

# Machine Learning

## Exploratory Analysis

---

Jan Platoš

November 22, 2023

Department of Computer Science  
Faculty of Electrical Engineering and Computer Science  
VŠB - Technical University of Ostrava

# Exploratory Analysis

---

- A matrix where rows are records and columns are features.
- The matrix usually contains many data types together.
- Combination of strings, numeric, categorical and other data.
- Analysis of such complex dataset is hard.
- Each feature has its own properties.

## Numeric values

- Numeric values may be integer, real or complex number.
- Integer values belong to a specific range, e.g. byte, short, int, long, etc.
- Real number are usually float, double or decimal.
- Complex numbers are rare case.
- Integer values are processed as real or as categorical type.

## Real numbers

- Real number columns has basic limits - minimum and maximum.
- Other important terms are mean, median, and quartile.
- Distribution of real values is the most important feature.
- Distribution depicts which values are present a how frequently.
- The overview of the feature may be taken from the histogram or a Box plot.

## Categorical data

- Categorical data represents a set of possibilities a feature may take.
- The data may be numeric or textual.
- Very frequent data type unusually together with numeric.
- The overview of the feature may be taken from the (discrete) histogram.
- Very frequent as a class/label for a data.

## Categorical data

- Processing is done using one of the following process:
  - binarization,
  - ordinal encoding (problem with sorting and distance),
  - one-hot encoding (dummy encoding),
  - embedding,
  - algorithmic encoding (cyclic feature).

## Textual data

- Textual data are in the form of a single word or an open text.
- Single word and short text may represent a categorical value (e.g. METAR).
- Open text columns are hard to process.
- Usually processed separately as a text data:
  - normalized, tokenized, encoded (embedding), ...



# Explorative Analysis

---

- What do you imagine?

- What do you imagine?
- What is the shape of the data?

- What do you imagine?
- What is the shape of the data?
- What is the distribution of the data?

- What do you imagine?
- What is the shape of the data?
- What is the distribution of the data?
- What is the distribution among the features?

- What do you imagine?
- What is the shape of the data?
- What is the distribution of the data?
- What is the distribution among the features?
- Is there any relation between features?

Exploratory Data Analysis refers to the critical process of performing initial investigations on data so as to **discover patterns**, to spot **anomalies**, to **test hypothesis** and to **check assumptions** with the help of **summary statistics** and **graphical representations**.<sup>1</sup>

---

<sup>1</sup><https://towardsdatascience.com/exploratory-data-analysis-8fc1cb20fd15>

- Dataset investigations
  - How many features it has?
  - How many records it has?
  - What is the content of the features?
  - What type of the features it has - Categorical, Numeric, Text?



- How many distinct features it has?
- Is it a class label or regular feature?
- What is the distribution of the values?

## Feature types - Numeric features

- What is the minimum and maximum?
- What are the quartiles?
- What is the distribution/histogram of the data?
- What are the properties of the data distribution?

- Mean  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$

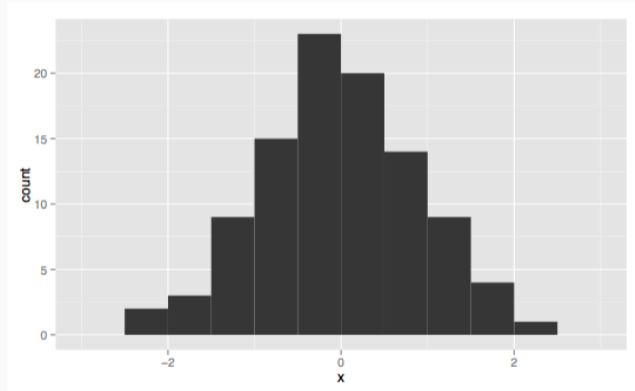
- Median - a middle value of sorted values.

- Variance  $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$

- Std. deviation  $\sigma = \sqrt{s^2}$

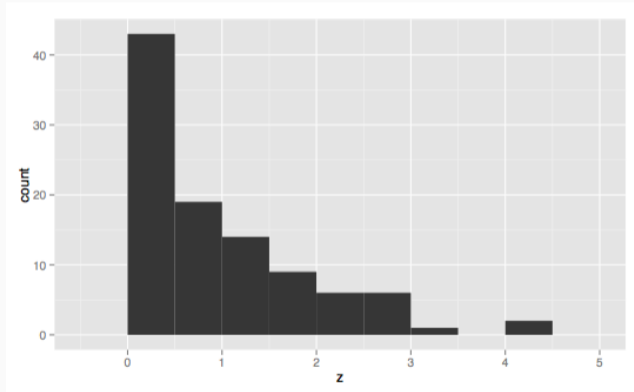
- Inter-quartile range  $IQR = Q_3 - Q_1$

- Histograms<sup>2</sup>



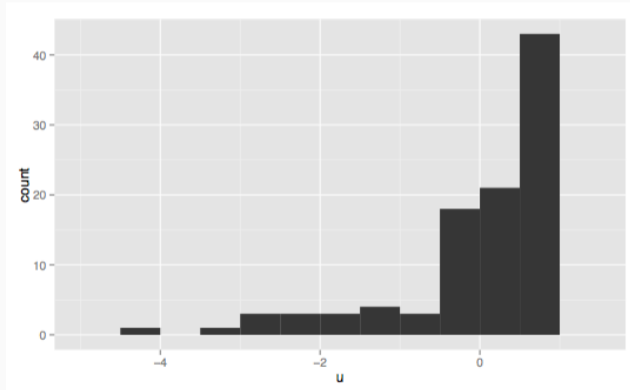
<sup>2</sup><https://en.wikipedia.org/wiki/Histogram>

- Histograms<sup>2</sup>



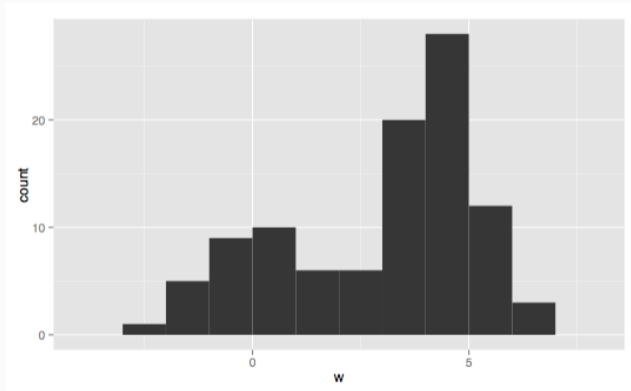
<sup>2</sup><https://en.wikipedia.org/wiki/Histogram>

- Histograms<sup>2</sup>



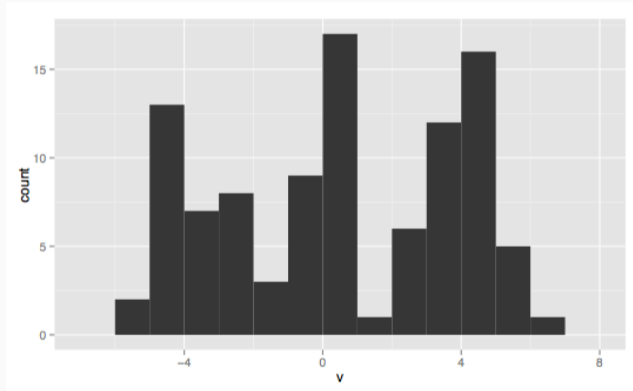
<sup>2</sup><https://en.wikipedia.org/wiki/Histogram>

- Histograms<sup>2</sup>



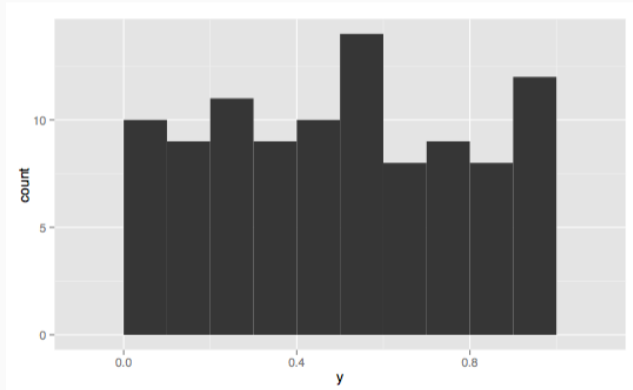
<sup>2</sup><https://en.wikipedia.org/wiki/Histogram>

- Histograms<sup>2</sup>



<sup>2</sup><https://en.wikipedia.org/wiki/Histogram>

- Histograms<sup>2</sup>



<sup>2</sup><https://en.wikipedia.org/wiki/Histogram>

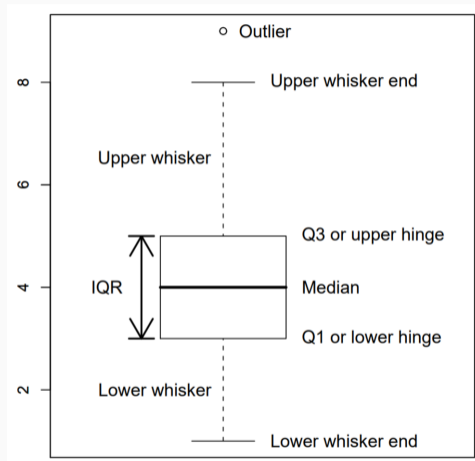


- Histograms
- Ideal for categorial data
- Numeric values requires another parameter - **Bin size**

# Feature types - Graphical analysis

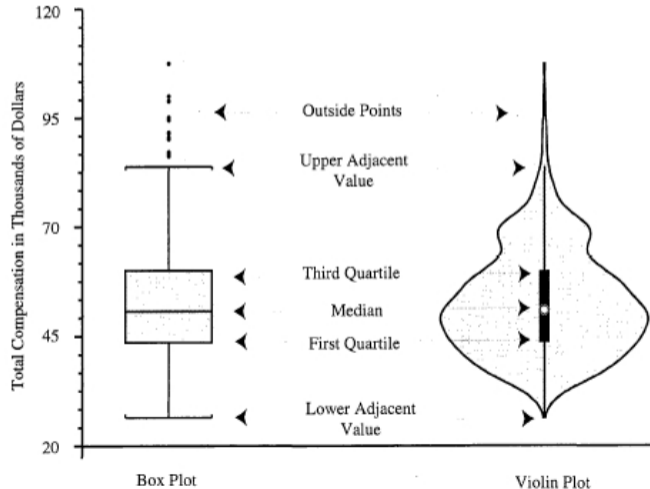
## Box plot (Box and Whiskers plot)

- Minimum : the lowest data point excluding any outliers.
- Maximum : the largest data point excluding any outliers.
- Median : the middle value of the dataset.
- First quartile : the median of the lower half of the dataset.
- Third quartile : the median of the upper half of the dataset.



# Feature types - Graphical analysis

- Violin plot<sup>3</sup>



## Relationship between features

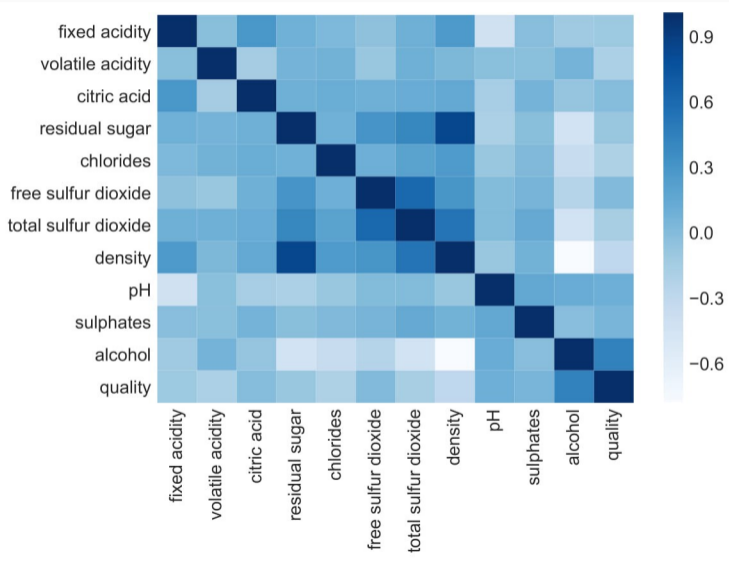
- There are many ways how to identify the relationships between features.
- **Covariance** - how much (and in what direction) should we expect one variable to change when the other changes.

$$\text{Cov}(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

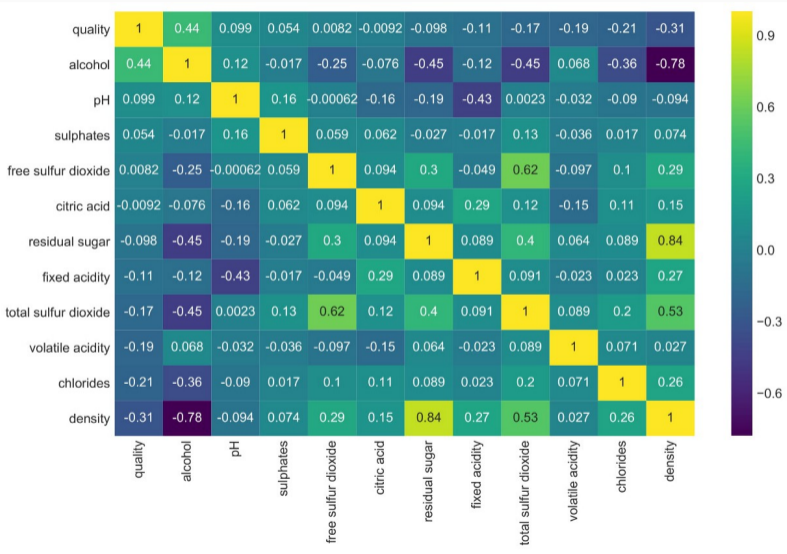
- **Correlation** - similarly to covariance, the power and direction of the relation.

$$\text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$$

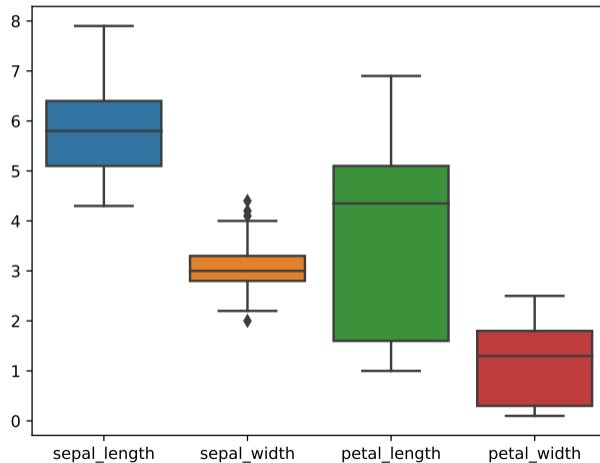
# Relationship between features - Correlation Heatmap



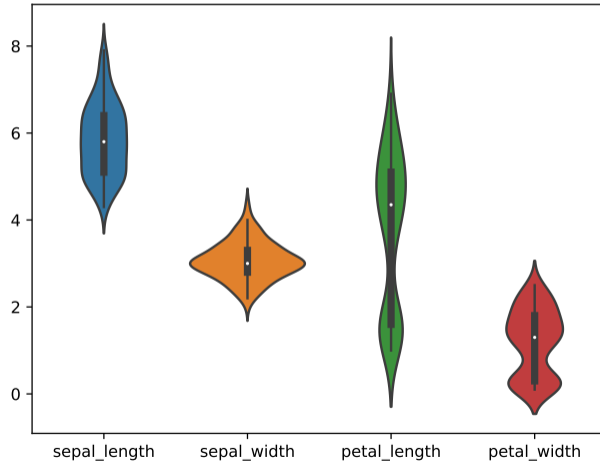
# Relationship between features - Correlation Heatmap



## Relationship between features - Box plot

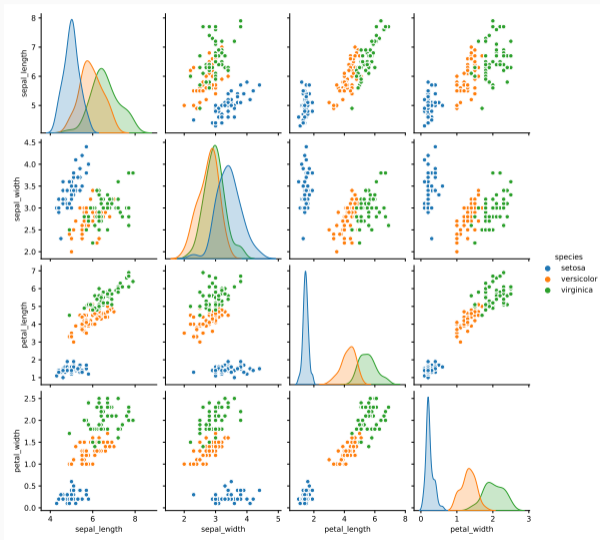


## Relationship between features - Violin plot

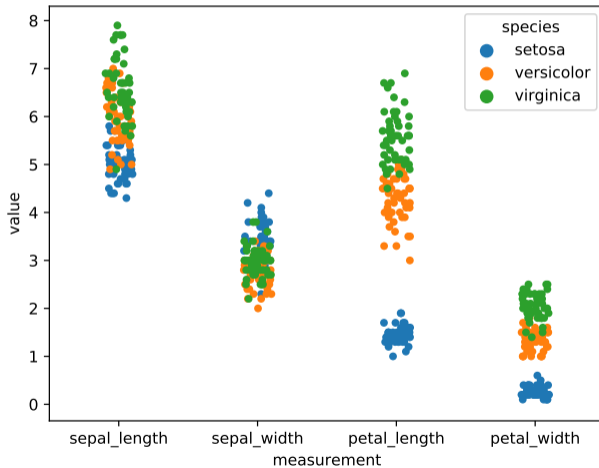




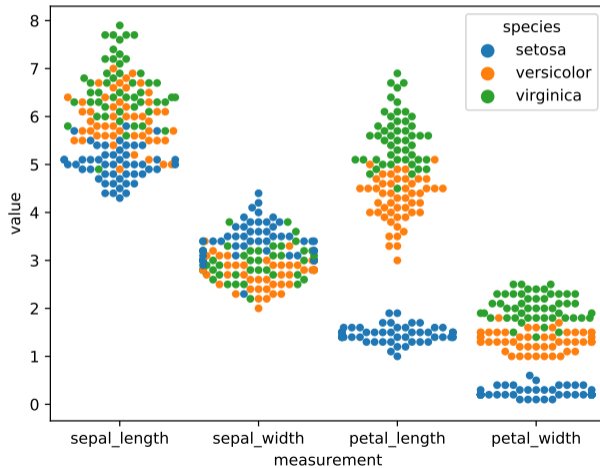
# Relationship between features - Pair plot



# Relationship between features - Strip plot



# Relationship between features - Swarm plot



## Feature scaling

---

- Feature scaling is a method used to normalize the range of independent variables or features of data.
- The process is known also as *data normalization*.
- The scaling is necessary to compare features between each other and use a metric to compute non-biased distance.
- The scaling have to respect the nature of the data.

## Min-Max Scaling

- The most common scaling principle.
- Unifies the *min* and *max* among the features.
- Normalize into range  $[0, 1]$

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

- Normalize into range  $[a, b]$

$$x' = a + \frac{(x - \min(x)) (b - a)}{\max(x) - \min(x)}$$

- Maximum Absolute Scaling where we normalize to  $|\max(x)|$

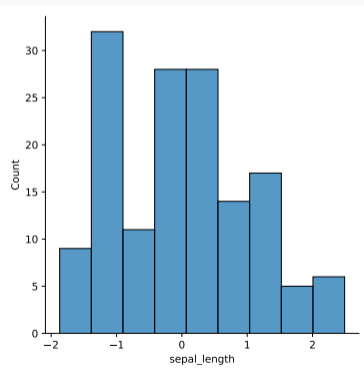
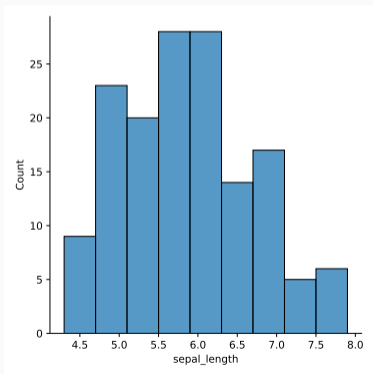
- Moves the mean of the data into 0.
- The final range is  $[-1, +1]$ .

$$x' = \frac{x - \bar{x}}{\max(x) - \min(x)}$$

# Standardizing

- Normalize the values to zero mean and unit variance.

$$X' = \frac{X - \bar{X}}{\sigma}$$





- Remapping features distribution into form close to Normal distribution
- Box-Cox transform or Yeo-Johnson transform

$$x^{(\lambda)} = \begin{cases} \frac{x_i^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \ln(x_i) & \text{if } \lambda = 0 \end{cases}$$

Questions?