# Exploratory Data Analysis
## Part B - Visualizing the data : The drill
Pavlos Protopapas

# Visualization

Visualization is incredibly important, both for EDA and for communicating our results to others.



"The greatest value of a picture is when it forces us to notice what we never expected to see."

John Tukey

(American mathematical statistician, best known for the development of the Fast Fourier Transform algorithm and box plot.)

What's the need to visualize?

# Anscombe's Quartet

| Set A | | Set B | | Set C | | Set D | |
|---|---|---|---|---|---|---|---|
| 10 | 8.04 | 10 | 9.14 | 10 | 7.46 | 8 | 6.58 |
| 8 | 6.95 | 8 | 8.14 | 8 | 6.77 | 8 | 5.76 |
| 13 | 7.58 | 13 | 8.74 | 13 | 12.74 | 8 | 7.71 |
| 9 | 8.81 | 9 | 8.77 | 9 | 7.11 | 8 | 8.84 |
| 11 | 8.33 | 11 | 9.26 | 11 | 7.81 | 8 | 8.47 |
| 14 | 9.96 | 14 | 8.1 | 12 | 8.84 | 8 | 7.04 |
| 6 | 7.24 | 6 | 6.13 | 6 | 6.08 | 8 | 5.25 |
| 4 | 4.26 | 4 | 3.1 | 4 | 5.39 | 19 | 12.5 |
| 12 | 10.84 | 12 | 9.11 | 12 | 8.15 | 8 | 5.56 |
| 7 | 4.82 | 7 | 7.26 | 7 | 6.42 | 8 | 7.91 |
| 5 | 5.68 | 5 | 4.74 | 5 | 5.73 | 8 | 6.89 |

# Anscombe's Quartet

Anscombe's Quartet is a set of four datasets, where each produces the same summary statistics (mean, standard deviation, and correlation), which could lead one to believe the datasets are quite similar. However, after visualizing (plotting) the data, it becomes clear that the datasets are markedly different.

**Summary Statistics**

$\mu_X$ = 9.0     $\sigma_X$ = 3.317

$\mu_Y$ = 7.5     $\sigma_Y$ = 2.03

**Linear Regression**

Y2 = 3 + 0.5 X

R2 = 0.67

# Make sure the statistics are not fooling you!



Fig 1. Anscombe's Quartet (left), and a "Unstructured Quartet" on the right, where the datasets have the same summary statistics as those in Anscombe's Quartet, but lack underlying structure or visual distinction.
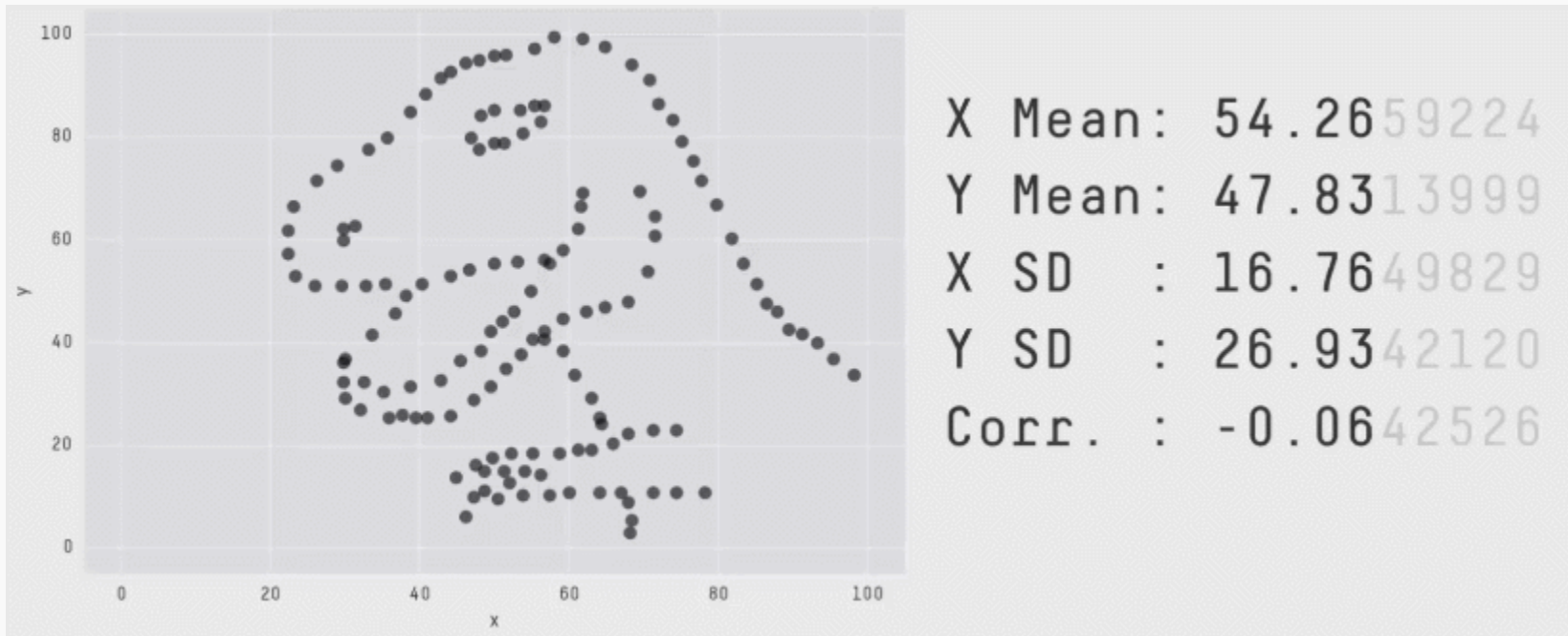
**Same stats** do not imply same graphs    **Same graphs** do not imply same stats

Credits: Same Stats, Different Graphs. (Autodesk Research)

# Same stats, different graphs!

# The Datasaurus Dozen

Datasets which are identical over a number of statistical properties, yet produce dissimilar graphs, are frequently used to illustrate the importance of graphical representations when exploring data.



X Mean: 54.2659224
Y Mean: 47.8313999
X SD  : 16.7649829
Y SD  : 26.9342120
Corr. : -0.0642526

THE DATASAURUS DOZEN

Credits: Same Stats, Different Graphs. (Autodesk Research)
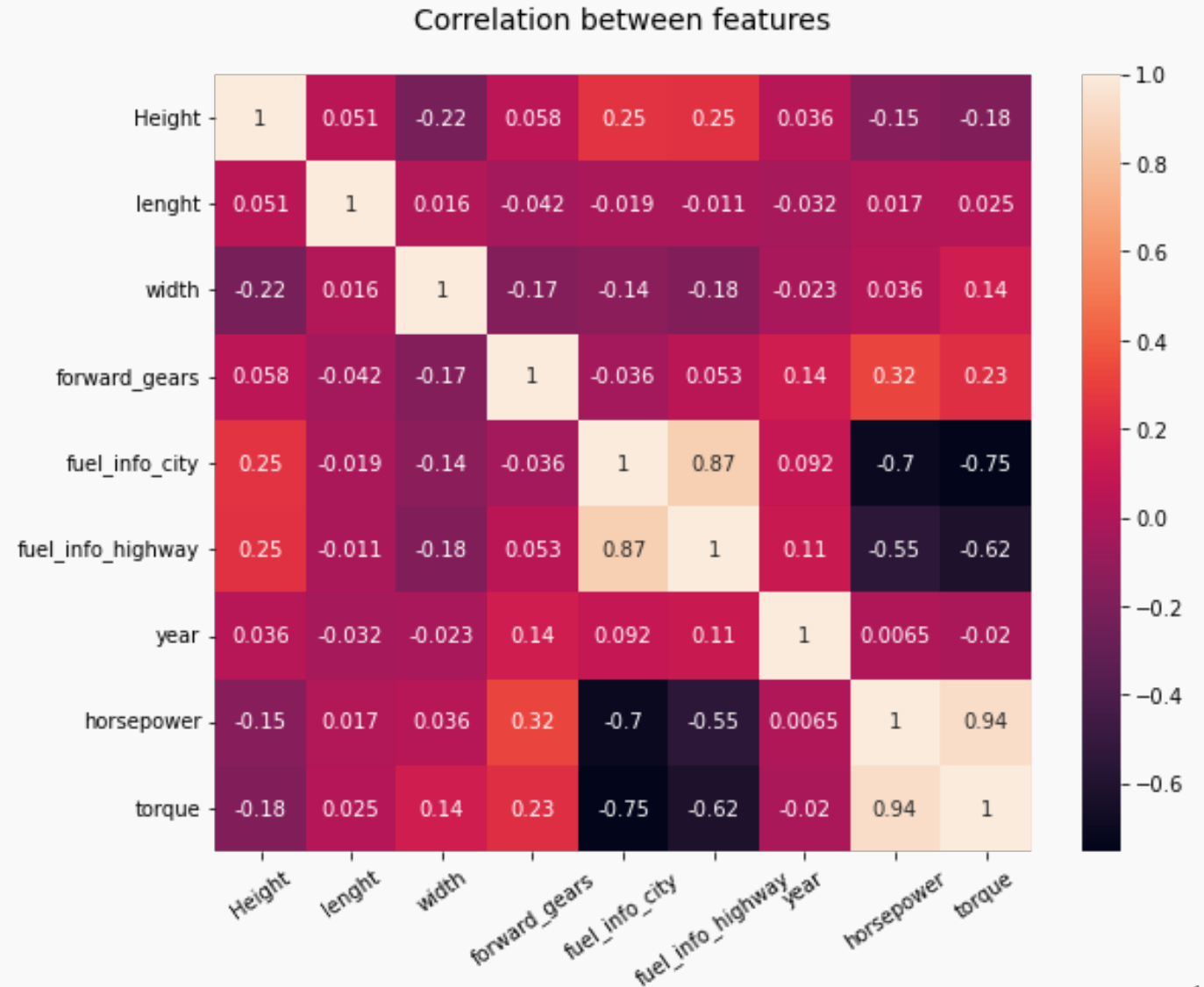
# Visualization Goals

**Analyze (Exploratory)**

1. Explore the data.

2. Assess the situation.

3. Determine how to proceed.

4. Decide what to do.

# Analyze (Exploratory)

## Exploring data

The figure illustrates the correlation plot of numerical variables using a heat map.

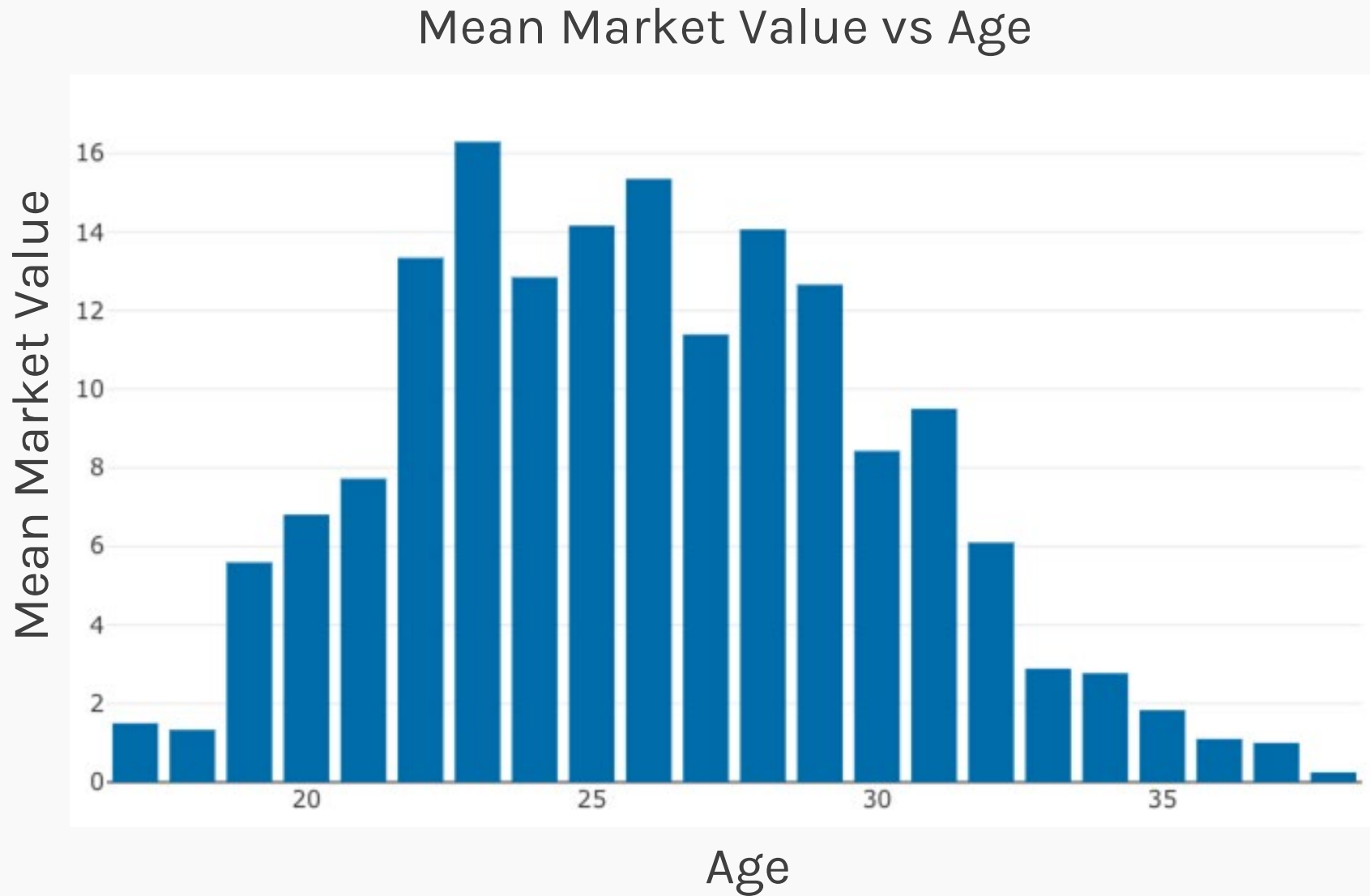The correlation plot is used to drop variables that are highly correlated..



Correlation between features

# Visualization Goals

**Analyze (Exploratory)**

1. Explore the data.

2. Assess the situation.

3. Determine how to proceed.

4. Decide what to do.

# Now, let's comeback to visualizing the data from English Premier League

# Visualization: The English Premier League

## Mean Market Value vs Age

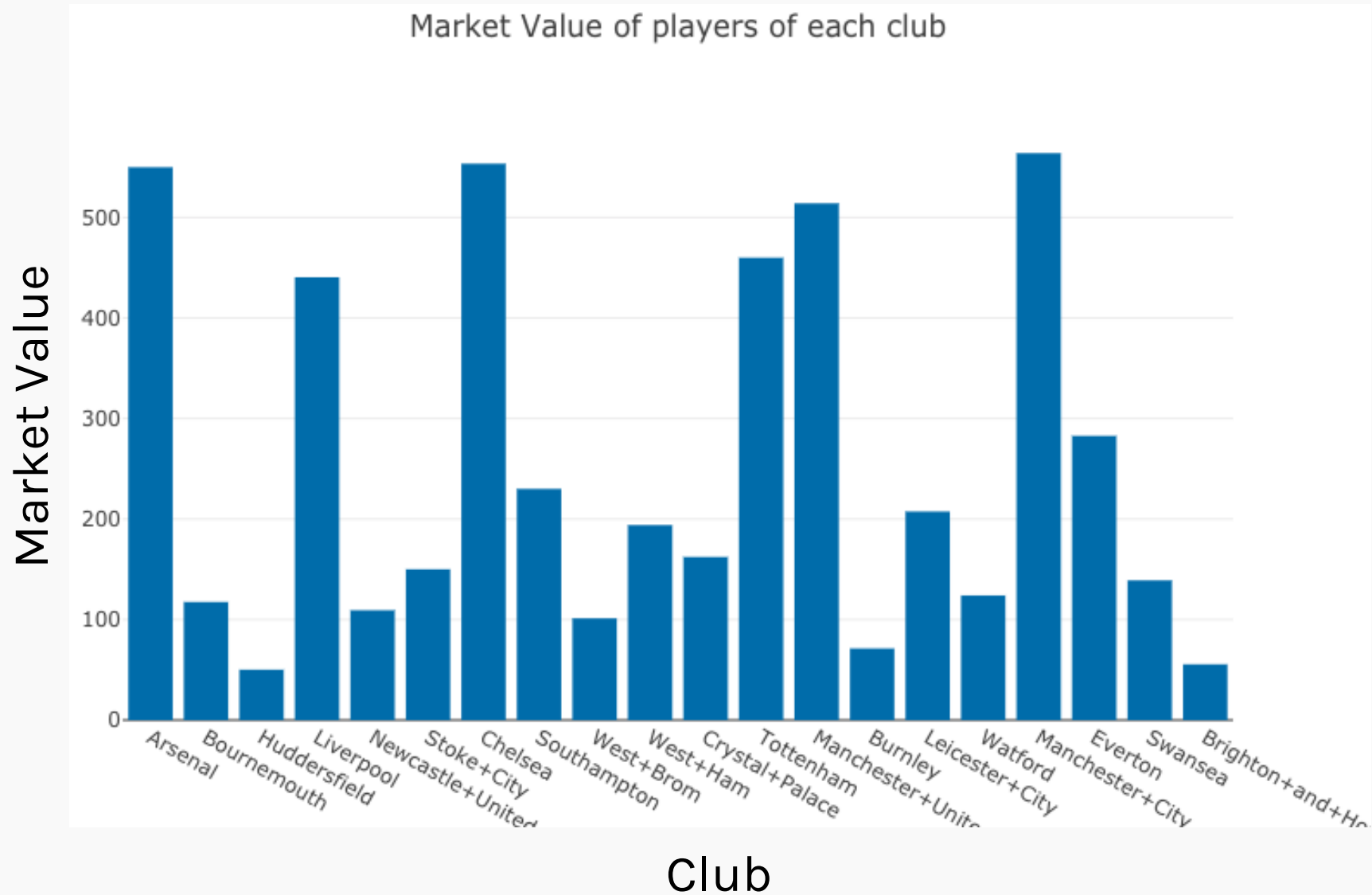# Visualization: The English Premier League

**Components of the graph**

Mean Market Value vs Age

Y label

Mean Market Value

Age

**Components of the graph**

Mean Market Value vs Age

X-label

Age

# Visualization: The English Premier League

**Components of the graph**



Mean Market Value vs Age

Keep in mind having adequate number of x-ticks & y-ticks and maintain good font density.

# Visualization: The English Premier League
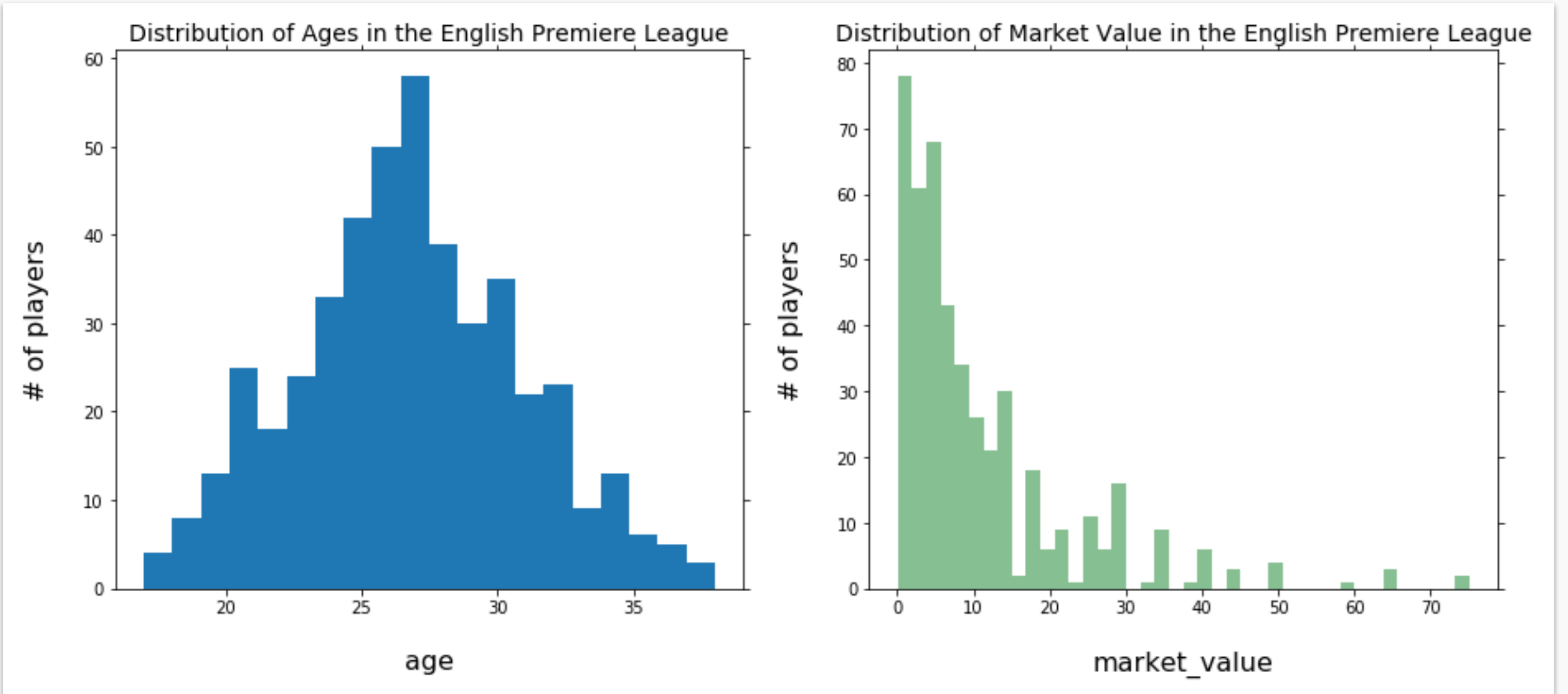


Market Value of players of each club

# Visualization: The English Premier League

# Visualization: The English Premier League



Distribution of Ages in the English Premiere League

Distribution of Market Value in the English Premiere League
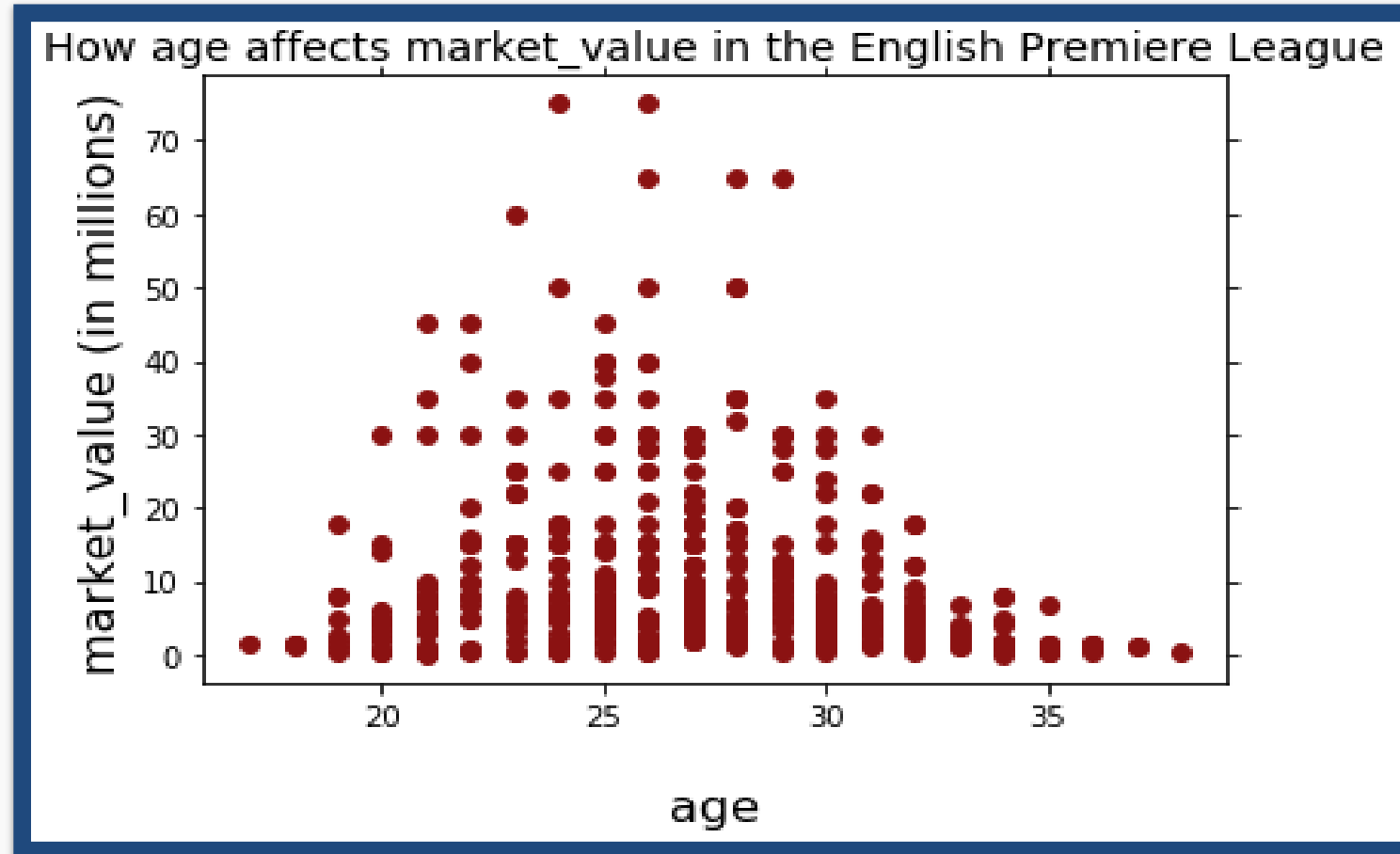
# Visualization: The English Premier League



Side-by-side plots are a lot helpful when we want to compare distributions of variables.
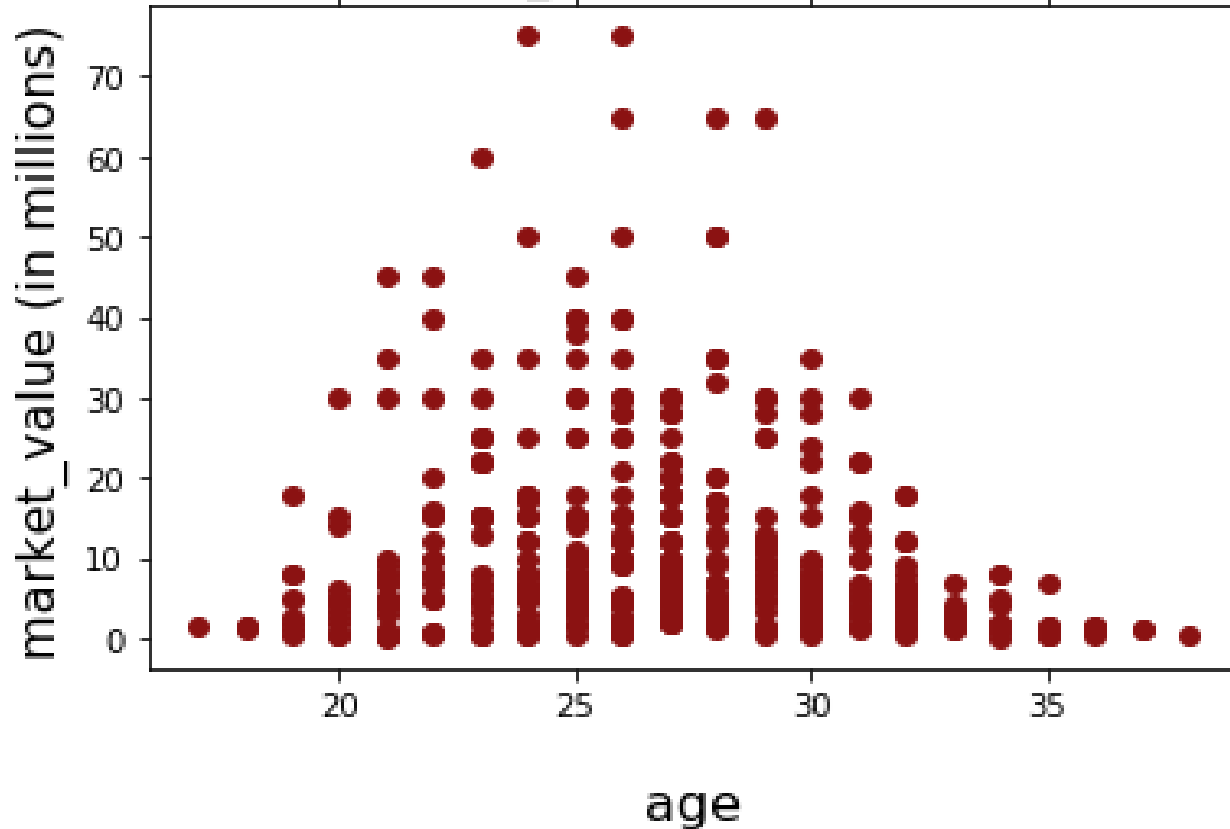
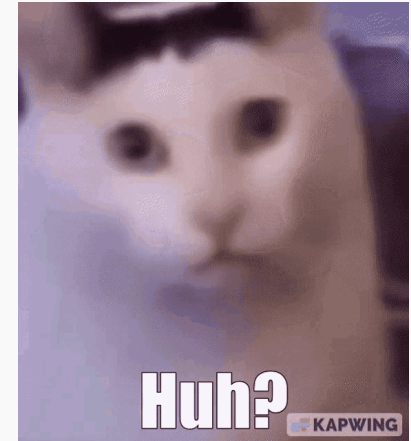# Visualization: The English Premier League



How age affects market_value in the English Premiere League

# Visualization: The English Premier League



How age affects market_value in the English Premiere League

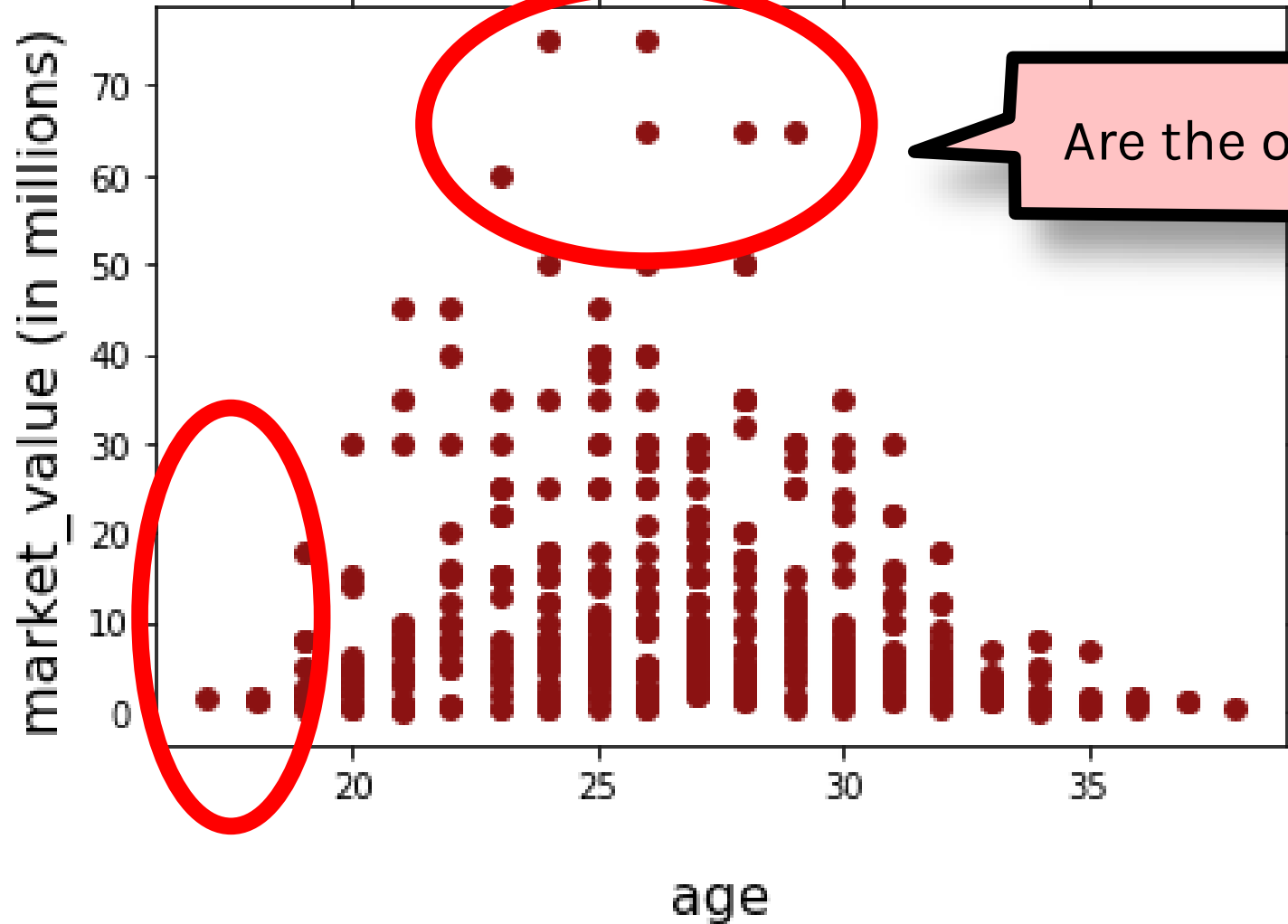**What do see in this graph? Is this sensible enough?**

**Are there any outliers?**

# Next steps:

- Ensure our data is expected/valid/appropriate for the task.
- Provide insights into the dataset.
- Extract/determine important variables/attributes/features.
- Detect outliers and anomalies.
- Test underlying assumptions.
- Make informed decisions in developing models.

# Next steps:

- Ensure our dat... ...ate for the task.
- Provide insight...
- Extract/determ... ...butes/features.
- Detect outliers...
- **Test underlying...**
- **Make informed...** ...dels.

⚠️**Fasten your seat belts for the upcoming sessions!**