# Intro to Inferential Statistics with R

Workshop 5

Course: VSK1004 Applied Researcher

# Workshop structure (draft)

| MONDAY<br>Intro to Statistic Inference | WEDNESDAY<br>More about inferential stats | TODAY<br>Linear & Logistic Regression |
|---|---|---|
| 1. Descriptive vs Inferential statistics<br>2. Population, sample and sampling distribution<br>3. Null Hypothesis testing<br>4. Correlation and interpretation | 1. Choosing a statistical test<br>2. t-test family<br>3. chi-squared<br>4. correlation<br>5. Chi-squared distribution | 1. Linear Regression<br>2. Multiple Linear Regression<br>3. Model Assumption<br>4. Logistic Regression |

# Our goal in the next 40 min

In this session, we will cover some of the **basic statistical models and its properties such as:**

1. Simple Linear Regression

2. Multiple Linear Regression

3. Linear Model Assumptions

4. Logistic Regression

# Simple Linear Regression

# The mathematical equation

Simple Linear Regression

$$y = b_0 + b_1 * x_1$$

# The mathematical equation

Simple
Linear
Regression

$$y = b_0 + b_1 * x_1$$

Dependent variable (DV)

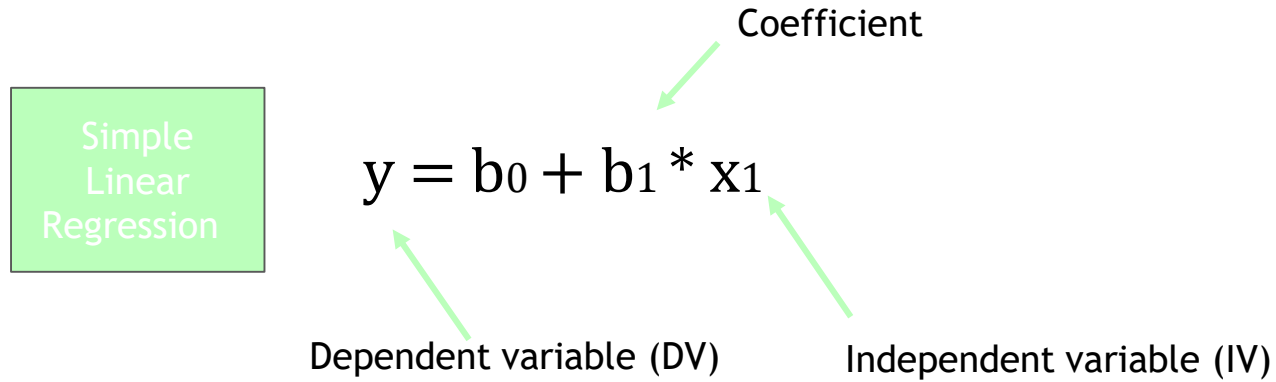# The mathematical equation

Simple Linear Regression

$$y = b_0 + b_1 * x_1$$
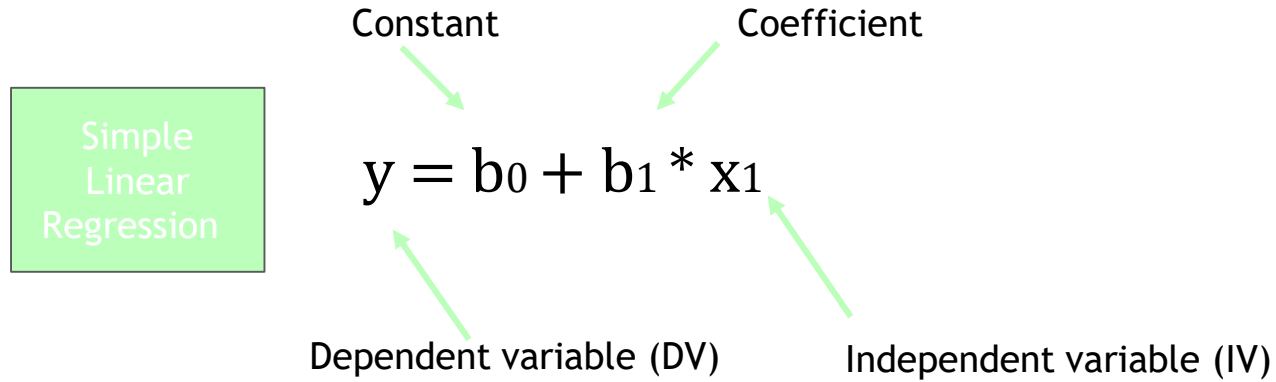
Dependent variable (DV)

Independent variable (IV)

# The mathematical equation

Coefficient

Simple Linear Regression

$$y = b_0 + b_1 * x_1$$
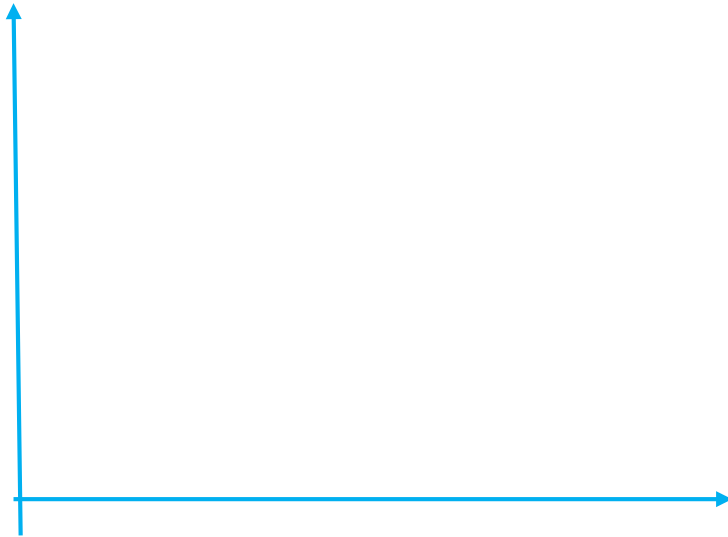
Dependent variable (DV)

Independent variable (IV)

# The mathematical equation

Constant

Coefficient

Simple Linear Regression

$$y = b_0 + b_1 * x_1$$

Dependent variable (DV)

Independent variable (IV)

# Look at the simple Linear regression

Simple Linear Regression:

# Representation in a graph: x and y axis
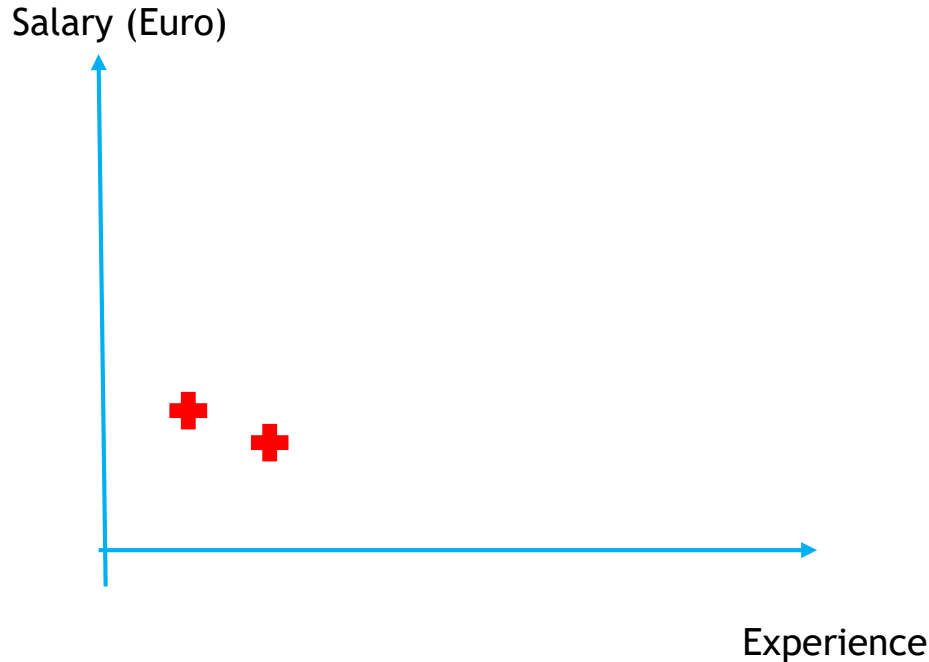
Simple Linear Regression:

Salary (Euro)

Experience

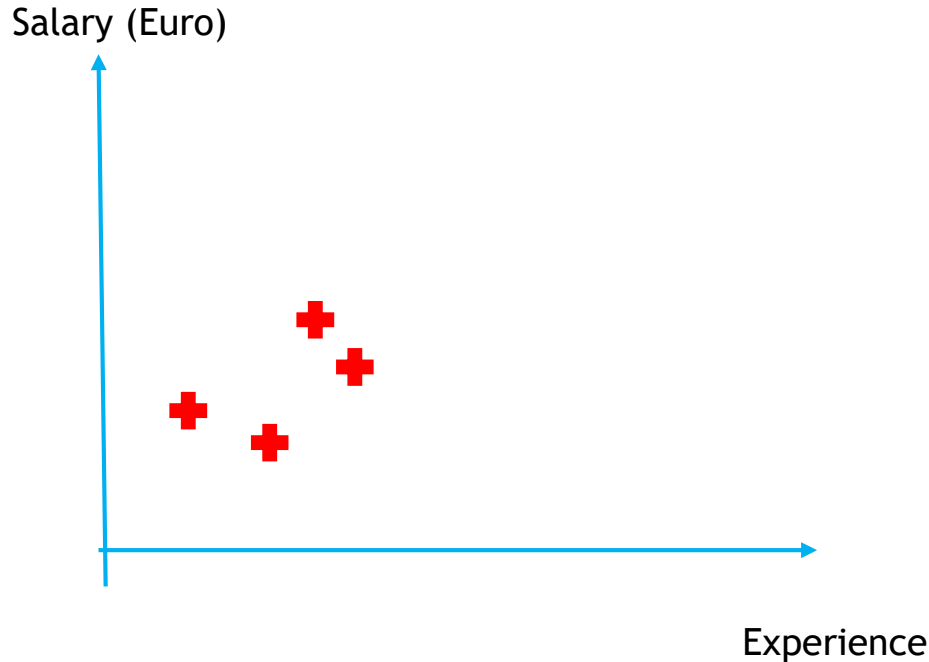# Repesentation in a graph: dots as observed data

Simple Linear Regression:

Salary (Euro)



Experience

# Repesentation in a graph: dots as observed data

Simple Linear Regression:

Salary (Euro)

Experience

*Intro to Inferential statistics with R* *c.utrillaguerrero@maastrichtuniversity.nl*

# Repesentation in a graph: dots as observed data

Simple Linear Regression:

Salary (Euro)

Experience

# Repesentation in a graph: dots as observed data

Simple Linear Regression:

Salary (Euro)



Experience

*Intro to Inferential statistics with R*  *c.utrillaguerrero@maastrichtuniversity.nl*

# how salary is distributed among people

## Simple Linear Regression:



Salary (Euro)

Experience

# Repesentation in a graph: observed data

Simple Linear Regression:

Salary (Euro)

$$y = b_0 + b_1 * x_1$$

Experience

# Repesentation in a graph with the equation

Simple Linear Regression:



Salary (Euro)

Experience

$$y = b_0 + b_1 * x_1$$

$$Salary = b_0 + b_1 * Experience$$

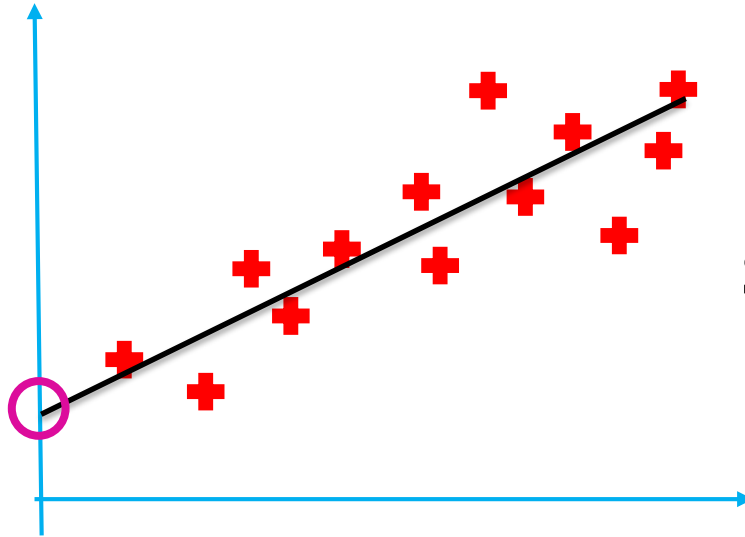# Add best fitting line for linear regression

## Simple Linear Regression:

Salary (Euro)



$$y = b_0 + b_1 * x_1$$

$$Salary = b_0 + b_1 * Experience$$

Experience

12th June 2020
Intro to Inferential statistics with R
c.utrillaguerrero@maastrichtuniversity.nl

# Identify parameters in the graphs: Constant

Simple Linear Regression:

Salary (Euro)

$$y = b_0 + b_1 * x_1$$

$$Salary = b_0 + b_1 * Experience$$

Experience

# Identify parameters in the graphs: Constant

Simple Linear Regression:
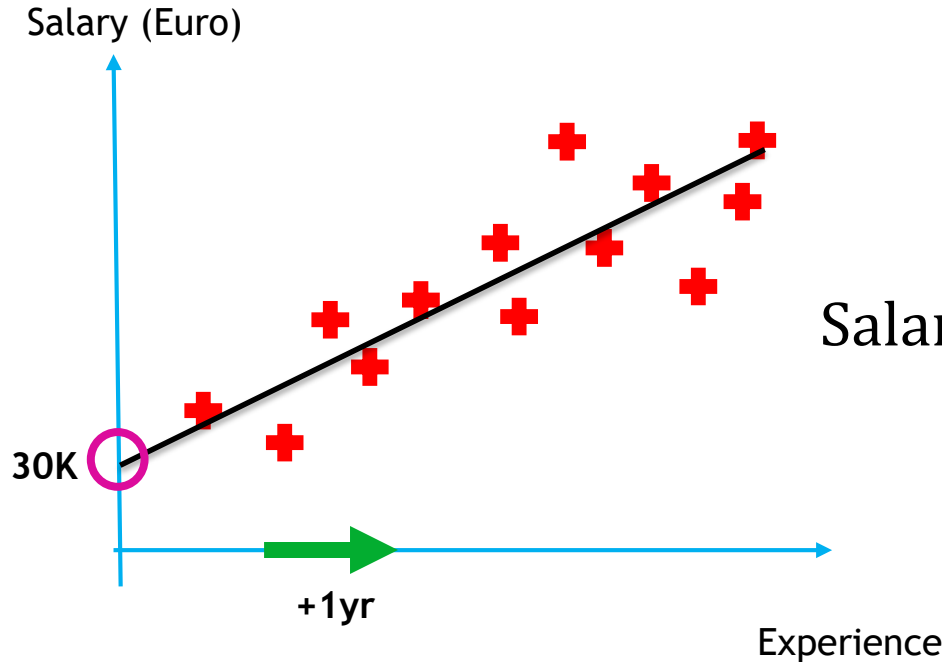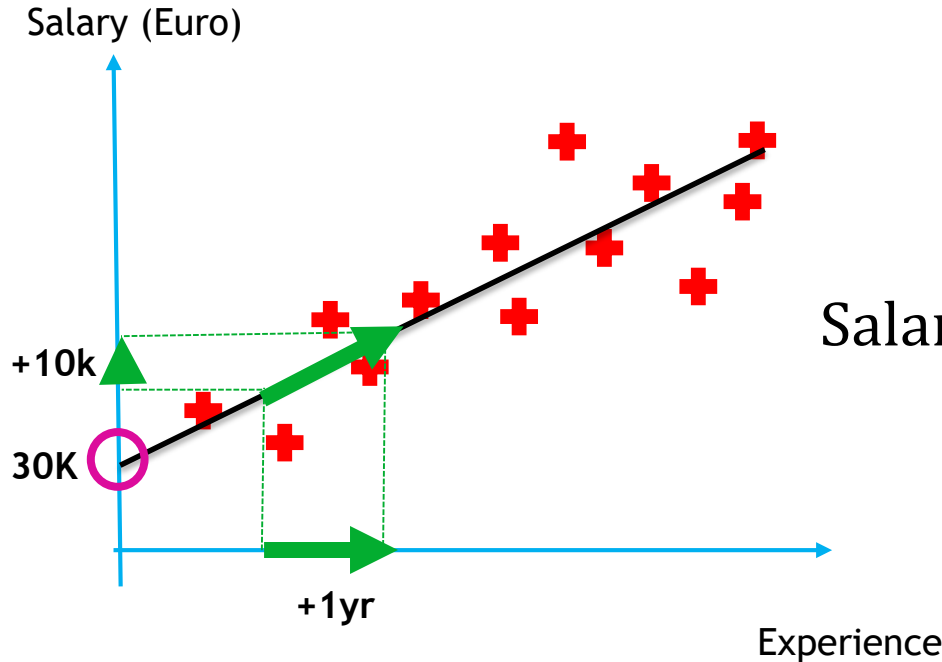
Salary (Euro)

$$y = b_0 + b_1 * x_1$$

$$Salary = b_0 + b_1 * Experience$$

30K

Experience

# Identify parameters in the graphs: Slope

Simple Linear Regression:



$$y = b_0 + b_1 * x_1$$

$$\text{Salary} = b_0 + b_1 * \text{Experience}$$

# Identify parameters in the graphs: Slope

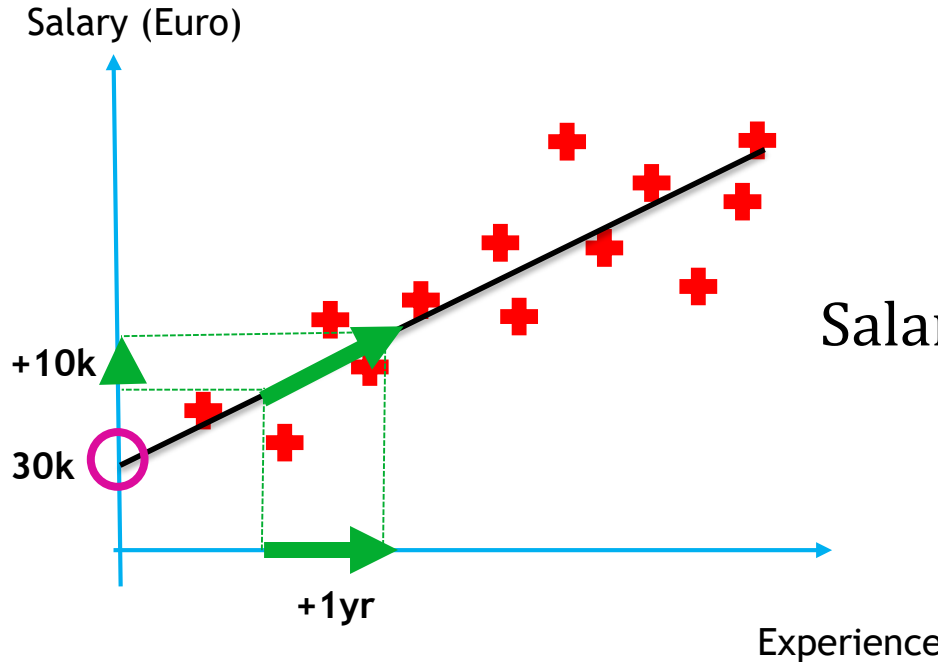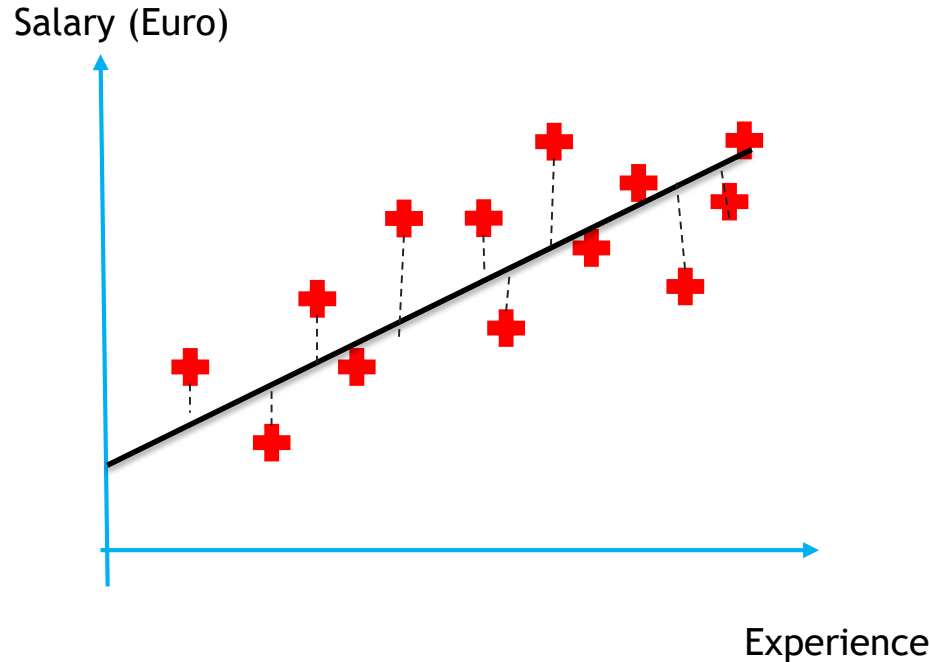## Simple Linear Regression:



$$y = b_0 + b_1 * x_1$$

$$Salary = b_0 + b_1 * Experience$$

# Identify parameters in the graphs: Slope

## Simple Linear Regression:

Salary (Euro)



$$y = b_0 + b_1 * x_1$$

$$\text{Salary} = b_0 + b_1 * \text{Experience}$$

+10k

30K

+1yr

Experience

# All linear parameters in the best fitted line

Simple Linear Regression:



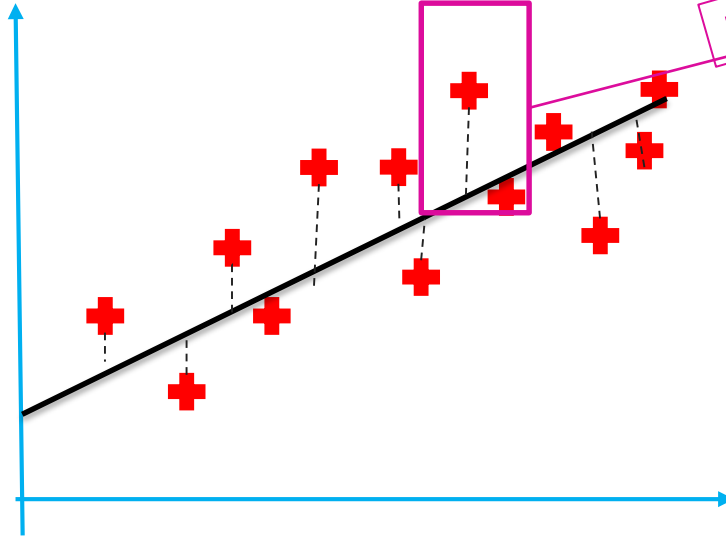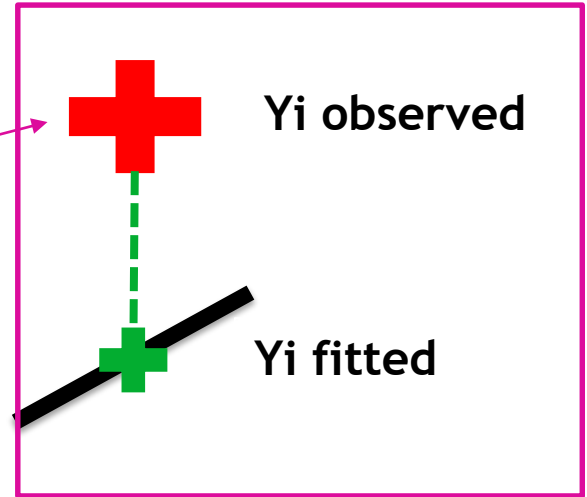$$y = b_0 + b_1 * x_1$$

$$\text{Salary} = b_0 + b_1 * \text{Experience}$$

# How linear regression find best fitting line?

Simple Linear Regression:

Salary (Euro)

Experience

# How linear regression find best fitting line
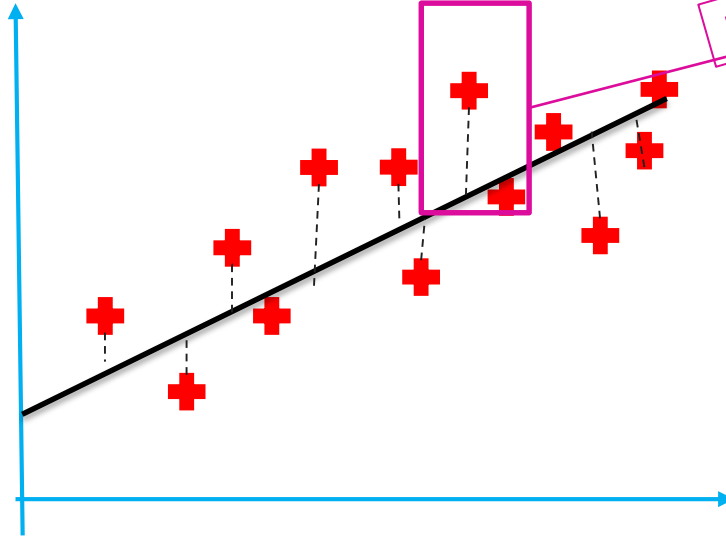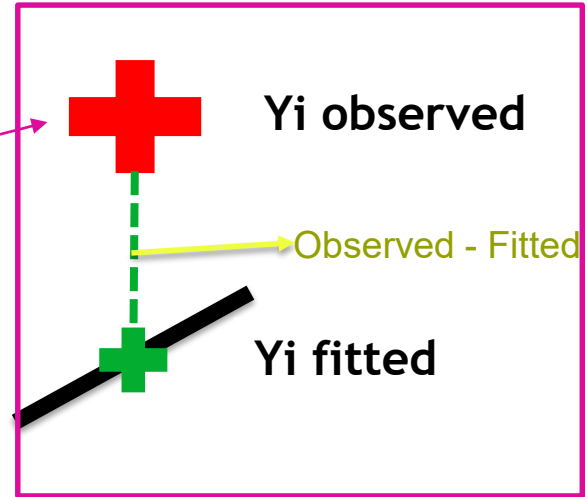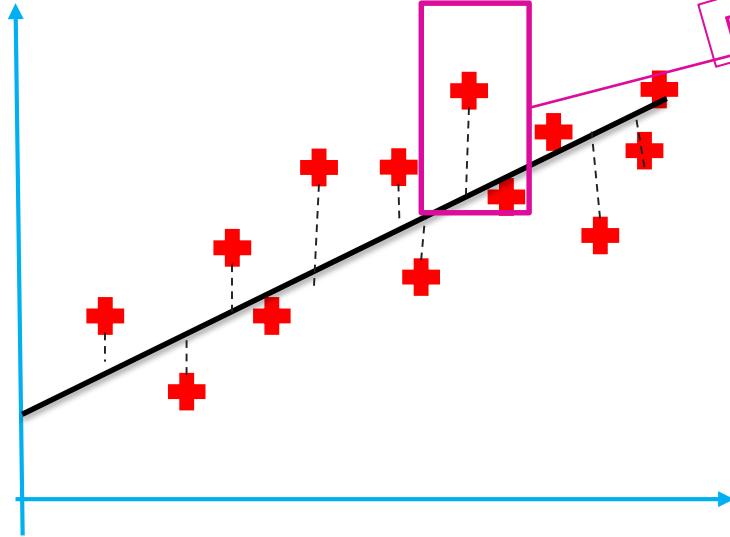
## Simple Linear Regression:

Salary (Euro)

Experience

Example
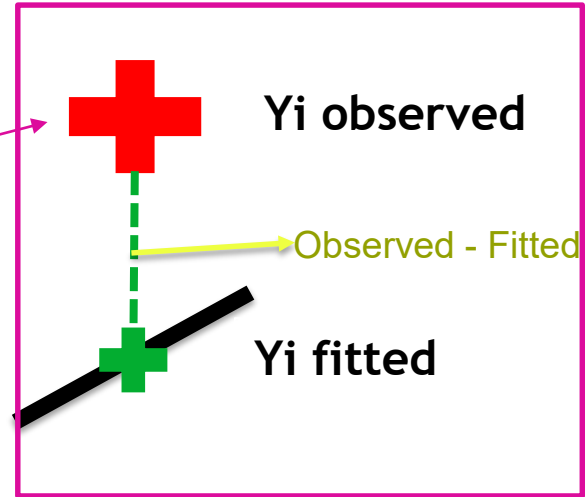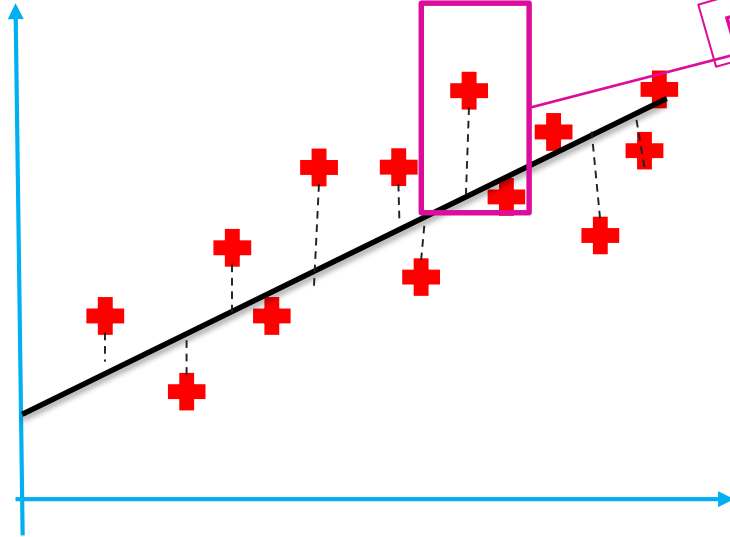
# How linear regression find best fitting line



Simple Linear Regression:

Salary (Euro)

Example

Yi observed

Yi fitted

Experience

# How linear regression find best fitting line

Simple Linear Regression:



Salary (Euro)
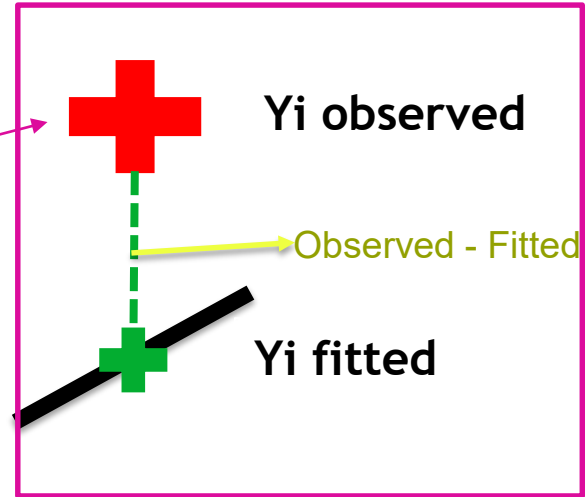
Experience

Example

Yi observed

Observed - Fitted

Yi fitted

# How linear regression find best fitting line

## Simple Linear Regression:

Salary (Euro)

Example

**Yi observed**

Observed - Fitted

**Yi fitted**

SUM (Yi observed - Yi fitted)$^2$

Experience

# How linear regression find best fitting line

## Simple Linear Regression:



Salary (Euro)

Example

Experience

**Yi observed**

Observed - Fitted

**Yi fitted**

**SUM (Yi observed - Yi fitted)²** ➡ *min*

# Linear Regression looks for min sum of squares to find the line which has the smallest sum squares possible, and its called, the best fitting line

*Intro to Inferential statistics with R* *c.utrillaguerrero@maastrichtuniversity.nl*

# Multiple Linear Regression

# Same thing but many variables into the model

| Simple Linear Regression |
|---|

$$y = b_0 + b_1 * x_1$$

| Multiple Linear Regression |
|---|

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + ... + b_n * X_n$$

# Same thing but many variables into the model

| | |
|---|---|
| Simple Linear Regression | $y = b_0 + b_1 * x_1$ |

Dependent variable (DV)

| | |
|---|---|
| Multiple Linear Regression | $y = b_0 + b_1 * x_1 + b_2 * x_2 + ... + b_n * X_n$ |

# Same thing but many variables into the model

Simple Linear Regression

$$y = b_0 + b_1 * x_1$$

Dependent variable (DV)     Independent variable (IVs)

Multiple Linear Regression

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + ... + b_n * x_n$$

# Regressions

| | |
|---|---|
| Simple Linear Regression | $y = b_0 + b_1 * x_1$ |

Dependent variable (DV)    Independent variable (IVs)

| | |
|---|---|
| Multiple Linear Regression | $y = b_0 + b_1 * x_1 + b_2 * x_2 \dots + b_n * X_n$ |

Constant                              Coefficients

# Model Assumptions

# A Caveat: Assumptions of Linear Regression

1. Linearity
2. Homoscedasticity
3. Multivariate Normality
4. Independence of errors
5. Lack of multicollinearity

# Logistic Regression

# What we know

## Linear Regression:

*Intro to Inferential statistics with R*            *c.utrillaguerrero@maastrichtuniversity.nl*

# What we know

## Linear Regression:
## - Simple
$$y = b_0 + b_1 * x$$

# What we know

**Linear Regression:**
**- Simple**

$$y = b_0 + b_1 * x$$

**- Multiple:**

$$y = b_0 + b_1 * x_1 + ... + b_n * X_n$$

# What we know

## We know this:

# What we know

## We know this:

# What we know

## We know this:

**Salary (Euro)**

$$y = b_0 + b_1 * x$$

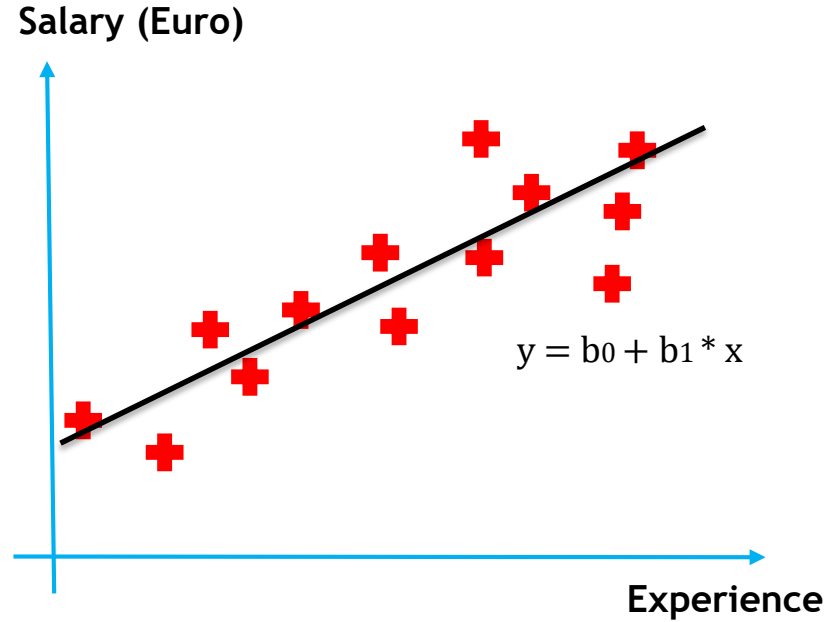**Experience**

# What is new: Logistic regression

## This is new:

## We know this:

**Salary (Euro)**



$$y = b_0 + b_1 * x$$

**Experience**

# Logistic regression

## This is new:

**Action (Y/N)**



1

0

**Age**

## We know this:

**Salary (Euro)**



$y = b_0 + b_1 * x$

**Experience**

# Logistic Regression

## This is new:

**Action (Y/N)**



**???**

**Age**

## We know this:

**Salary (Euro)**



$$y = b_0 + b_1 * x$$

**Experience**

# Logistic Regression

## This is new:

Action (Y/N)



???

Age

## We know this:

Salary (Euro)



$y = b_0 + b_1 * x$

Experience

# Logistic Regression: what if predict the probability or likelihood a person taking the offer?

# Logistic Regression: this part make sense to get probabilities

# Logistic Regression: this part does not make sense to get probabilities (<0>1)

# Logistic Regression

# Logistic Regression

$$y = b_0 + b_1 * x$$

# Logistic Regression

$$y = b_0 + b_1 * x$$

Sigmoid Function
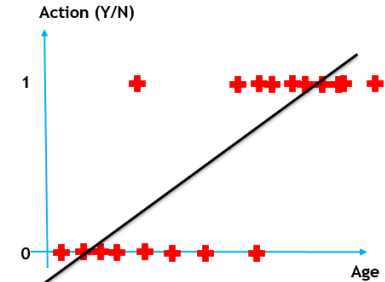
$$p = \frac{1}{1 + e^{-y}}$$



Action (Y/N)

Age

# Logistic Regression

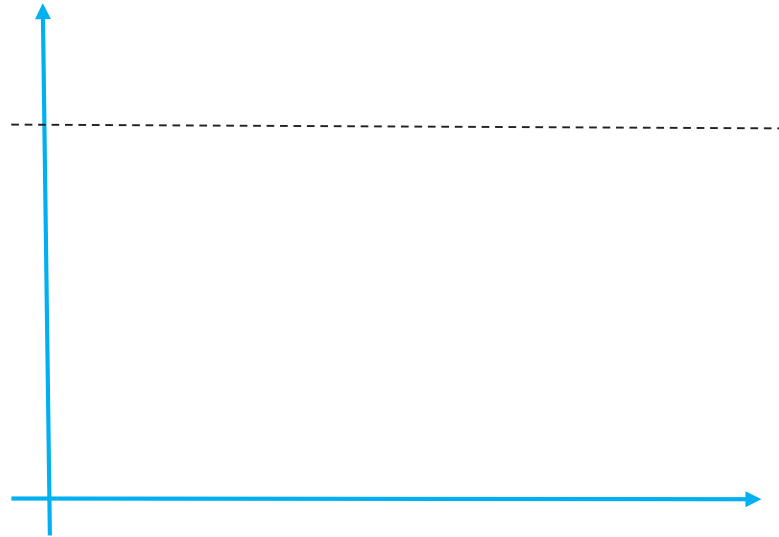$$y = b_0 + b_1 * x$$

Sigmoid Function

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

$$\ln\left(\frac{p}{1-p}\right) = b_0 + b_1 * x$$

# Logistic Regression

$$y = b_0 + b_1 * x$$

Sigmoid Function

$$p = \frac{1}{1 + e^{-y}}$$
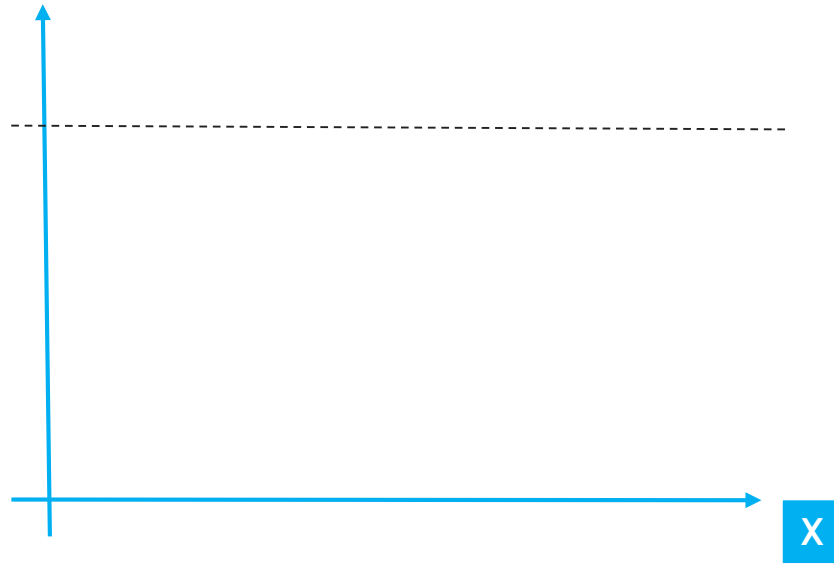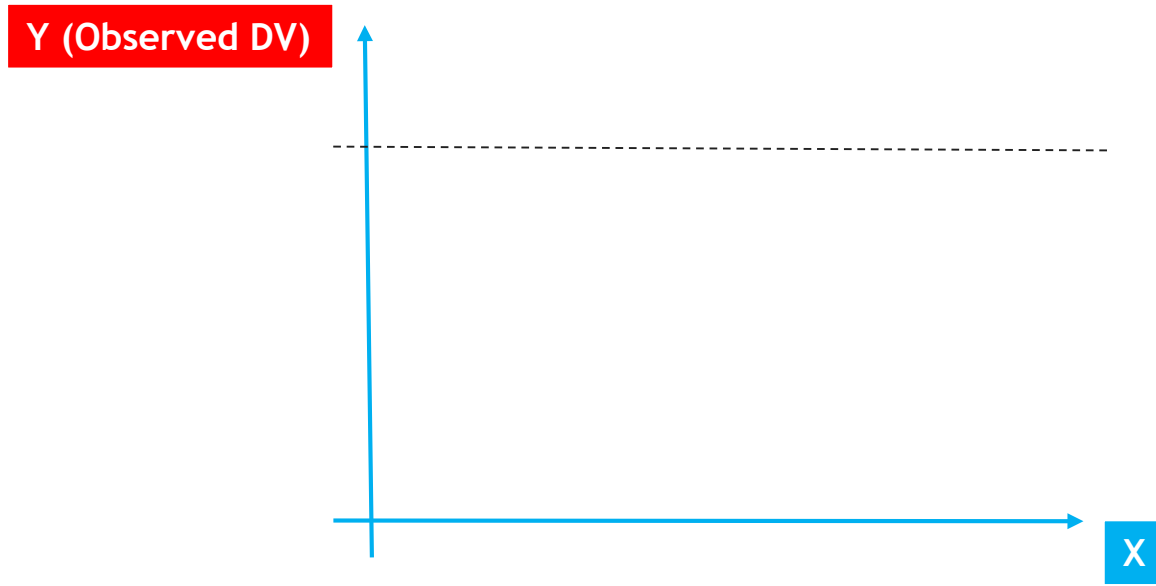
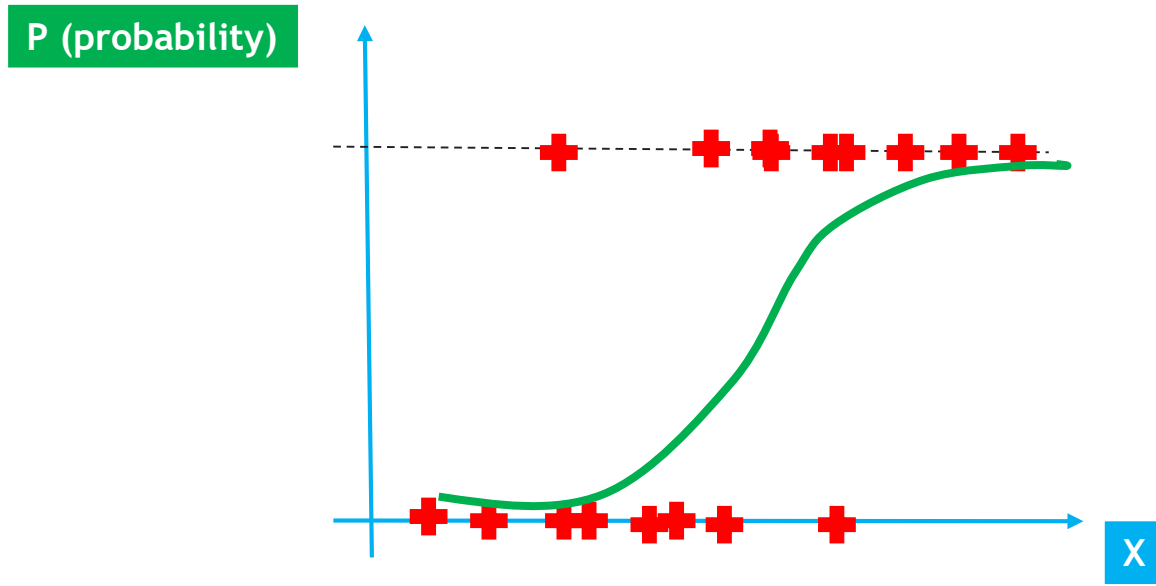$$ln\left(\frac{p}{1-p}\right) = b_0 + b_1 * x$$

# Logistic Regression

# Logistic Regression

# Logistic Regression

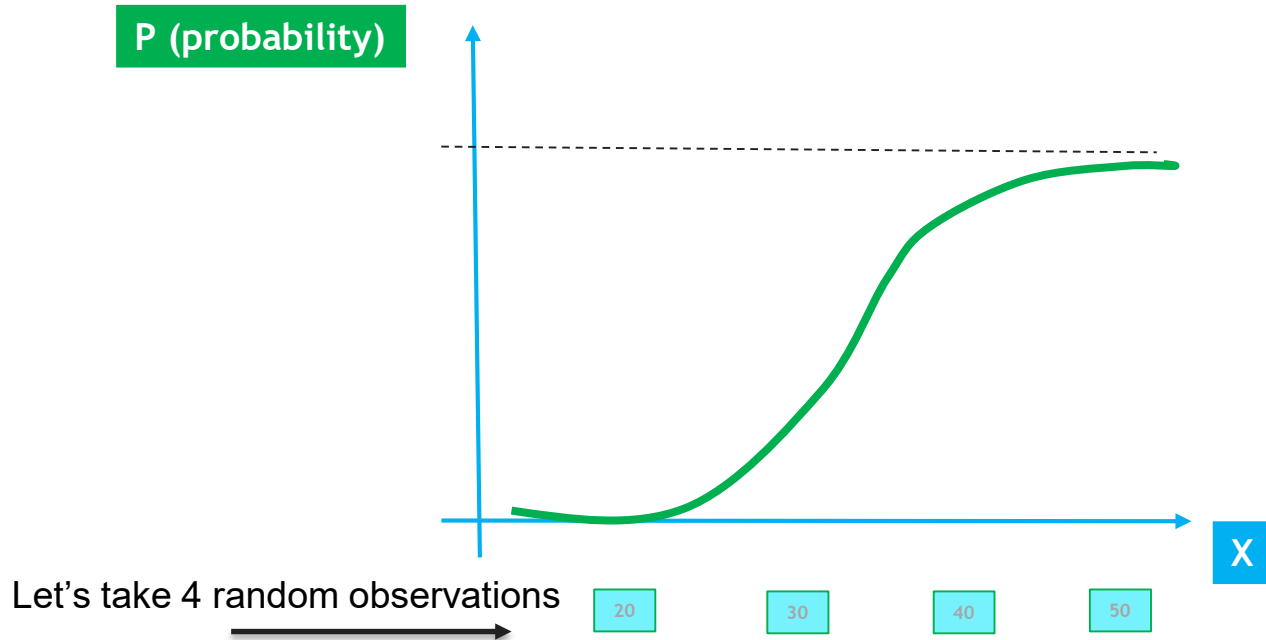Y (Observed DV)

X

# Logistic Regression: best fitting line that fits our data

Y (Observed DV)

$$ln\left(\frac{p}{1-p}\right) = b_0 + b_1 * x$$

X

# Logistic Regression: we created the model

# Logistic Regression: we predict this probability

P (probability)



Let's take 4 random observations

X

20   30   40   50

*Intro to Inferential statistics with R*

*c.utrillaguerrero@maastrichtuniversity.nl*

# Logistic Regression: what we need to do to get the probability of this 4 random variables?

*Intro to Inferential statistics with R*

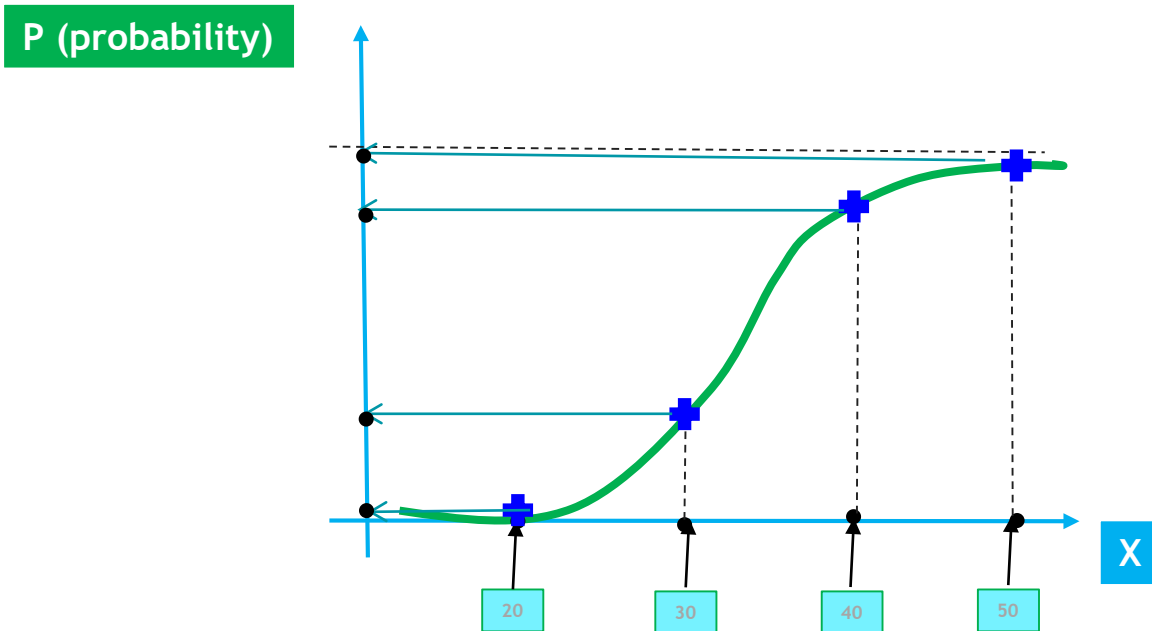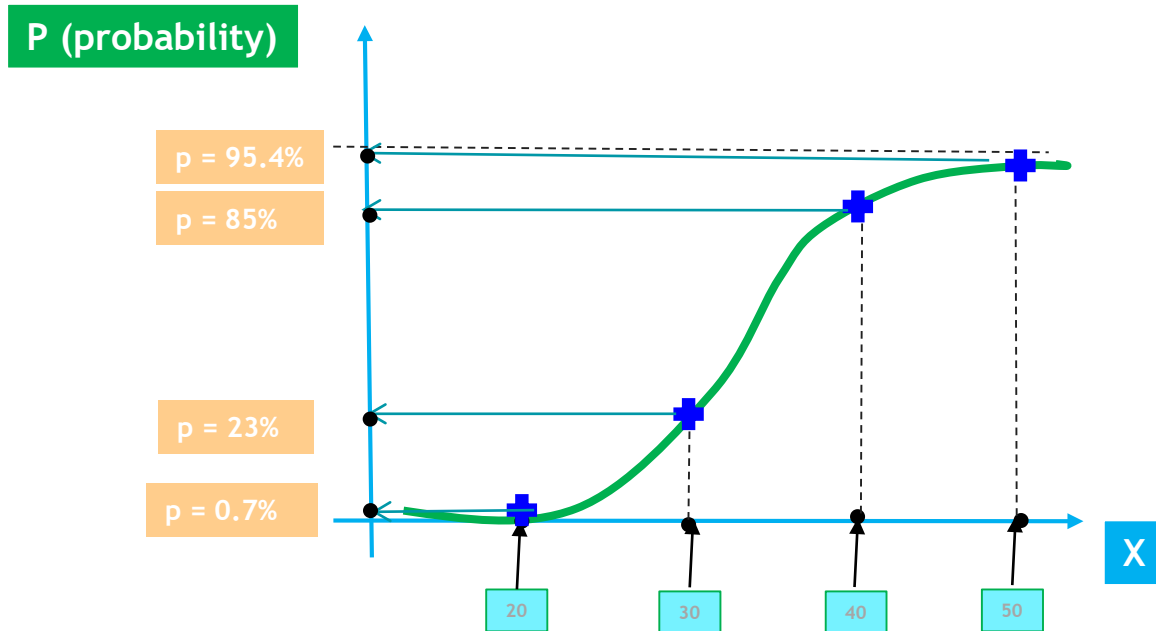*c.utrillaguerrero@maastrichtuniversity.nl*

# Logistic Regression: we need to project them in the curve

# Logistic Regression: if you need the probabilities, then it has to be projected to the left
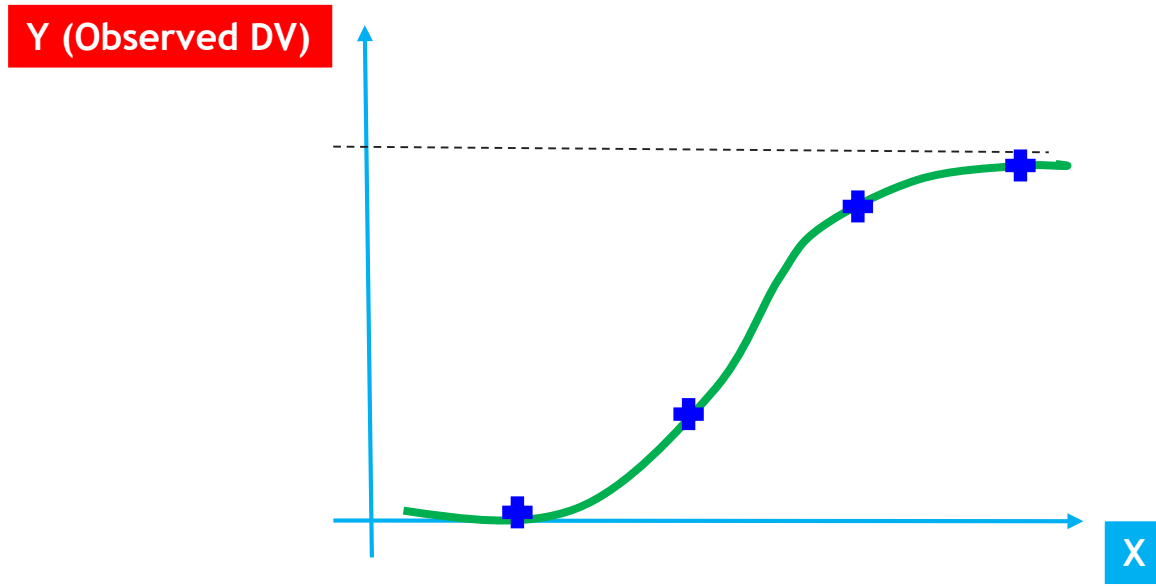
**P (probability)**



X

20   30   40   50

# Logistic Regression: who is the least/most likely to take the offer?

# Logistic Regression: I want a prediction instead probabilities

**P (probability)**
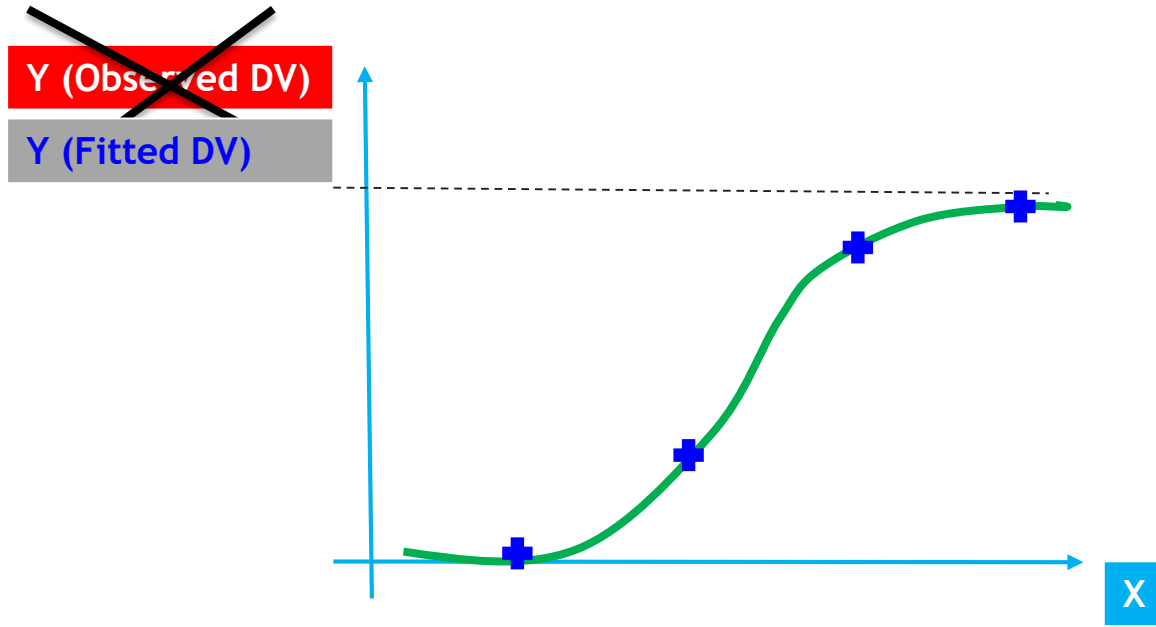


X

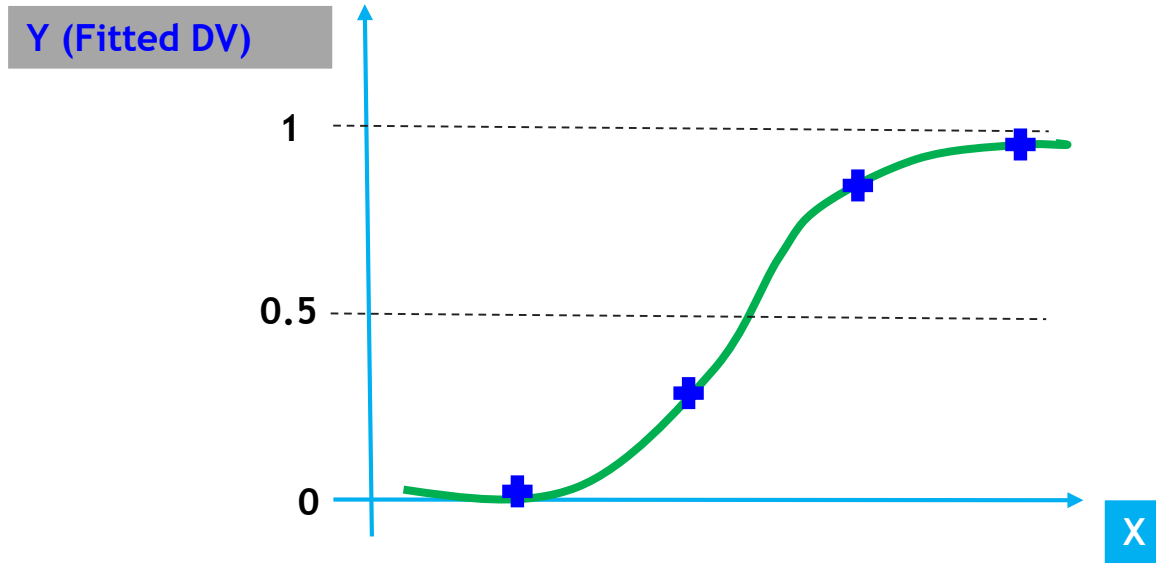# Logistic Regression: get predicted (fitted) values

**Y (Observed DV)**

X

# Logistic Regression: get predicted (fitted) values
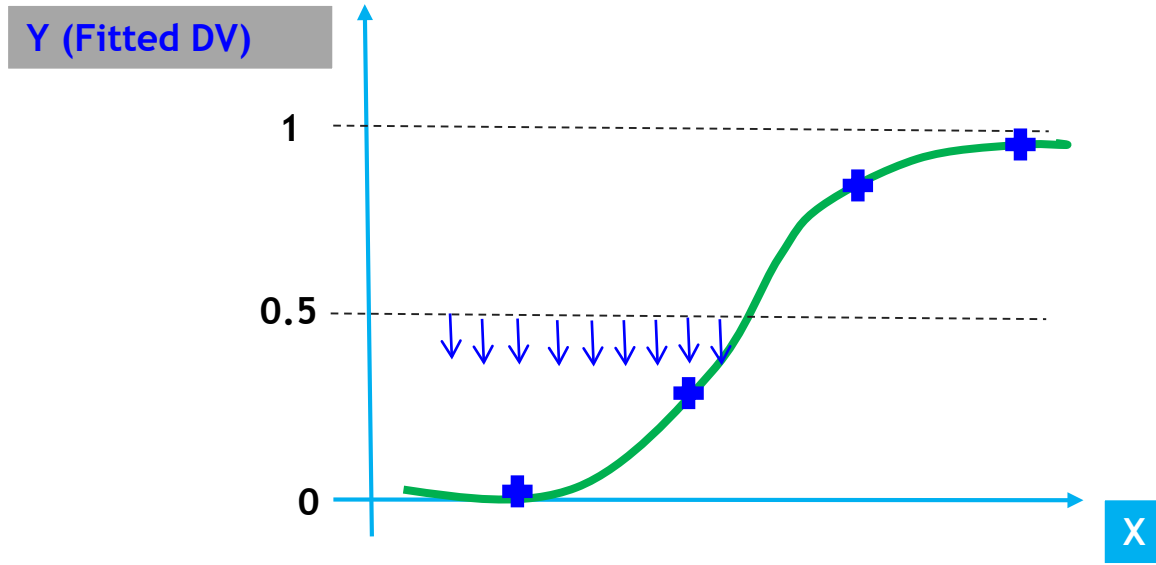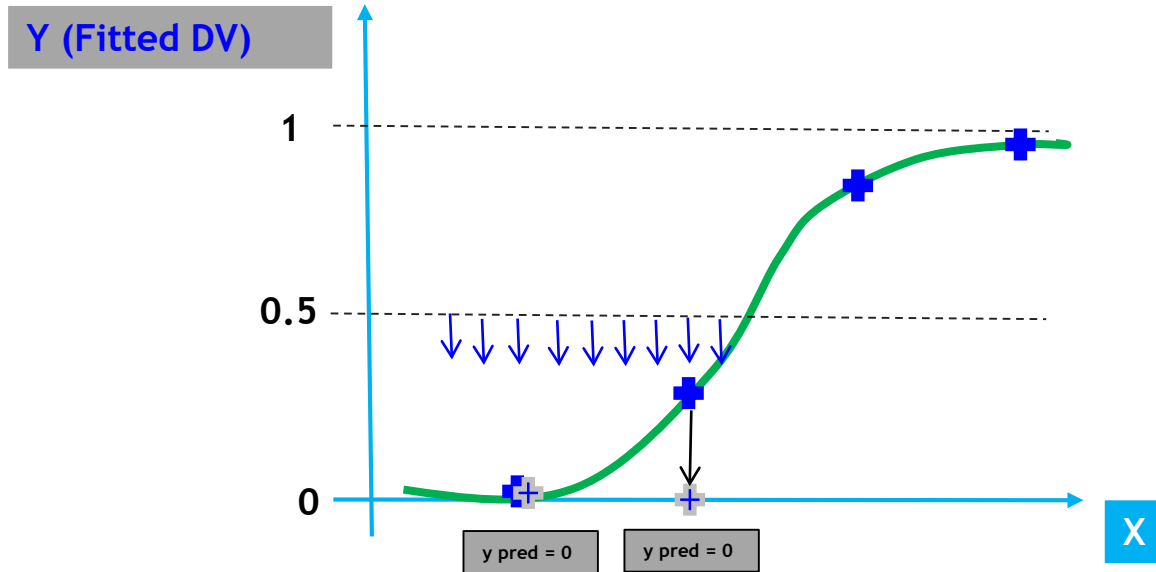
# Logistic Regression: get predicted (fitted) values

Y (Observed DV)

Y (Fitted DV)

X

# Logistic Regression: select a threshold line



**Y (Fitted DV)**

1 - - - - - - - - - - - - - - - - - - - - - - -

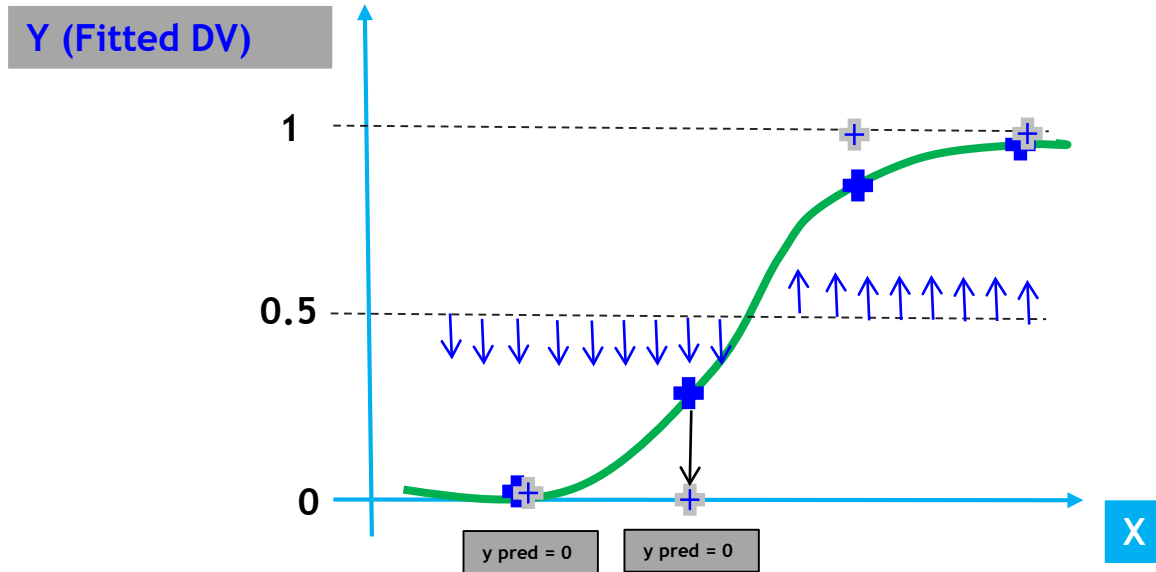0.5 - - - - - - - - - - - - - - - - - - - - - - -

0

X

# Logistic Regression: if predicted probability is less than 50% then we predict 0

Y (Fitted DV)

# Logistic Regression: if predicted probability is less than 50% then we predict 0
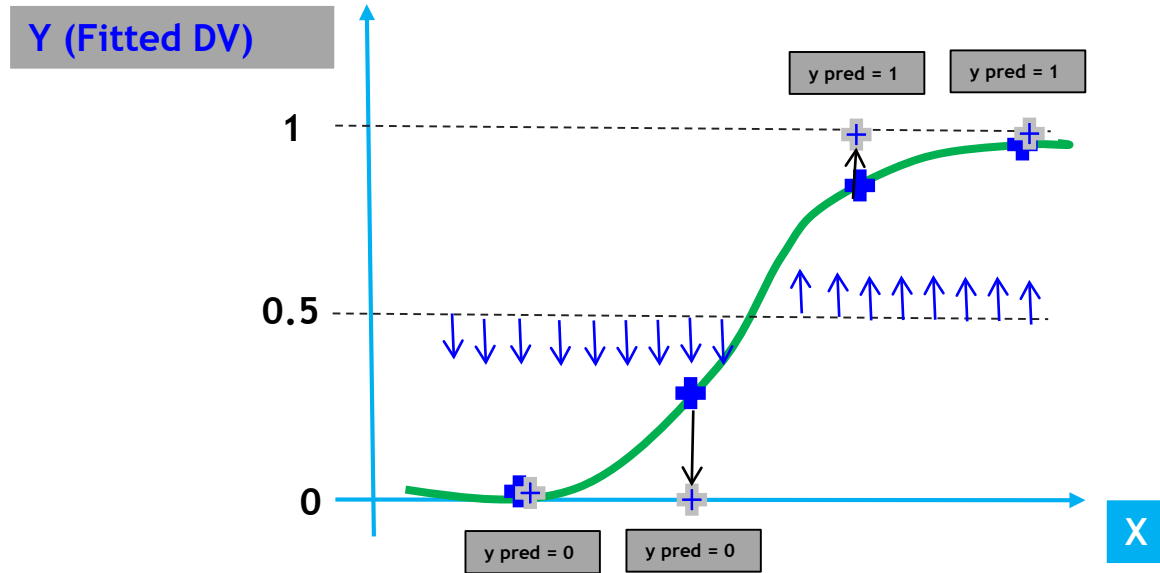
# Logistic Regression: anything above this threshold are predicted YES

# Logistic Regression: anything above this threshold are predicted YES

# Let's do it in R!!