# Data Analysis 2: Foundations of Statistics

## Instructor: Arieda Muço, Fall 2017

# Types of Random Variables

- Continuous, takes values in any interval
    - i.e prices, temperature, grades…
- Discrete, countable number of values
    - i.e hotel stars, gender, number of rooms…
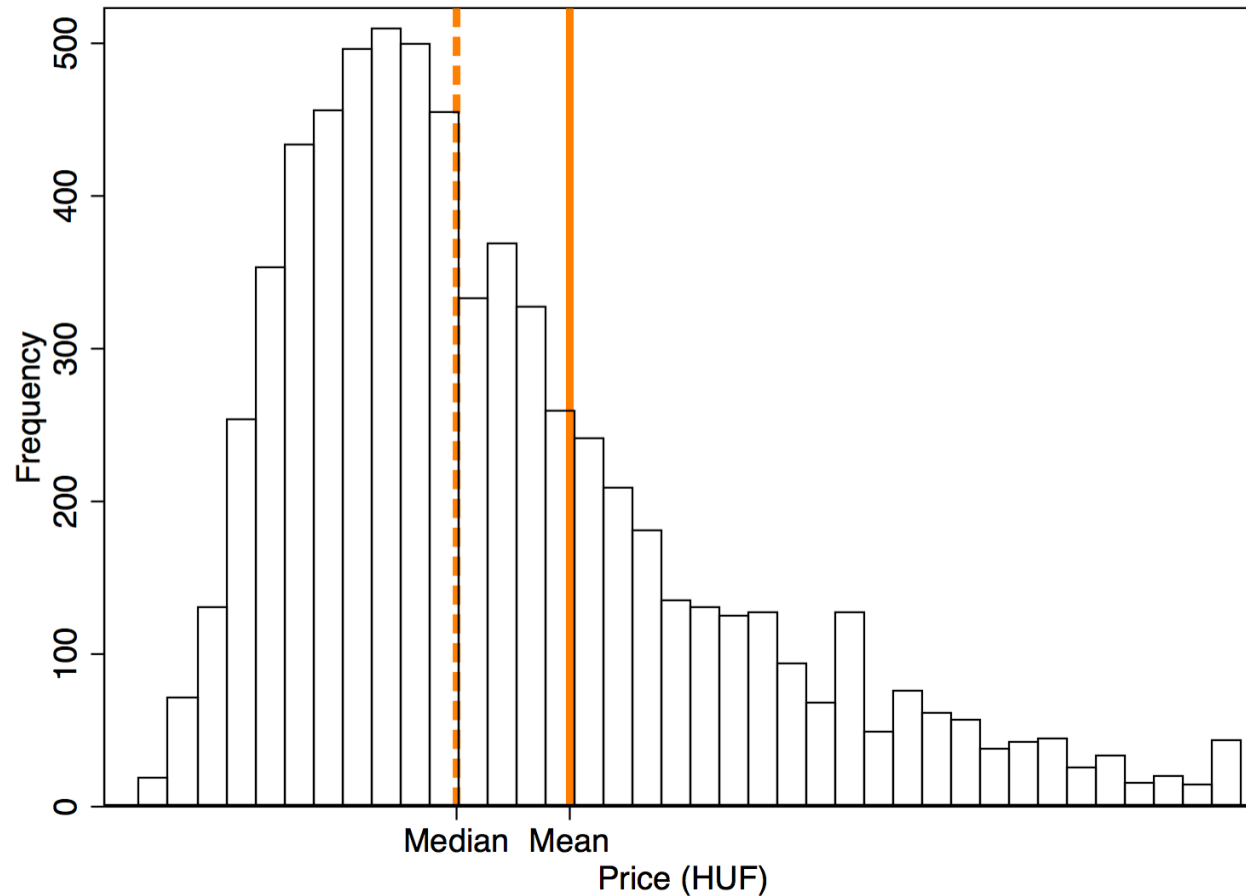
Each have their own distributions

# Distributions in Practice

- The distribution of a variable tells the number of times each possible value of the variable occurs in the data
  - Can be expressed as frequency, percentage

- It does so in isolation from other variables.
  - It does not tell if certain values are more likely to occur when some other variable, or variables, take certain values
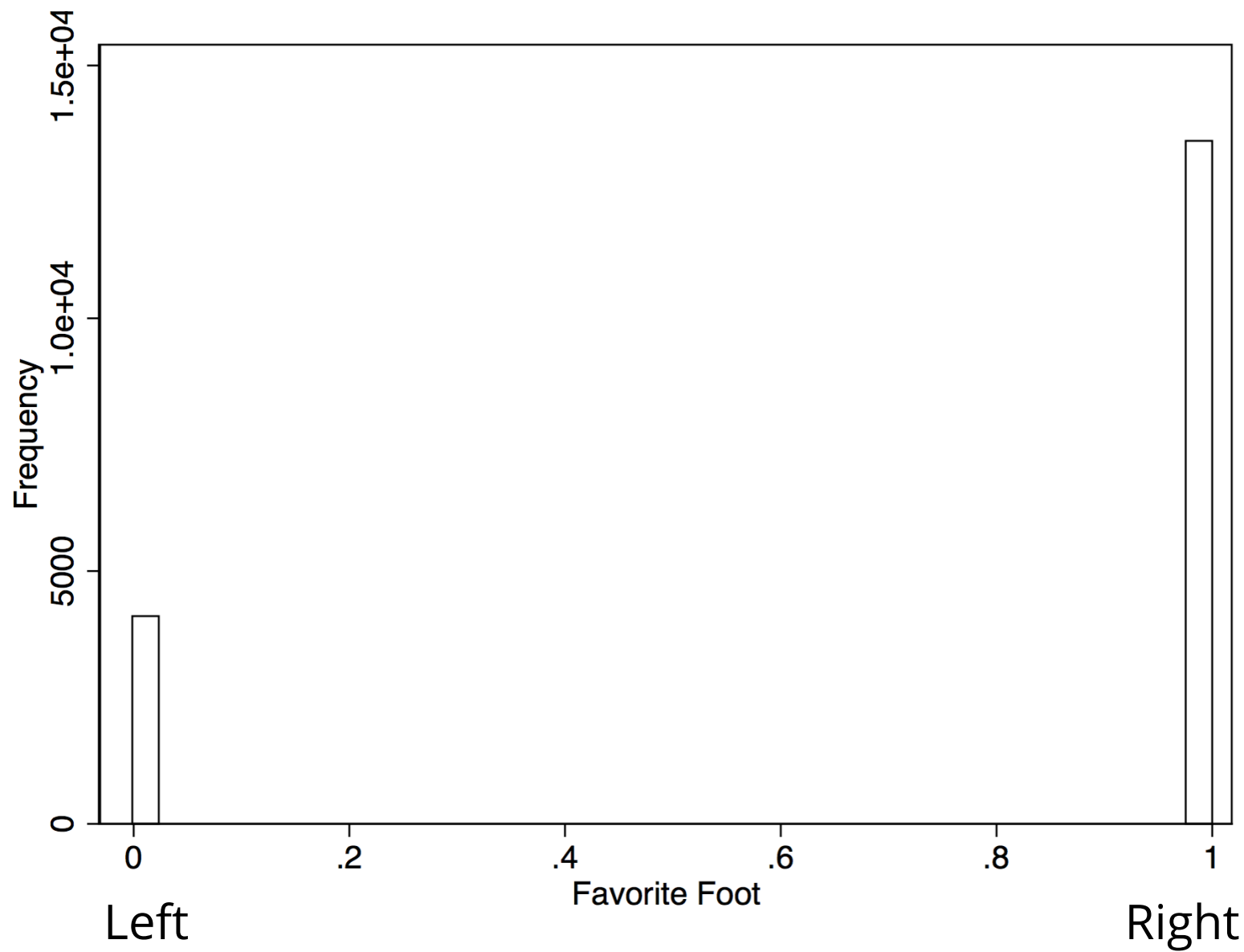
- The simplest way to visualize a distribution is through an histogram

- The histogram may takes on as many bars as the number of possible values

- From visual inspection of histograms we find out many interesting properties
    - the peaks, the immediate neighbourhood, if they have tails etc

Let's go back to our hotel price example.
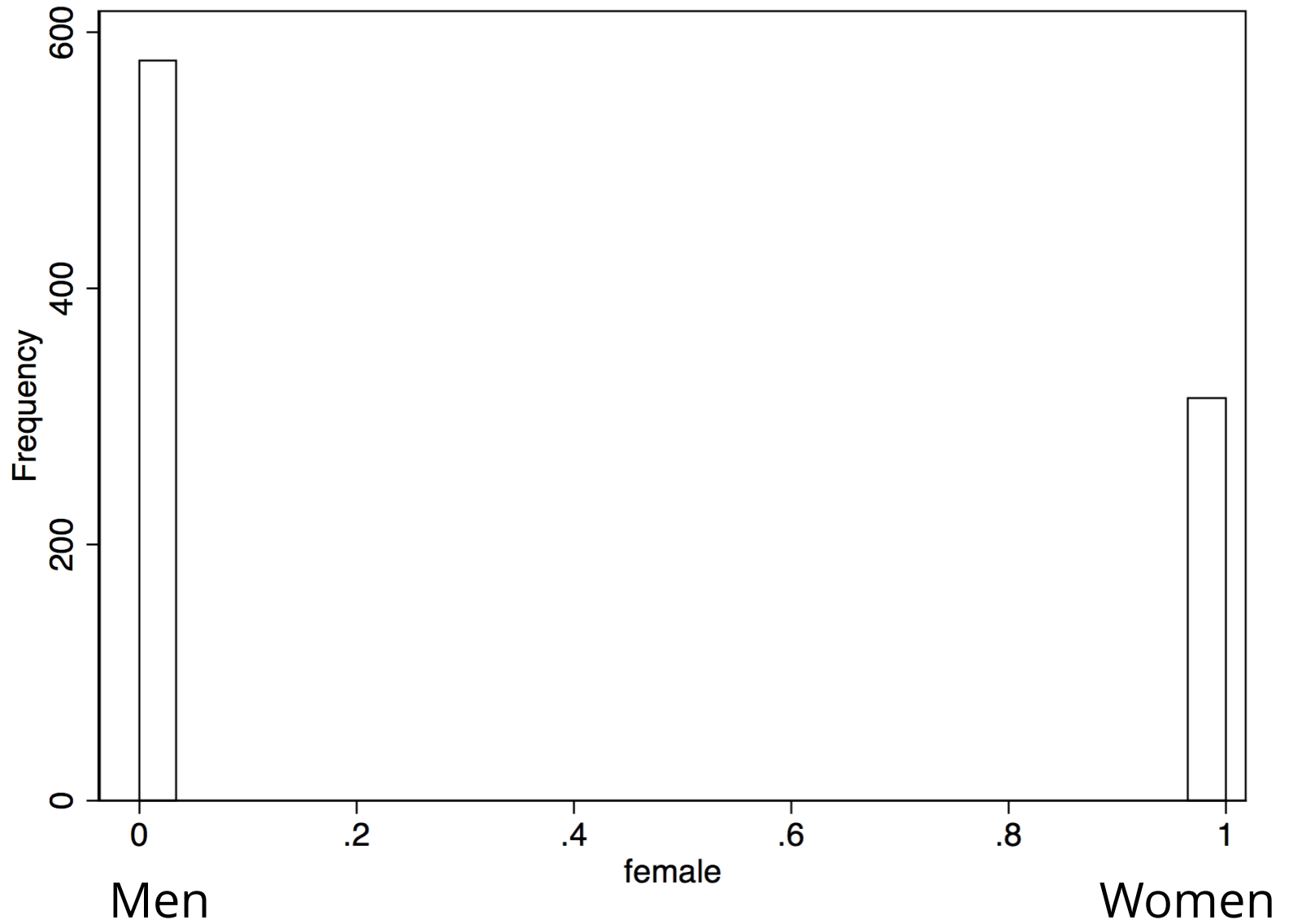What do we see?

# Histograms in practice

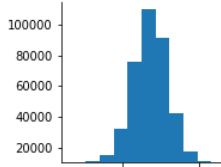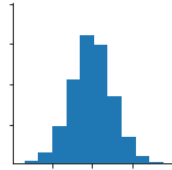# If I tell you that football players have a favorite foot?
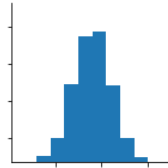
# Titanic passengers data

Men

Women

- Some of the properties of the distributions change if we change their bin size
  - Very wide bins may lump together multiple modes.
  - The statistical softwares you use, will compute the histogram with the default bin size
    - Try to figure out on your own the bin size that Stata or R use for their histograms.
    - Play with the simulations of distributions I provided you with by changing mean, variance, number of observarions, and binsize
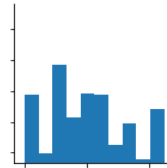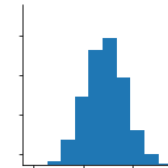    - Is this also the case for discrete variables?
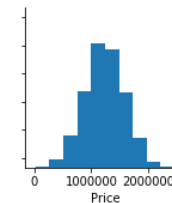
Average Area Income

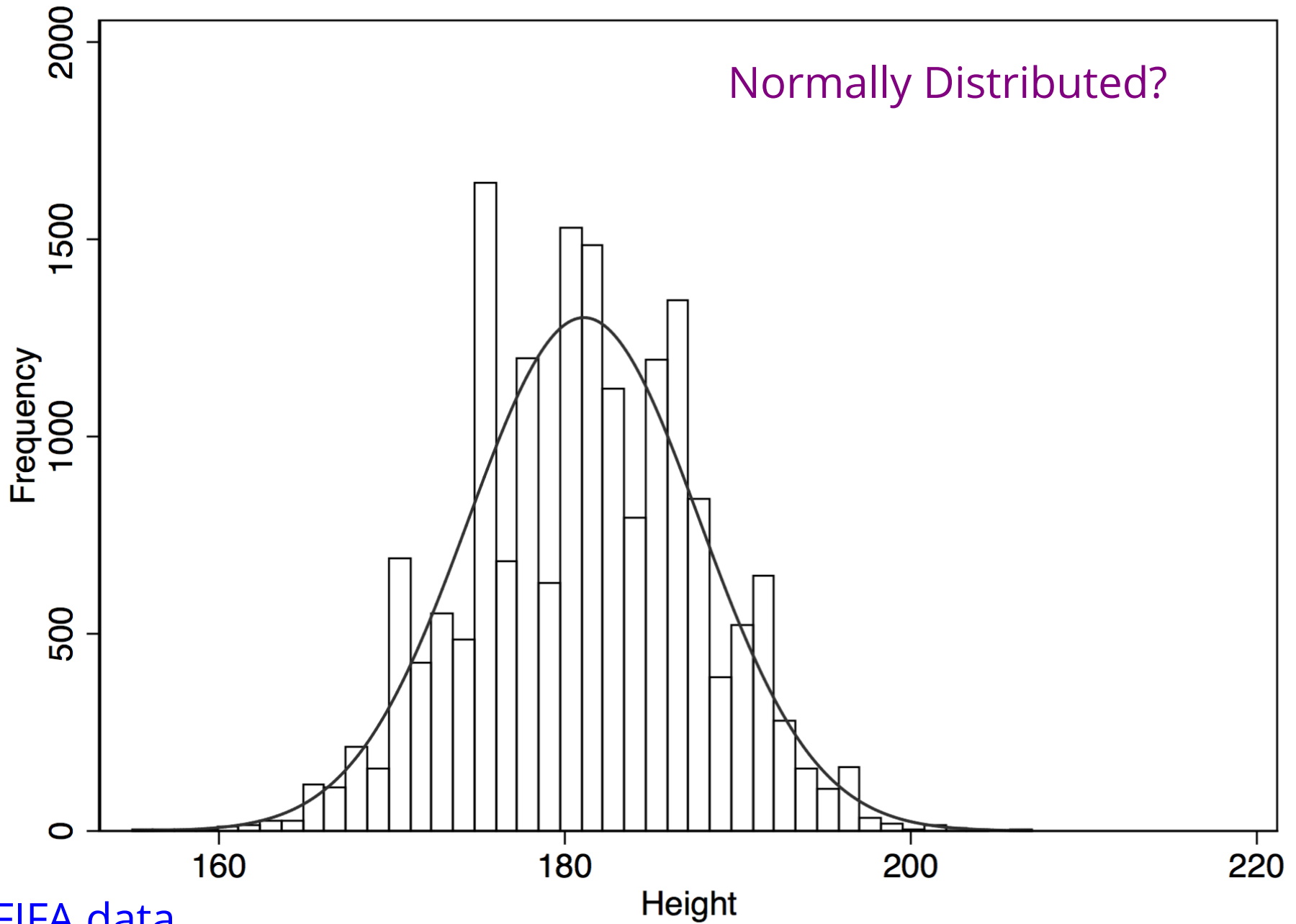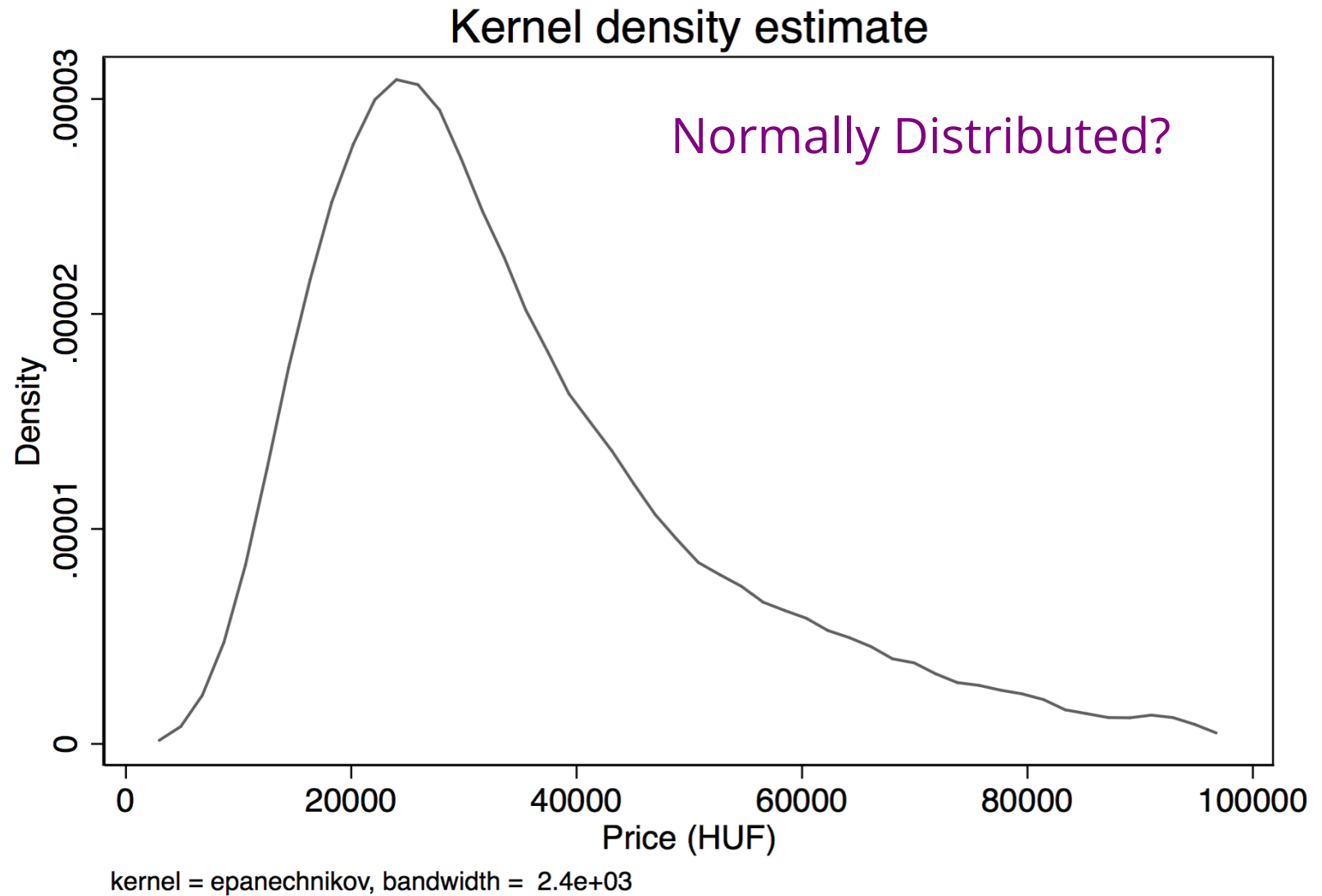Average Area House age

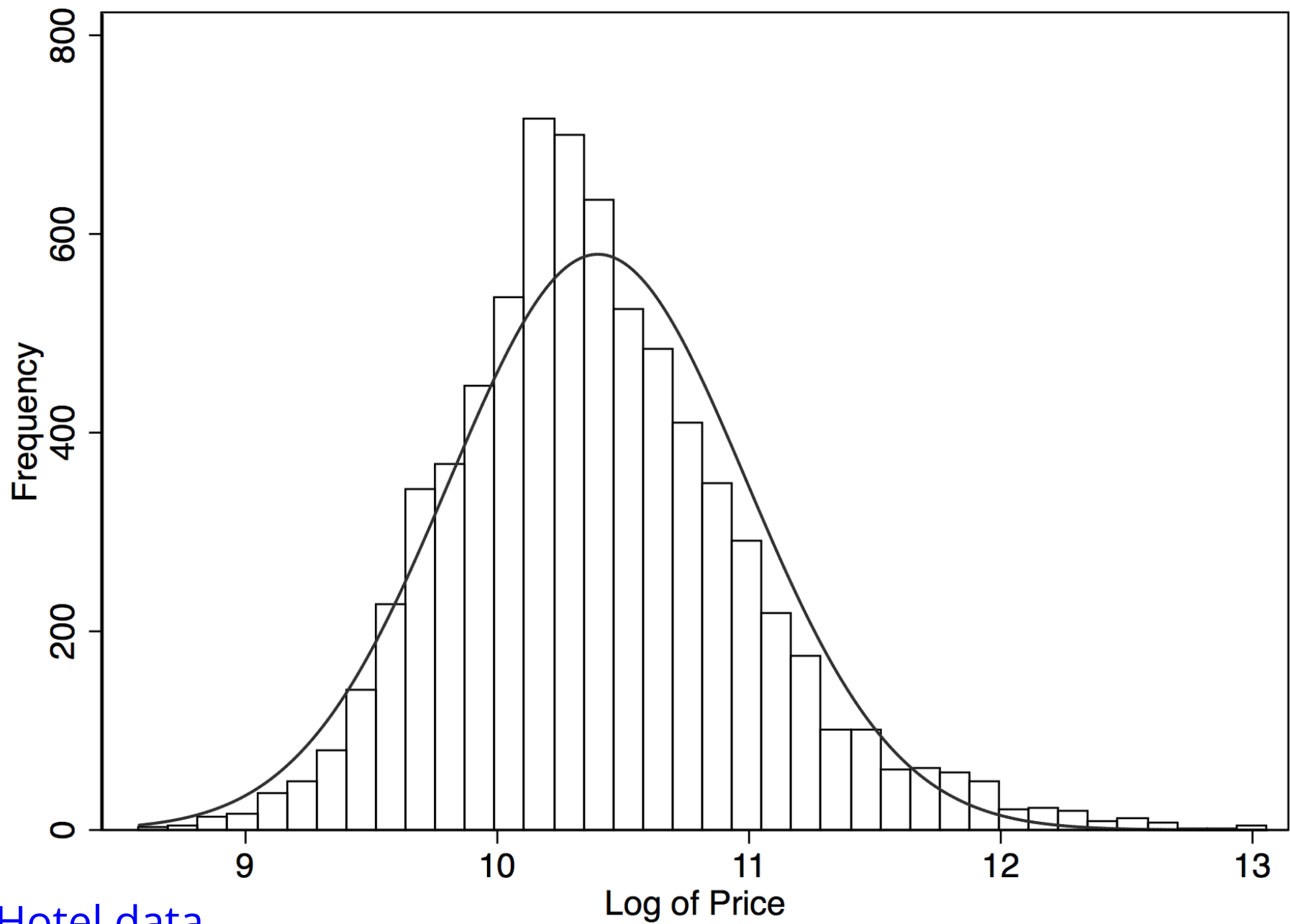Number of rooms

Number of bedrooms

Area Population

Price

- Kernel densities are an alternative way to histograms
  - for variables with many potential values

- A way to think about them is like curves that wrap around the corresponding hisogram

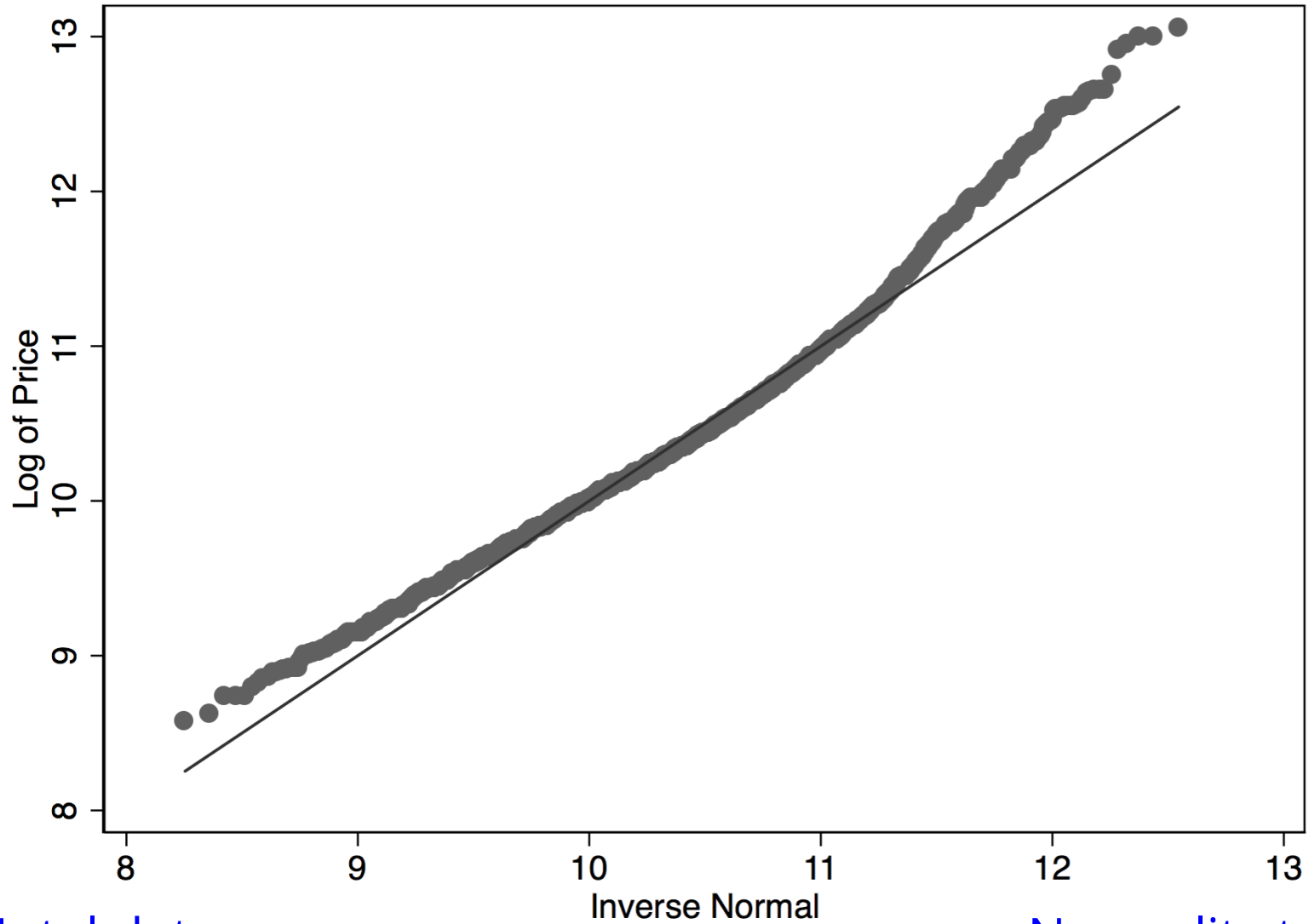- The most important parameter to set is the bandwidth, which is similar to setting the binsize in the histogram

Normally Distributed?

FIFA data

14

# Kernel density estimate



Normally Distributed?

kernel = epanechnikov, bandwidth = 2.4e+03

Hotel data

Hotel data

16

Hotel data                                    Normality test

# Joint Distributions

- In real life, we are often interested how variables that are related to each other. For example, number of rooms or bedrooms and average price of the house
  - The joint distribution shows the probabilities of each value combination of these variables we are interested in

# Conditional Distributions

- Conditional distributions are distributions of one variable for one (or more) values of the other variable

# Covariance

- It provides an indication of the dependence between two variables

$$Cov(x, y) = \frac{\sum_i (x_i - \overline{x})(y_i - \overline{y})}{n}$$

# Covariance

- In case: $\quad y_i = a + bx_i$

$$Cov(x, y) = \frac{b(\sum_i (x_i - \overline{x})(y_i - \overline{y}))}{n}$$

# Correlation

- Is computed by dividing the covariance by the standard deviation of each variable

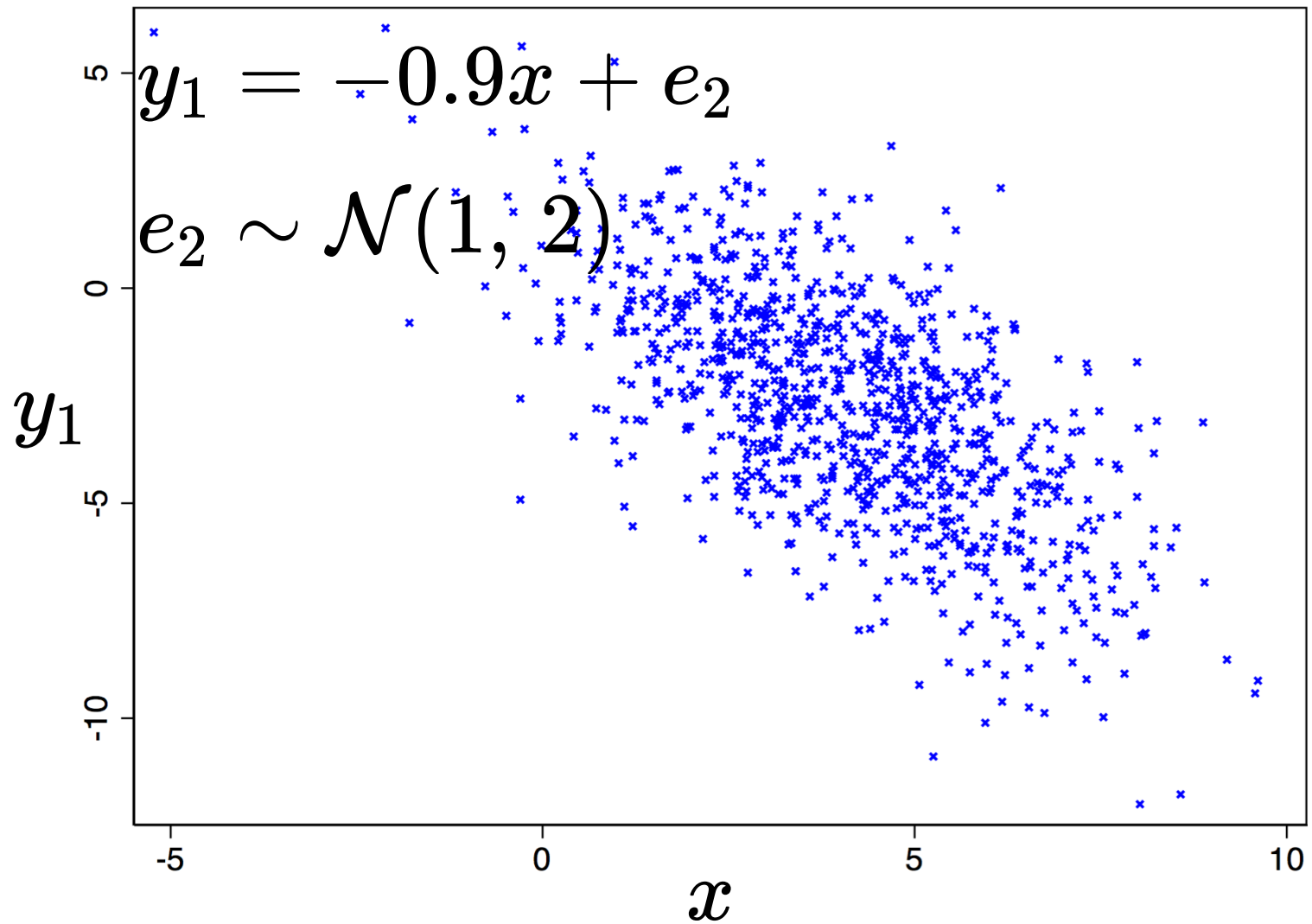$$\rho = Corr(x, y) = \frac{Cov(x,y)}{Std(x) \cdot Std(y)}$$
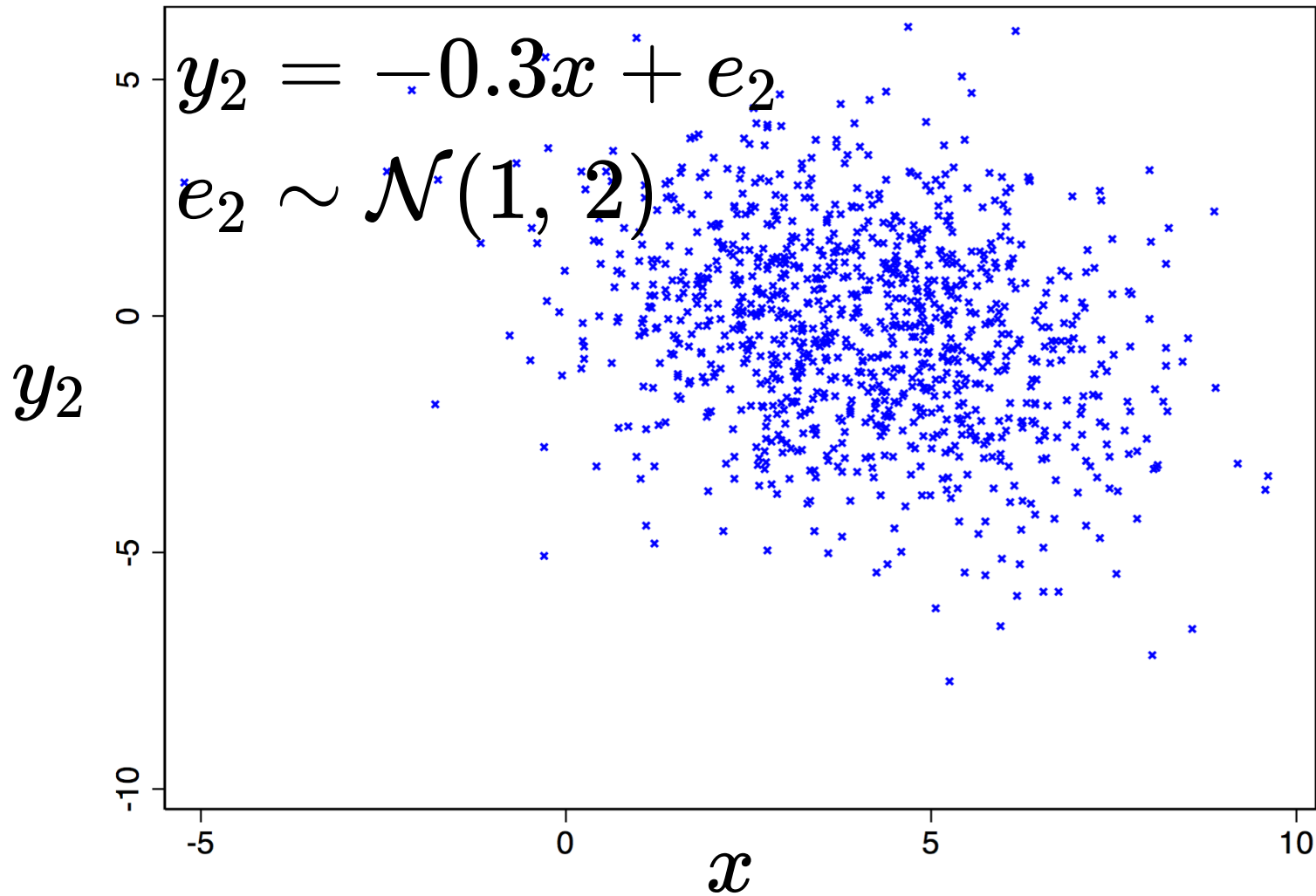
$$\rho \in [-1, 1]$$

# Scatterplot

- Allows us to tell if there is a relationship among pairs of variables under consideration

  - We investigate if there is a linear relationship, nonlinear, or no relationship
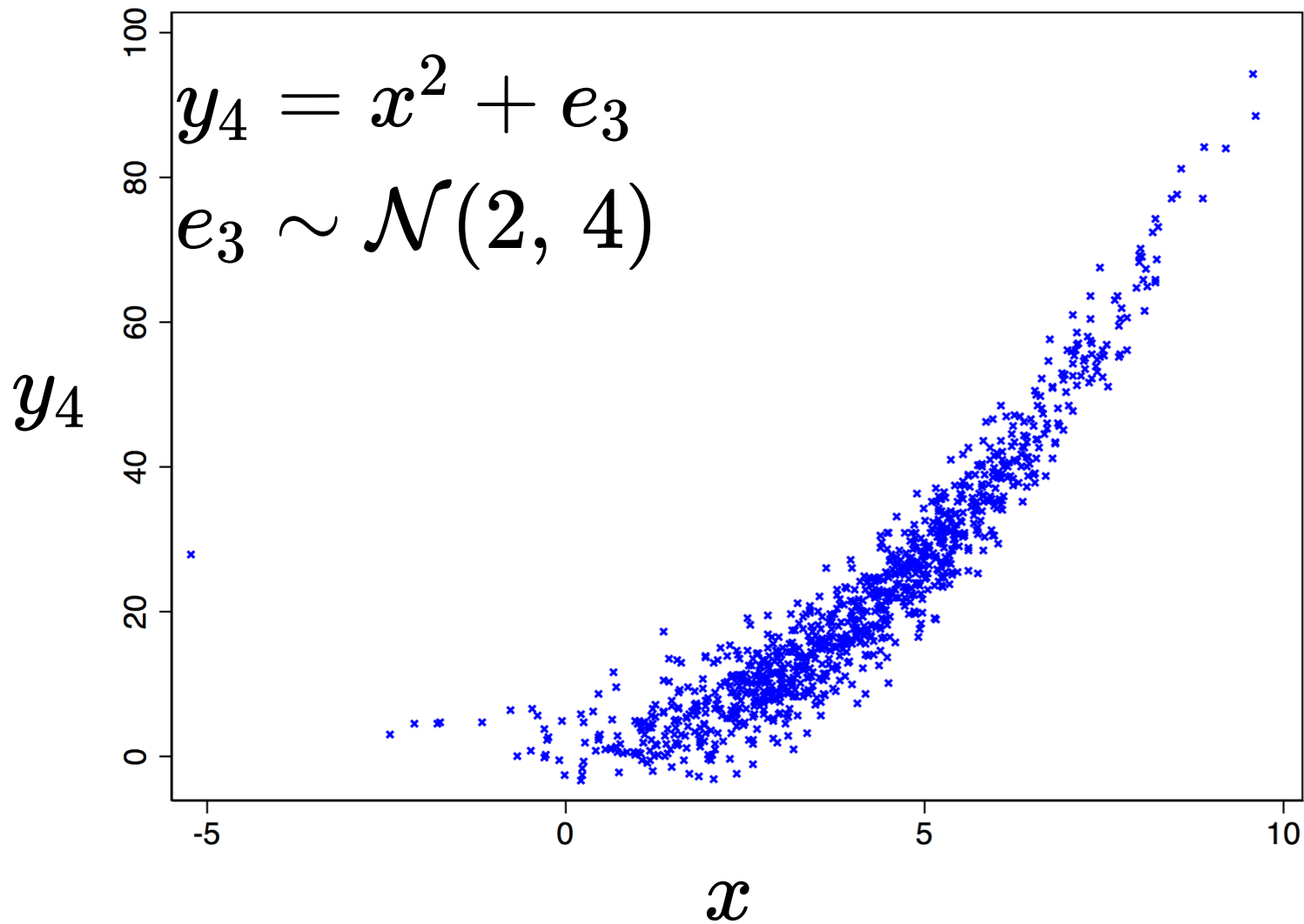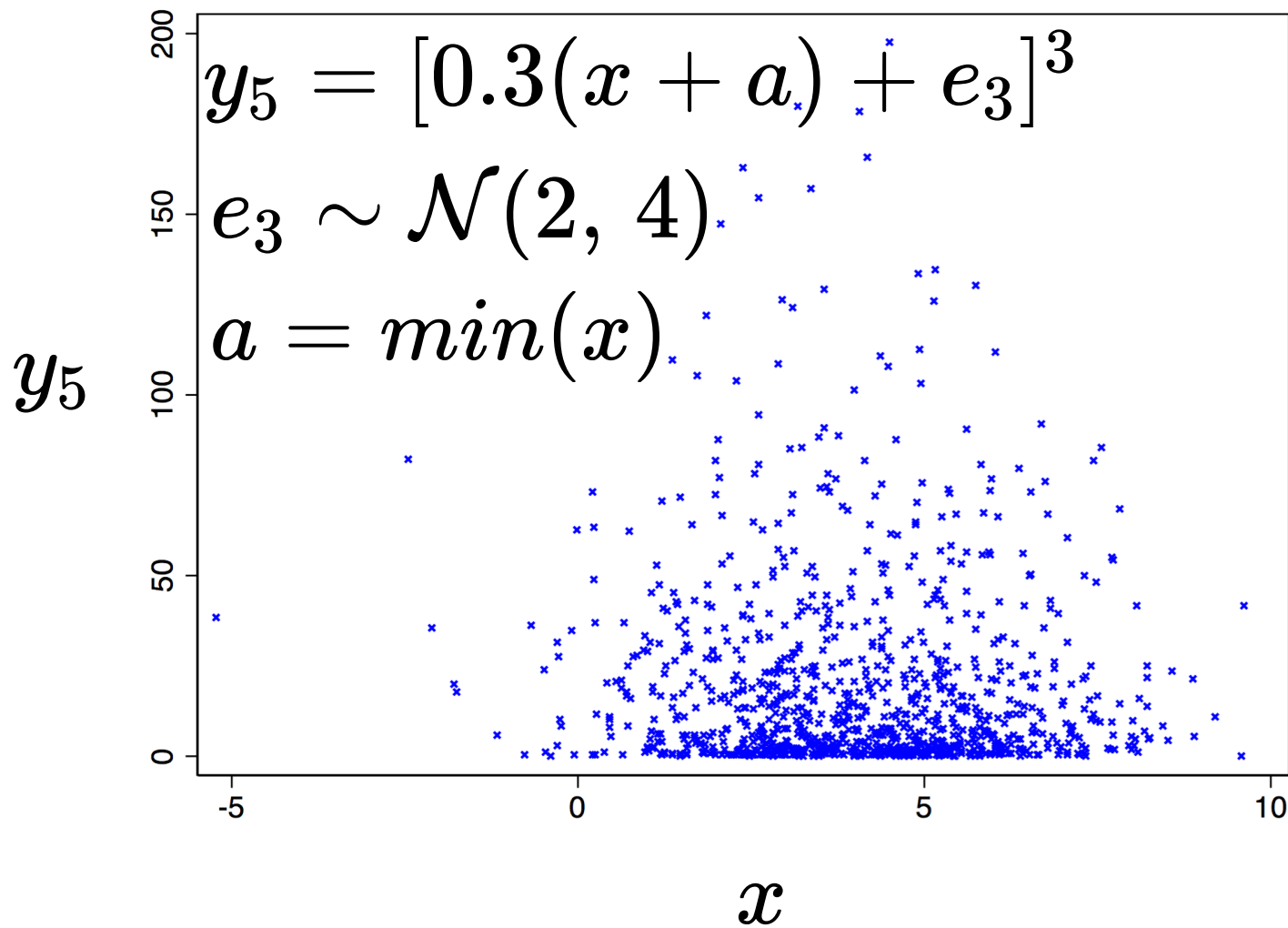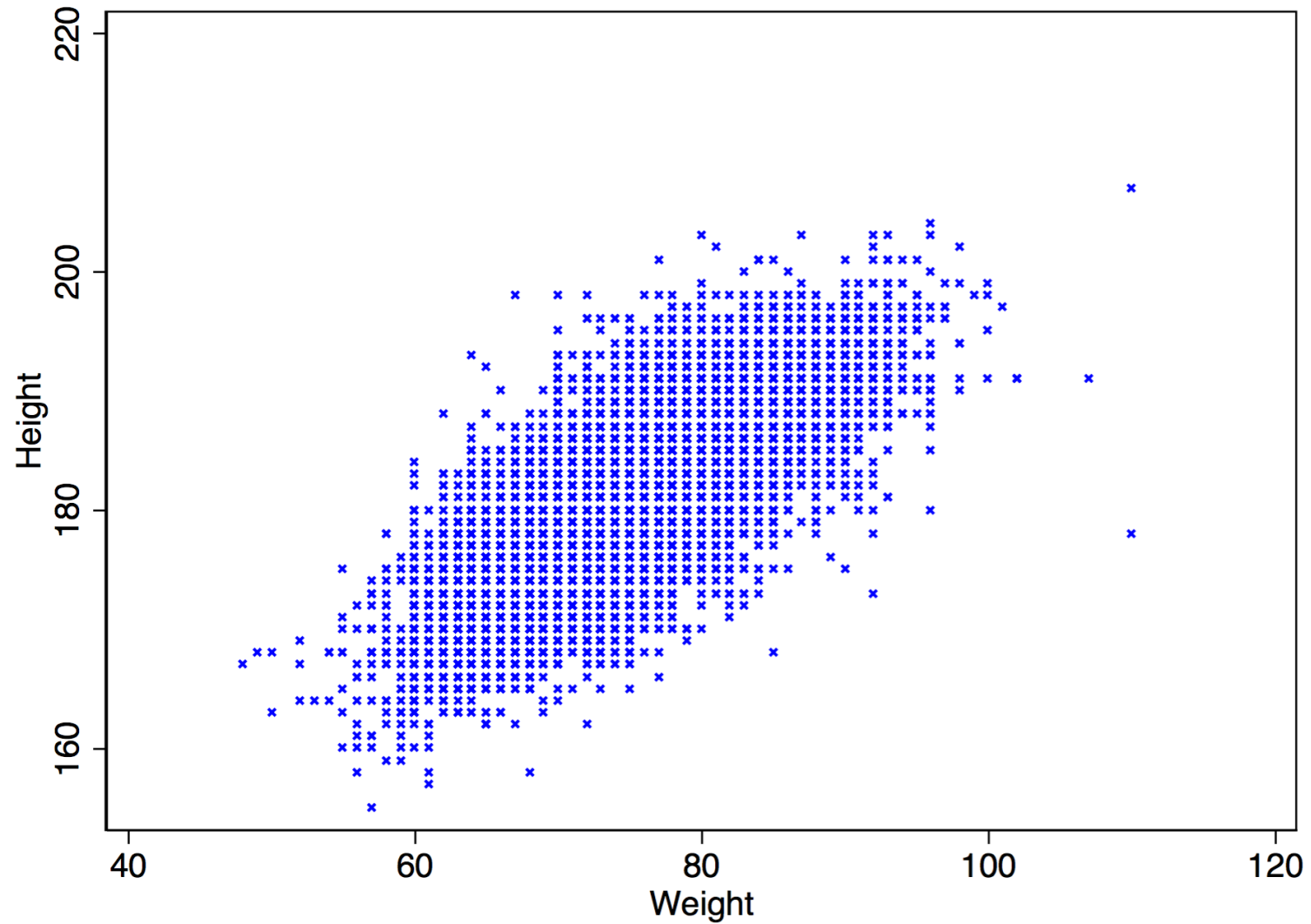
# Simulated data

CEU CENTRAL EUROPEAN UNIVERSITY

$$y = 0.9x + e_1$$

$$e_1 \sim \mathcal{N}(0, 1)$$

$$y_1 = -0.9x + e_2$$

$$e_2 \sim \mathcal{N}(1, 2)$$

$$y_2 = -0.3x + e_2$$
$$e_2 \sim \mathcal{N}(1, 2)$$

$$y_3 = -0.02x + e_2$$

$$e_2 \sim \mathcal{N}(1, 2)$$

$$y_4 = x^2 + e_3$$

$$e_3 \sim \mathcal{N}(2, 4)$$

$$y_5 = [0.3(x + a) + e_3]^3$$

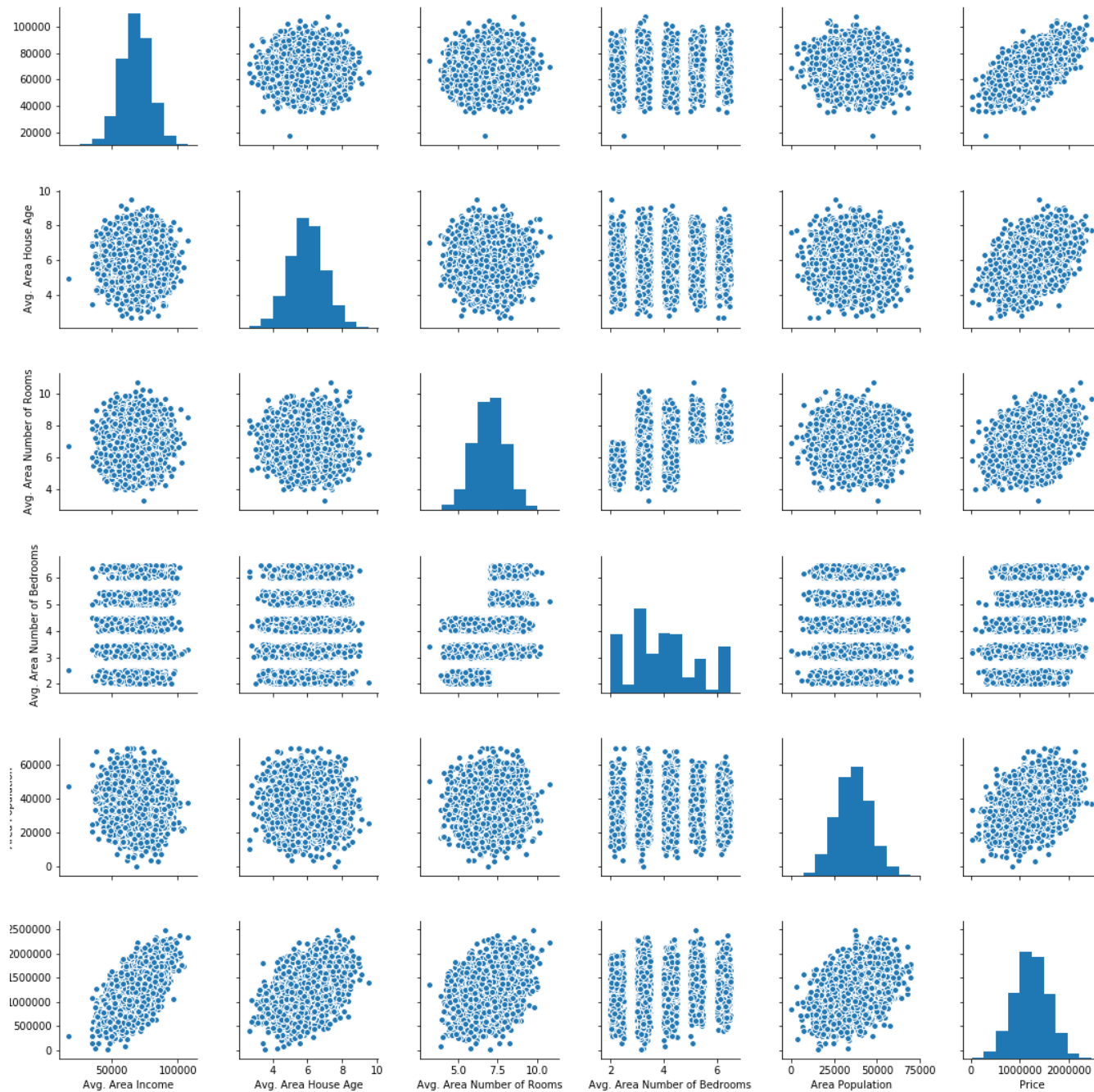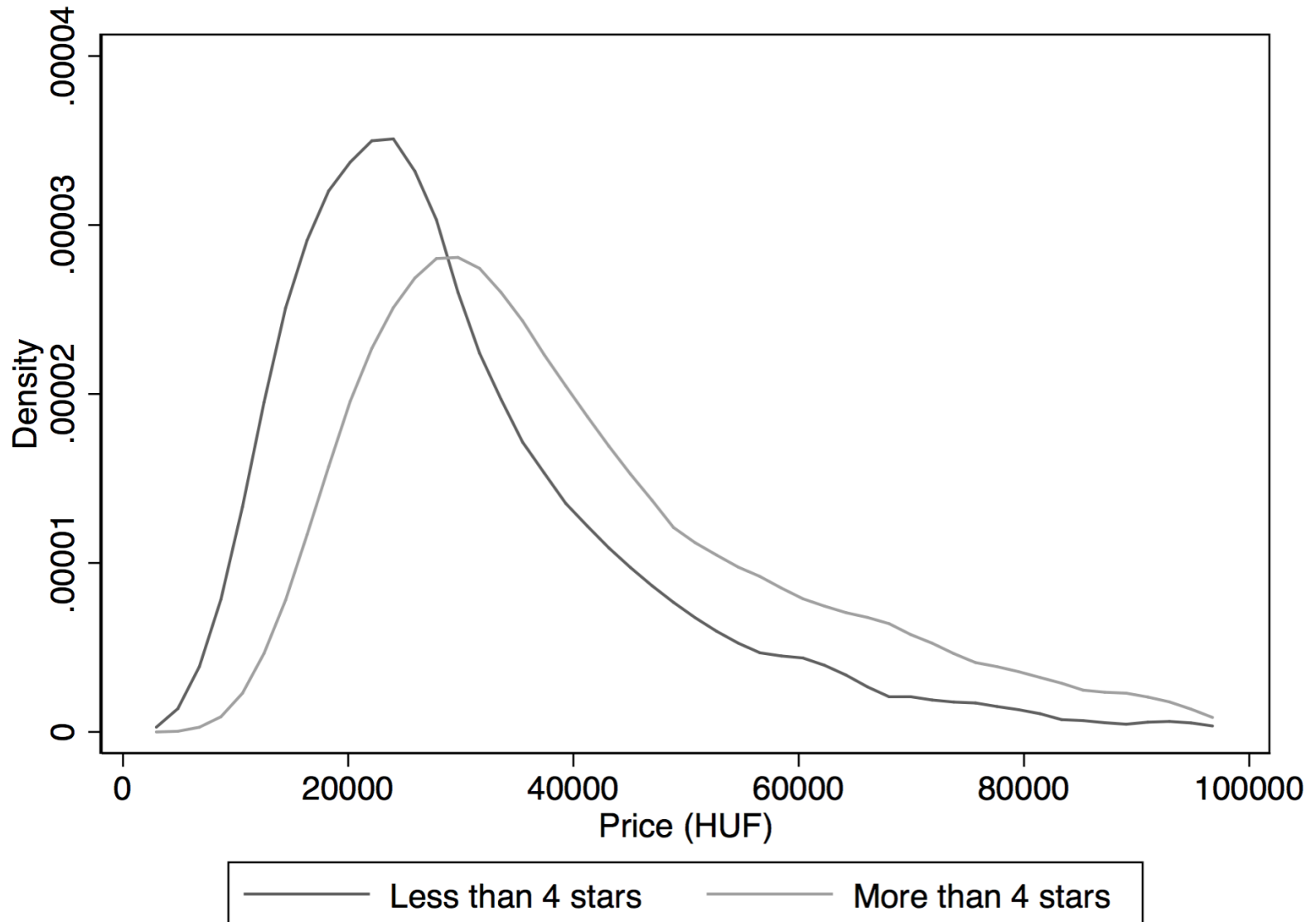$$e_3 \sim \mathcal{N}(2, 4)$$

$$a = min(x)$$

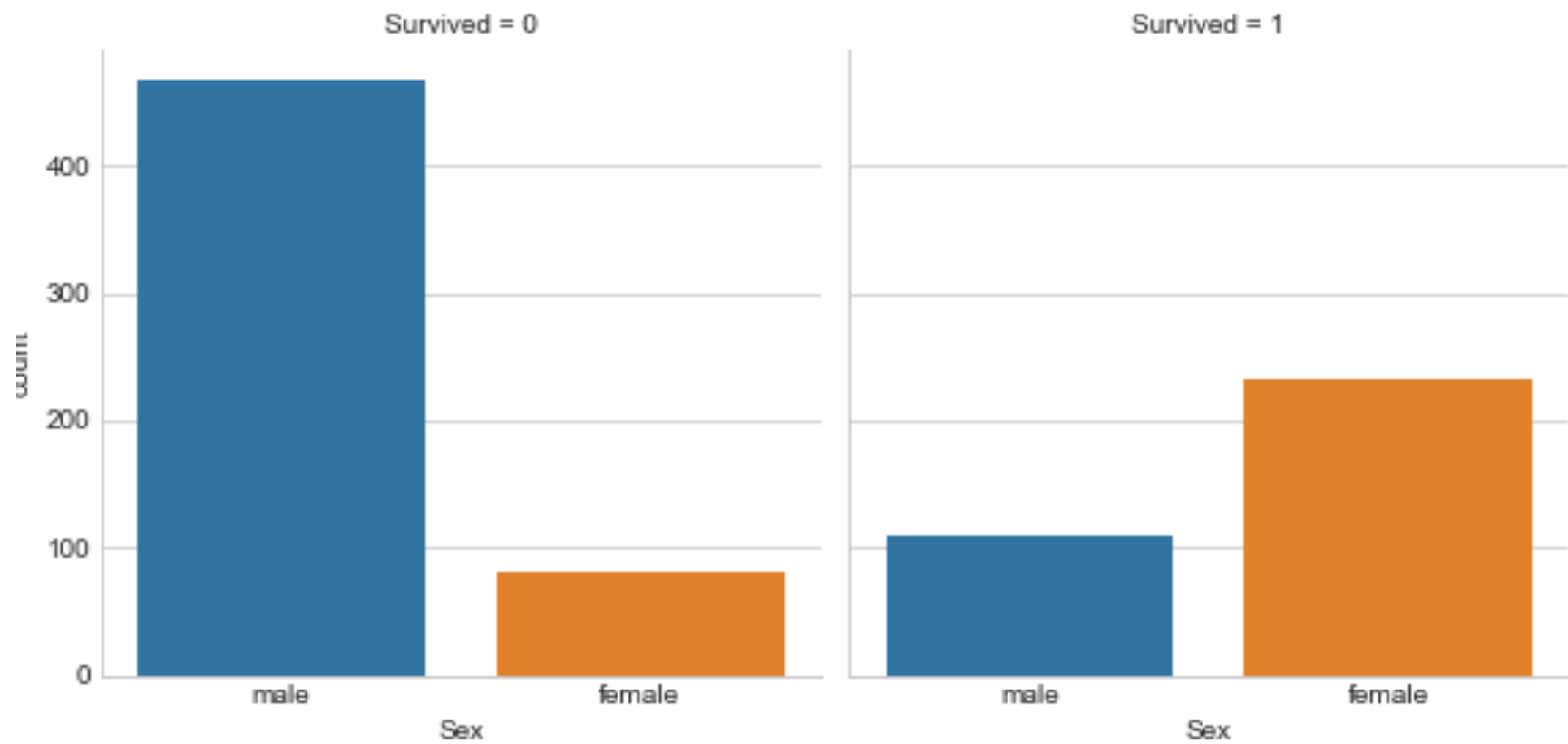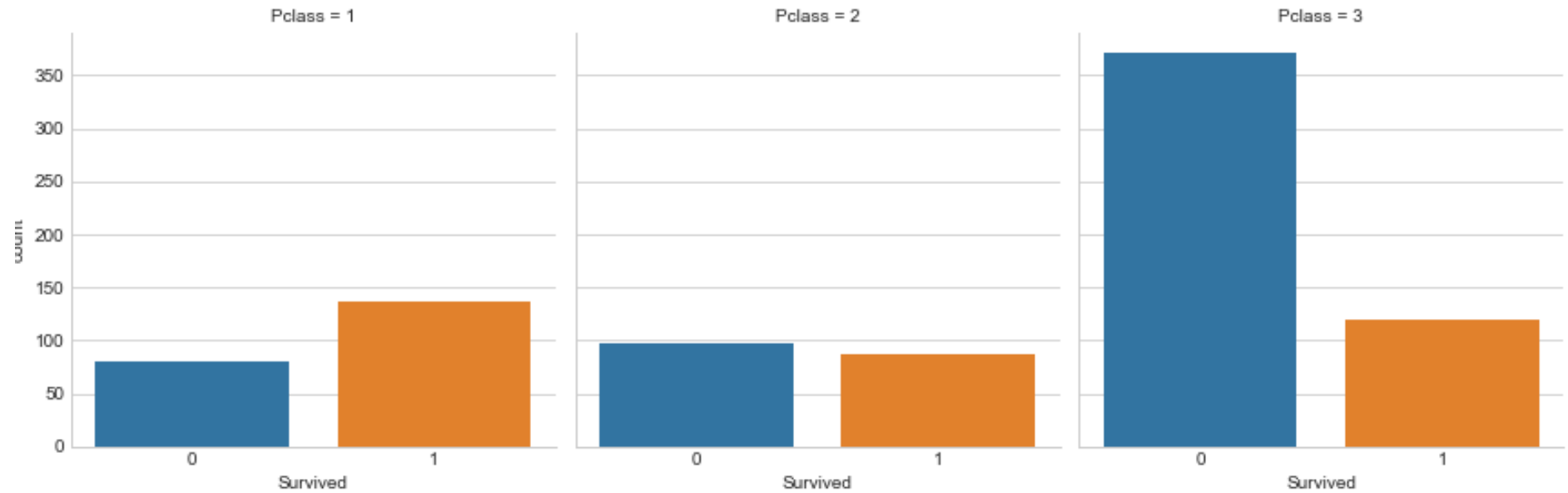# FIFA data

CEU CENTRAL EUROPEAN UNIVERSITY
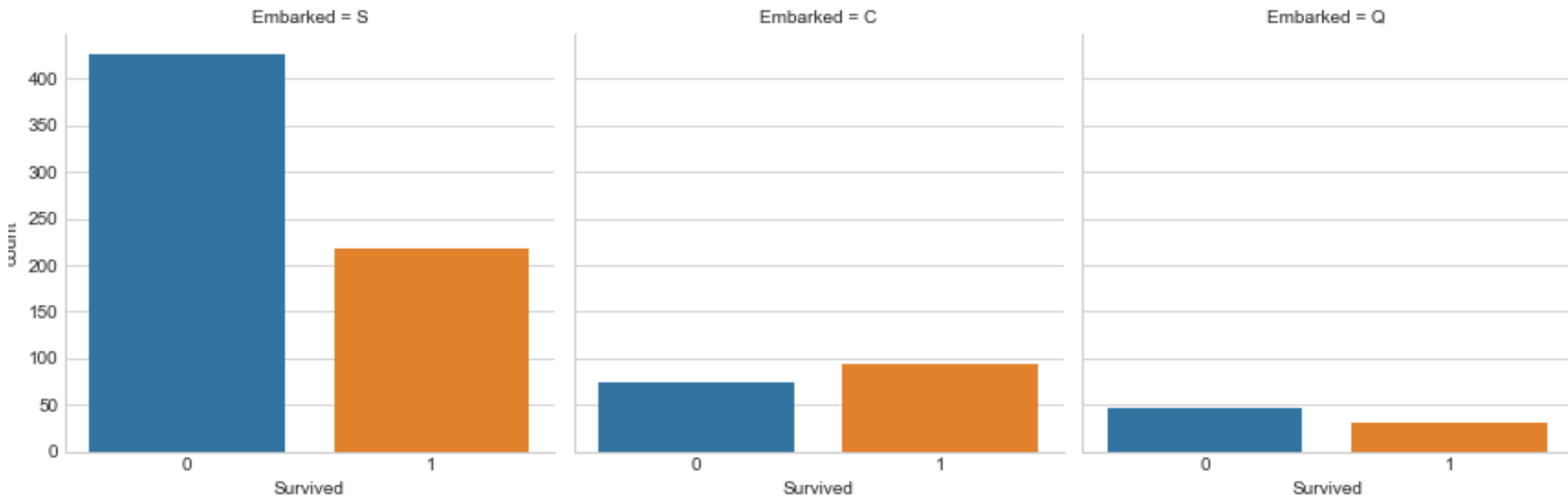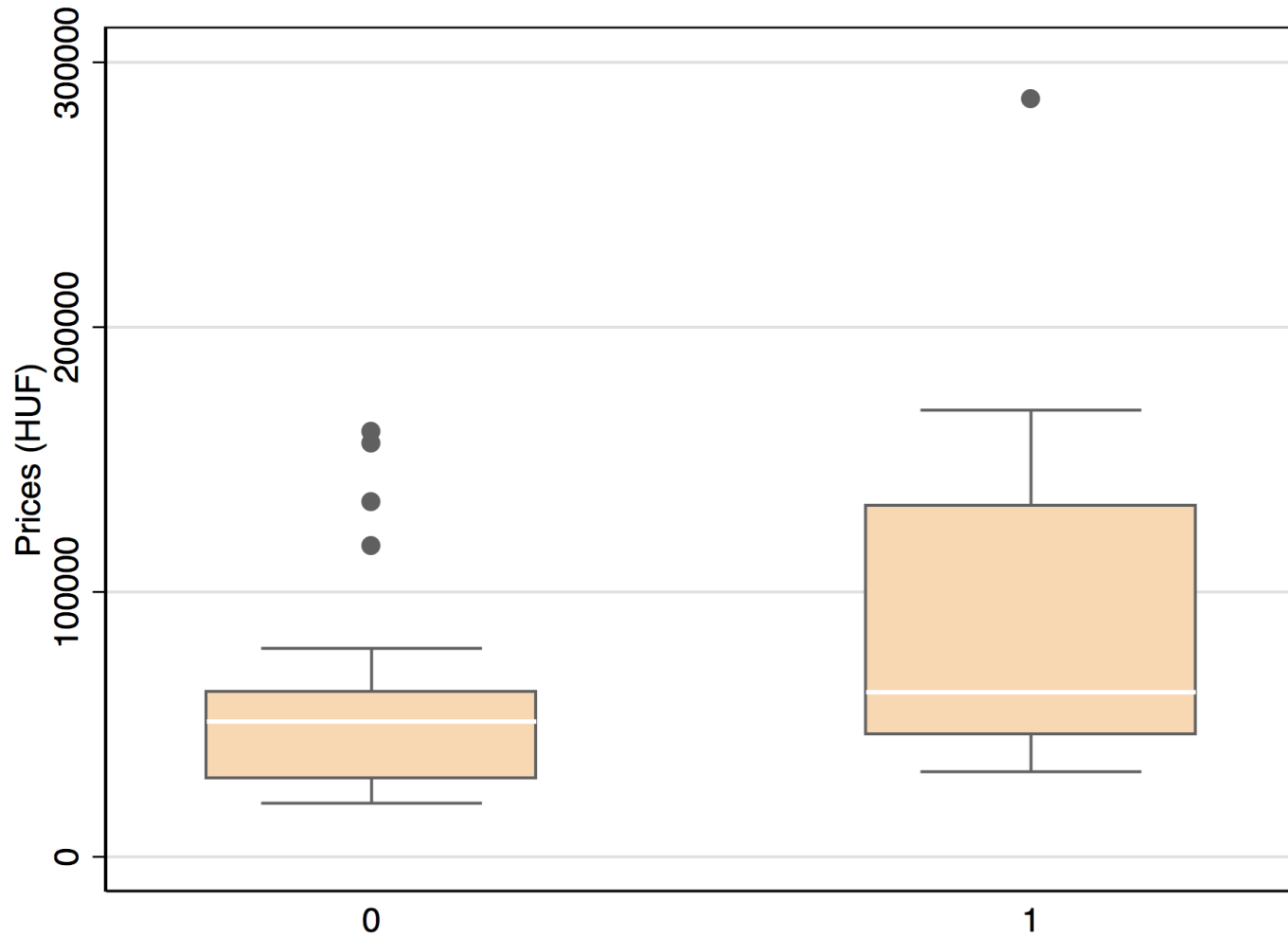
# Conditional Distributions in Practice

pclas=passanger class

Embarked=port of embarkation

# Vienna Hotel Prices (again)

Sources of data used in these slides

- Hotel data, already available (only to be shared in Moodle)

- Simulated variables will be posted (Moodle and GitHub), try it on yourself, now you have the codes

- Fifa database and Titanic passengers database are going to also be posted (Moodle and GitHub).
  - Downloaded from Kaggle

CEU | CENTRAL
EUROPEAN
UNIVERSITY

45