vitoratou

[silia.vitoratou@kcl.ac.uk](mailto:silia.vitoratou@kcl.ac.uk)

Dr Raquel Iniesta

[raquel.iniesta@kcl.ac.uk](mailto:raquel.iniesta@kcl.ac.uk)

Department of Biostatistics and Health Informatics (BHI)

IoPPN