



# Eksploracyjna analiza danych

## Wykład 3

October 14, 2021



Co-funded by the  
Erasmus+ Programme  
of the European Union



- 1 Terminologia danych
- 2 Eksploracja danych
- 3 Eksploracja danych za pomocą statystyk podsumowujących
- 4 Eksploracja danych poprzez wykresy
- 5 Przykładowe rozwiązania z użyciem różnych narzędzi
  - Description of Daily Weather Dataset
  - Exploring Data with KNIME Plots
  - Eksploracja danych w pandas
  - Data Exploration in Spark

# Terminologia danych I

Po tej sekcji będziesz mógł...

- Opisać czym jest dana cecha (feature) i jaki jest jej związek z próbką
- Nazwać kilka alternatywnych terminów dla „feature”
- Podsumować, jak kategoriowa cecha różni się od cechy liczbowej

	number	air_pressure_9am	air_temp_9am	avg_wind_direction_9am	avg_wind_speed_9am	max_wind_direction_9am
0	0	918.060000	74.822000	271.100000	2.080354	295.400000
1	1	917.347688	71.403843	101.935179	2.443009	140.471548
2	2	923.040000	60.638000	51.000000	17.067852	63.700000
3	3	920.502751	70.138895	198.832133	4.337363	211.203341
4	4	921.160000	44.294000	277.800000	1.856660	136.500000

Figure: **Samples** (Próbki) i **Variables** (Zmienne) są prezentowane odpowiednio jako wiersze i kolumny tabeli



## Inne nazwy dla 'Sample' (Próbka)

- record (rekord)
- example (przykład)
- row (wiersz)
- instance (instancja)
- observation (obserwacja)

## Inne nazwy dla 'Variable' (Zmienna)

- attribute (atrybut)
- field (pole)
- feature (cecha)
- column (kolumna)
- dimension (wymiar)

## Typy danych

- Najczęściej używane
  - Numeric (Liczbowy)
  - Categorical (Kategoryjny)
- Inne
  - String (Ciąg znaków)
  - Date (Data)
  - ...

## Liczbowe zmienne

- Wartości są liczby
- Nazywany również „ilościowym”
- Przykłady: 1

163.92

$7 \times 10^5$



## Przykłady zmiennych liczbowych

- Wysokość
- Wynik na egzaminie
- Liczba transakcji na godzinę
- Zmiana ceny akcji

## Kategoryjne Zmienne

- Wartości to etykiety, nazwy lub kategorie
- Nazywany również „jakościowym” lub „nominalnym”

Kolor
Red
Silver
Blue
White
Black

**Table:** Kolor jest zmienną kategoryczną; Wartości są etykietami

## Przykłady zmiennych kategorialnych

- Płeć
- Stan cywilny
- Typ klienta
- Kategorie produktów

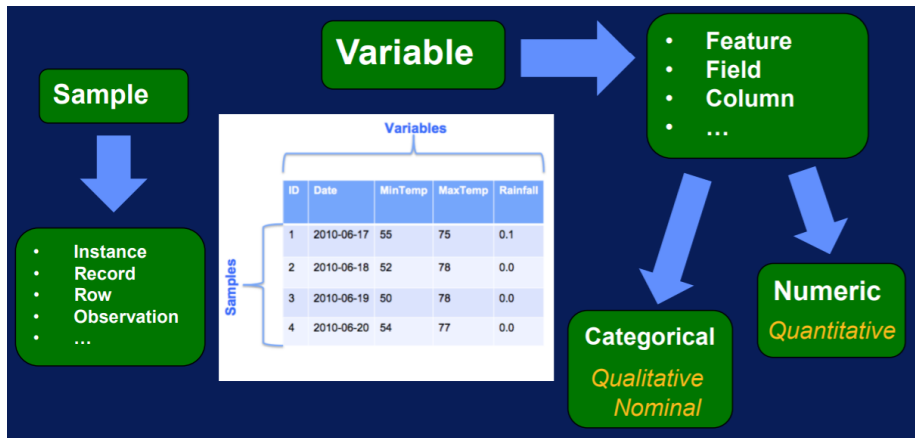
• Kolor elementów



Co-funded by the  
Erasmus+ Programme  
of the European Union



# Terminologia danych VI





# Eksploracja danych I

Po tej sekcji będziesz mógł...

- Wyjaśnić, dlaczego eksploracja danych jest konieczna
- Przedstawić cele eksploracji danych
- Wymienić kategorie technik eksploracji danych

"Experts often possess more data than judgment"  
(by Colin Powell)

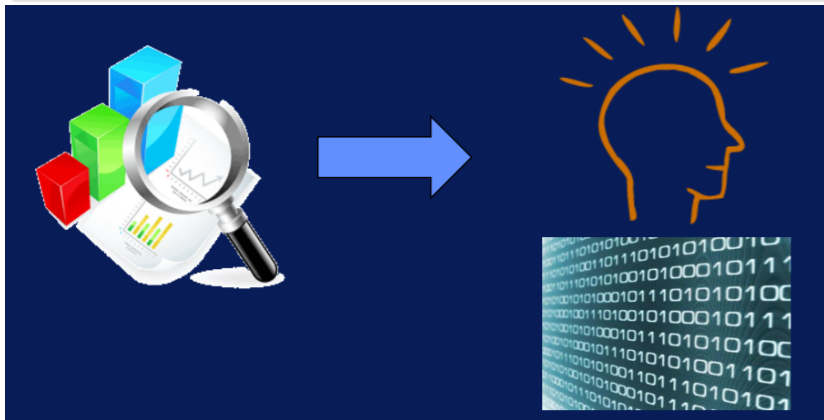
## Eksploracja Danych

Gdy już zidentyfikujesz pytania, na które próbujesz uzyskać odpowiedź, i zdobędziesz pewne dane, możesz ulec pokusie, aby od razu zanurkować i zacząć budować modele i uzyskiwać odpowiedzi. Musisz jednak oprzeć się tej pokusie. Pierwszym krokiem powinna być eksploracyjna analiza danych.



## Dlaczego eksplorować dane?

**Cel:** Aby zrozumieć swoje dane



Eksploracyjna Analiza Danych (EDA)

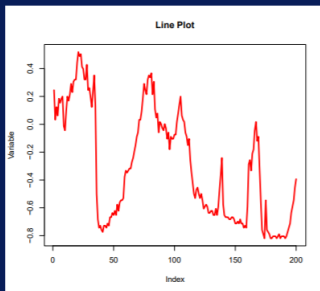
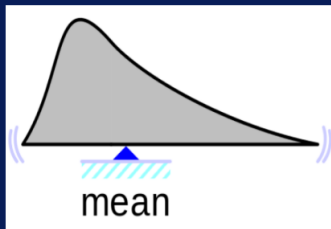


Co-funded by the  
Erasmus+ Programme  
of the European Union



# Eksploracja danych III

## Sposoby eksploracji danych



Summary  
Statistics

Visualization

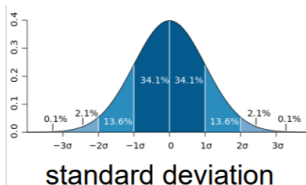
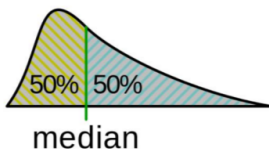
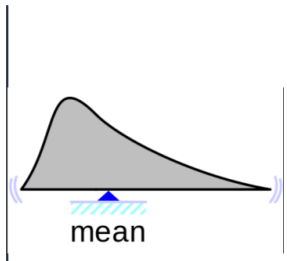


Co-funded by the  
Erasmus+ Programme  
of the European Union



## Statystyki podsumowujące

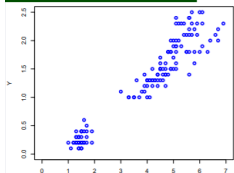
- Informacje podsumowujące zbiór danych



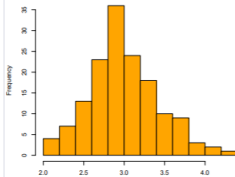
## Wizualizacja danych

- Spójrz na dane graficznie

Scatter Plot

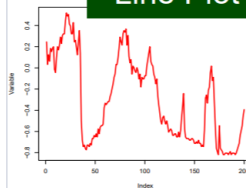


Histogram



Histogram

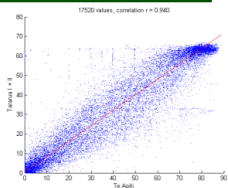
Line Plot



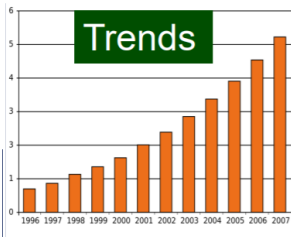
# Eksploracja danych VI

Kilka rzeczy, na które należy zwrócić uwagę

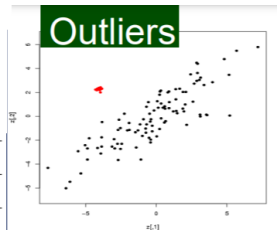
## Correlations



## Trends



## Outliers

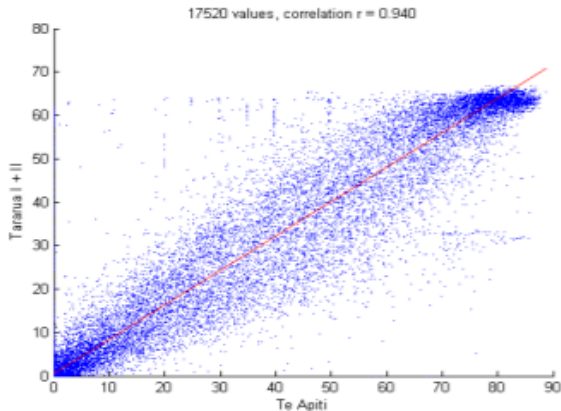


Co-funded by the  
Erasmus+ Programme  
of the European Union



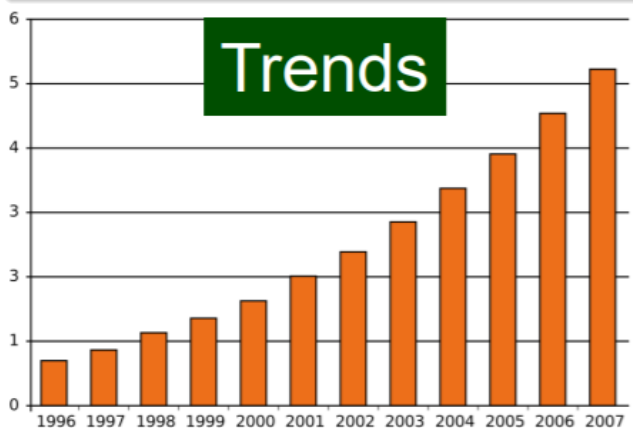
## Korelacje

- Podają informacje o relacji między zmiennymi



## Trendy

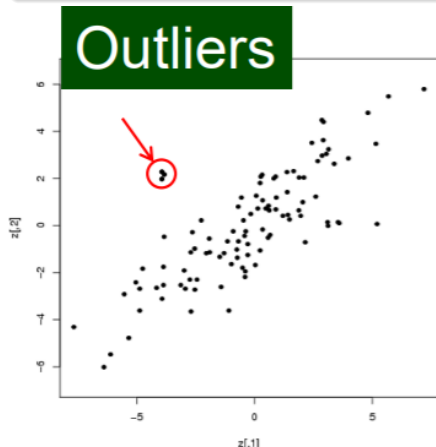
- Wskazują ogólną charakterystykę danych



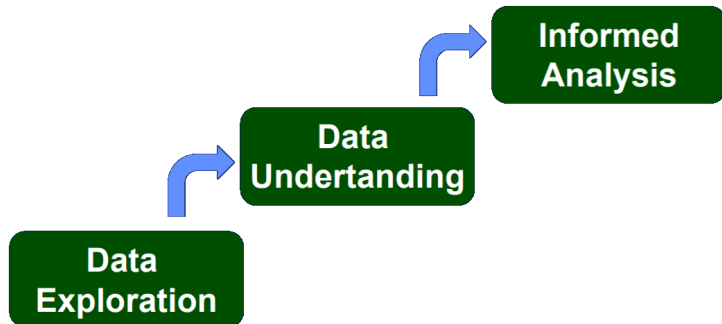


## (Outliers) Odstające

- Wskazują potencjalne problemy z danymi



Eksploracja danych



Po tej sekcji będziesz mógł...

- Zdefiniować, czym jest statystyka podsumowująca
- Wymienić trzy popularne statystyki podsumowujące
- Wyjaśnić, w jaki sposób statystyki podsumowujące są przydatne w badaniu danych

## Definition

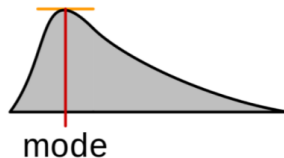
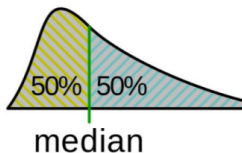
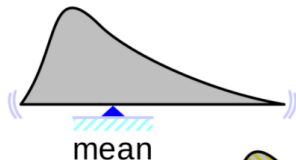
Co to są statystyki podsumowujące?

- Ilości podsumowujące i opisujące zbiór wartości danych
- Miary
  - Lokalizacja: średnia, mediana
  - Rozrzut: odchylenie standardowe
  - Kształt: skośność

# Data Exploration through Summary Statistics II

## Miary lokalizacji

Opisują centralną lub typową wartość zbioru danych



# Data Exploration through Summary Statistics III

## Miry lokalizacji - przykład

Wiek	Wiek (posortowane)
35	21
42	22
78	35
22	42
56	42
50	50
42	56
78	78
21	78
87	87

Mean = 51.1 Median =  $(42+50)/2 = 46$  Mode = 42 & 78

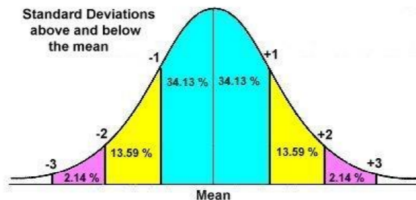


# Data Exploration through Summary Statistics IV

## Miary rozprzestrzeniania się

Opisz, jak rozproszone lub zróżnicowane są dane

- minimum
- maximum
- range
- standard deviation
- variation



# Data Exploration through Summary Statistics V

## Miary rozprzestrzeniania się – Przykład

$$\text{Range} = 87 - 21 = 66$$

$$\text{Variance} = 548.767$$

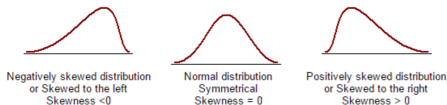
$$\text{Standard deviation} = 23.426$$



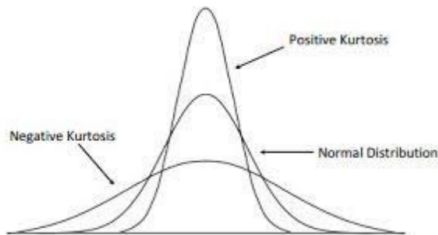
# Data Exploration through Summary Statistics VI

## Miary kształtu

skewness (skośność)



kurtosis (kurtoza)





# Data Exploration through Summary Statistics VII

## Miary kształtu - Przykład

Age
35
42
78
22
56
50
42
78
21
87

Skewness = 0.2995

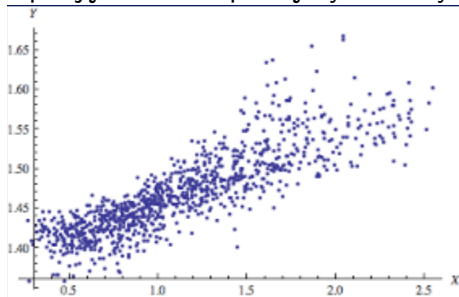
Kurtosis = -1.2028



# Data Exploration through Summary Statistics VIII

## Miary zależności

Opisują zależność pomiędzy zmiennymi



korelacja



# Data Exploration through Summary Statistics IX

Height	Weight
180	68
153	70
204	84
133	44
208	81
142	53
122	40
168	50
175	64
200	72

Correlation = 0.8906

Miary zależności – Przykład



Co-funded by the  
Erasmus+ Programme  
of the European Union



# Data Exploration through Summary Statistics X

## Statystyki dotyczące zmiennych kategorialnych

Opisują liczbę kategorii i częstotliwość każdej kategorii

Color/Pet	White	Brown	Black	Orange	Total
Dog	34	44	32	0	110
Cat	25	2	43	0	70
Fish	1	0	5	33	39
<b>Total</b>	60	46	80	33	219

contingency table (tablica kontyngencji)



# Eksploracja danych poprzez wykresy I

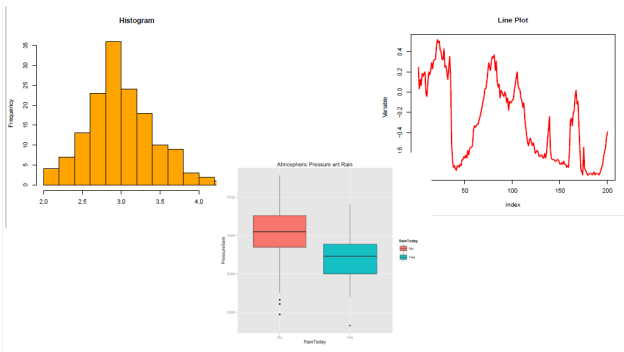
Po tej sekcji będziesz mógł ...

- omówić, w jaki sposób wykresy mogą być przydatne w badaniu danych
- opisać, jak wykorzystałbyś wykres punktowy
- podsumować, co pokazuje wykres pudełkowy



# Eksploracja danych poprzez wykresy II

## Wizualizacja danych



# Eksploracja danych poprzez wykresy III

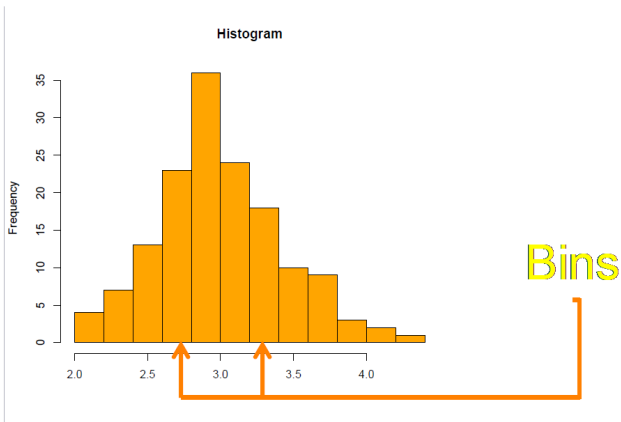
## Typy wykresów

- Histogram
- Line plot
- Scatter plot (punktowy)
- Bar plot (słupkowy)
- Box plot (pudełkowy)
- inne



## Histogram

pokazuje rozkład zmiennej numerycznej





# Eksploracja danych poprzez wykresy V

## Co pokazuje histogram

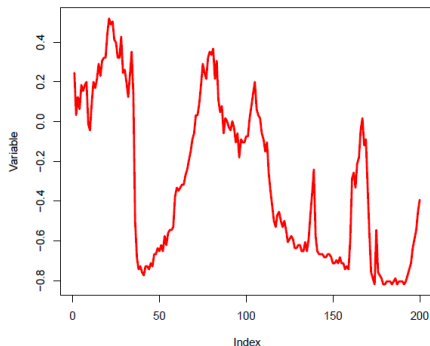


# Eksploracja danych poprzez wykresy VI

## Wykres liniowy

Pokazuje zmiany danych w czasie

Line Plot

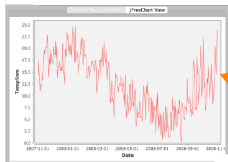


Co-funded by the  
Erasmus+ Programme  
of the European Union



# Eksploracja danych poprzez wykresy VII

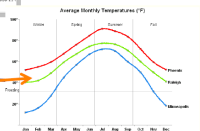
## Co pokazuje wykres liniowy



Trend  
Cyclical  
pattern



Compare  
variables

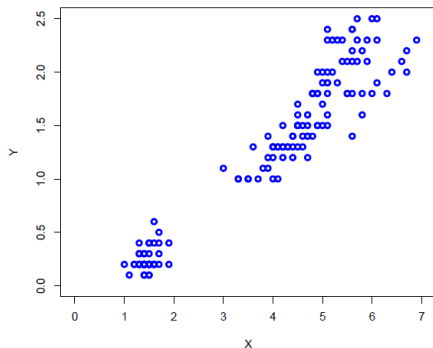


# Eksploracja danych poprzez wykresy VIII

## Wykres punktowy

Pokazuje związek między dwiema zmiennymi

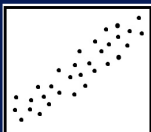
Scatter Plot



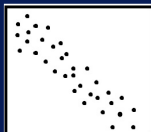
# Eksploracja danych poprzez wykresy IX

Co pokazuje wykres punktowy?

Positive  
Correlation



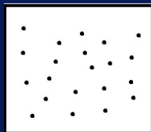
Negative  
Correlation



Non-  
Linear  
Correlation

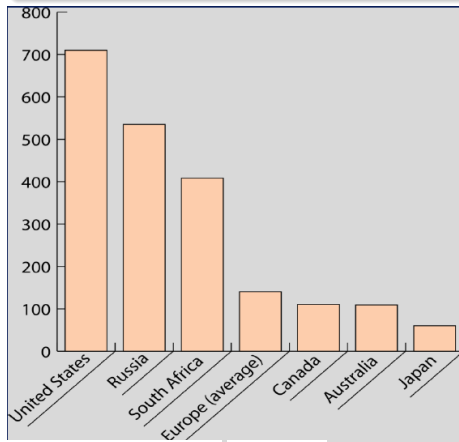


No Correlation



## Wykres słupkowy (Bar plot)

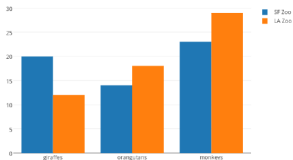
Pokazuje rozkład zmiennej kategoryjnej



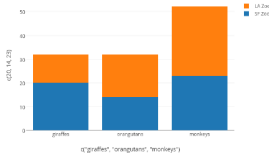
# Eksploracja danych poprzez wykresy XI

Co pokazuje wykres słupkowy

## Grouped Bar Chart



## Stacked Bar Chart



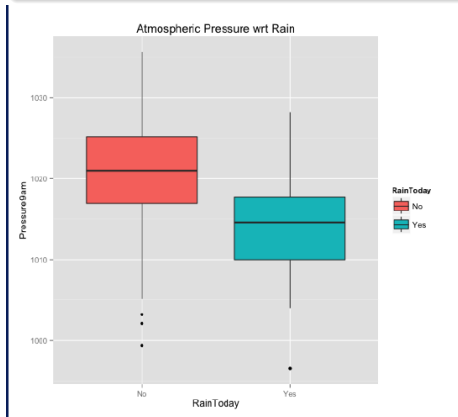
Co-funded by the  
Erasmus+ Programme  
of the European Union



# Eksploracja danych poprzez wykresy XII

## Wykres pudełkowy (Box plot)

Porównuje rozkłady zmiennych



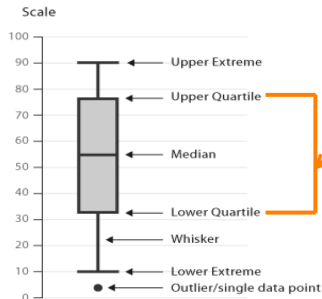
Co-funded by the  
Erasmus+ Programme  
of the European Union





# Eksploracja danych poprzez wykresy XIII

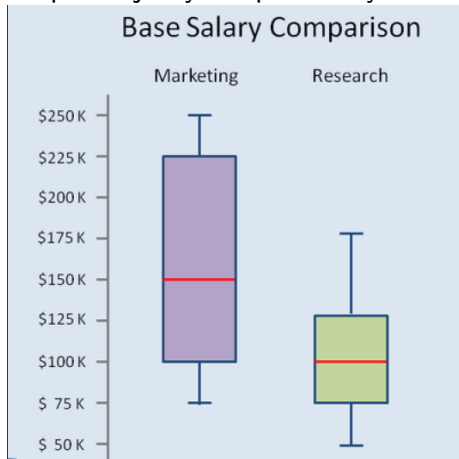
## Składniki wykresu pudełkowego



The middle 50% of data are in this region

# Eksploracja danych poprzez wykresy XIV

Co pokazuje wykres pudełkowy

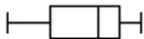
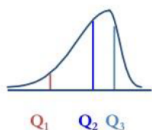


# Eksploracja danych poprzez wykresy XV

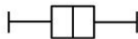
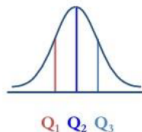
Co pokazuje wykres pudełkowy

## Distribution Shape and The Boxplot

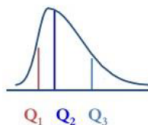
Negative Skew



Symmetric



Positive Skew



## Wizualizacja danych

- Zapewnia intuicyjny sposób przeglądania danych
- Powinna być używana ze statystykami podsumowującymi do eksploracji danych
- jest również przydatna do komunikowania wyników

- 1 Terminologia danych
- 2 Eksploracja danych
- 3 Eksploracja danych za pomocą statystyk podsumowujących
- 4 Eksploracja danych poprzez wykresy
- 5 Przykładowe rozwiązania z użyciem różnych narzędzi**
  - Description of Daily Weather Dataset
  - Exploring Data with KNIME Plots
  - Eksploracja danych w pandas
  - Data Exploration in Spark

# Description of Daily Weather Dataset I

## Description of Daily Weather Dataset



Co-funded by the  
Erasmus+ Programme  
of the European Union



- 1 Terminologia danych
- 2 Eksploracja danych
- 3 Eksploracja danych za pomocą statystyk podsumowujących
- 4 Eksploracja danych poprzez wykresy
- 5 Przykładowe rozwiązania z użyciem różnych narzędzi**
  - Description of Daily Weather Dataset
  - Exploring Data with KNIME Plots**
  - Eksploracja danych w pandas
  - Data Exploration in Spark

# Exploring Data with KNIME Plots I

The screenshot displays the KNIME Analytics Platform interface. On the left, the **KNIME Explorer** shows a project structure with 'LOCAL (Local Workspace)' selected, containing 'Example Workflows' and 'Plots\_of\_weather'. Below it, the **Workflow Coach** shows 'Recommended Nodes' and 'Color Manager' (set to 11%). The **Node Repository** lists various node categories like IO, Manipulation, Views, Analytics, DB, etc.

The central canvas shows a workflow with five nodes:

- File Reader (Node 1)**: Connected to 'Plots\_of\_weather'.
- Numeric Binner (Node 4)**: Connected to File Reader.
- Scatter Plot (Node 3)**: Connected to File Reader.
- Histogram (local) (Node 2)**: Connected to File Reader.
- Histogram (local) (Node 5)**: Connected to Numeric Binner.

On the right, the **Scatter Plot** configuration panel is visible, showing a description: 'A scatter plot using a JavaScript based charting library. The view can be accessed either via the "interactive view" action on the executed node or in KNIME Server web portal page.' It also mentions that the configuration of the node lets you choose the size of a sample to display and enable certain controls.

The bottom section includes an **Outline** view, a **Console** window showing the KNIME welcome message and copyright information, and a **Node Monitor** window.



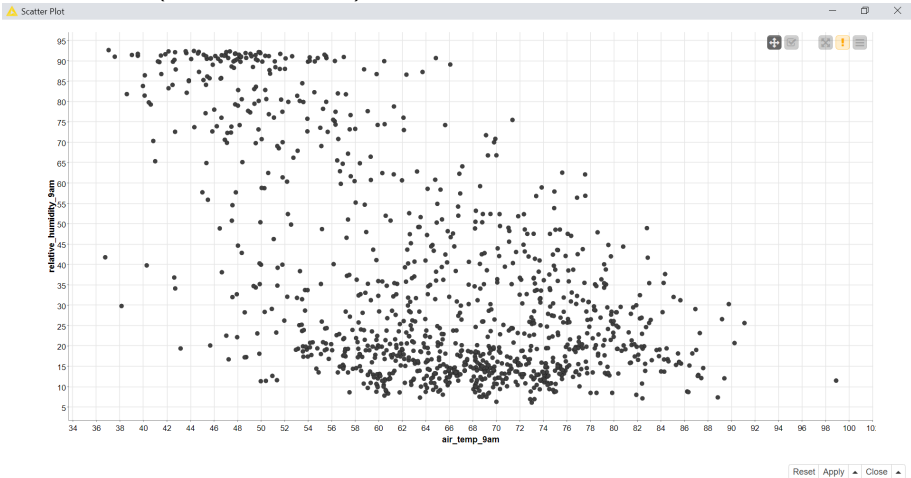
Co-funded by the  
Erasmus+ Programme  
of the European Union





# Exploring Data with KNIME Plots II

## Scatter plot (interactive view)



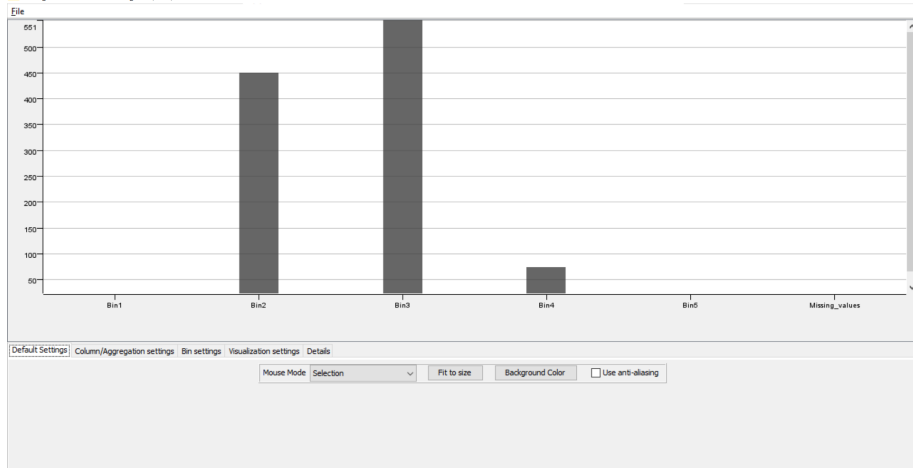
Co-funded by the  
Erasmus+ Programme  
of the European Union



# Exploring Data with KNIME Plots III

## Histogram

▲ Histogram View - 0.5 - Histogram (local)



Co-funded by the  
Erasmus+ Programme  
of the European Union



- 1 Terminologia danych
- 2 Eksploracja danych
- 3 Eksploracja danych za pomocą statystyk podsumowujących
- 4 Eksploracja danych poprzez wykresy
- 5 Przykładowe rozwiązania z użyciem różnych narzędzi
  - Description of Daily Weather Dataset
  - Exploring Data with KNIME Plots
  - **Eksploracja danych w pandas**
  - Data Exploration in Spark

```
http://localhost:8889/lab/tree/OneDrive/Documents/  
PythonProjects/Tutorials/Jupyter_Polish_python_lessons-main/  
Wyklad3/Wyklad3_Exploratory_data_Analysis.ipynb
```



- 1 Terminologia danych
- 2 Eksploracja danych
- 3 Eksploracja danych za pomocą statystyk podsumowujących
- 4 Eksploracja danych poprzez wykresy
- 5 Przykładowe rozwiązania z użyciem różnych narzędzi**
  - Description of Daily Weather Dataset
  - Exploring Data with KNIME Plots
  - Eksploracja danych w pandas
  - Data Exploration in Spark**

# Data Exploration in Spark I

In order to **install pyspark** you need to run:

```
! pip install pyspark

# Import PySpark
import pyspark
from pyspark.sql import SparkSession

#Create SparkSession
spark = SparkSession.builder.master("local[1]")
    .appName("SparkByExamples.com").getOrCreate()
sc=spark.sparkContext

from pyspark.sql import SQLContext

# creating sqlContext
sqlContext = SQLContext(sc)
```



# Data Exploration in Spark II

```
creating spark dataframe
df = sqlContext.read.load('file:///C:/Users/marce/OneDrive
/Documents/PythonProjects/Tutorials/DataExploration
/coursera-master/coursera-master/big-data-4
/daily_weather.csv',
format = 'com.databricks.spark.csv',
header = 'true', inferSchema='true')

# display the names of columns
df.columns

# display scheme of data (types of data and hierarchy)
df.printSchema()

# display summary statistics
df.describe().toPandas().transpose()
```

# Data Exploration in Spark III

	0	1	2	3	4
summary	count	mean	stddev	min	max
number	1095	547.0	316.24357700987383	0	1094
air_pressure_9am	1092	918.8825513138094	3.184161180386833	907.99000000000024	929.32000000000012
air_temp_9am	1090	64.93300141287072	11.175514003175877	36.7520000000000685	98.90599999999992
avg_wind_direction_9am	1091	142.2355107005759	69.13785928889189	15.500000000000046	343.4
avg_wind_speed_9am	1092	5.50828424225493	4.5528134655317185	0.69345139999974	23.554978199999763
max_wind_direction_9am	1092	148.95351796516923	67.23801294602953	28.899999999999991	312.19999999999993
max_wind_speed_9am	1091	7.019513529175272	5.598209170780958	1.1855782000000479	29.84077959999996
rain_accumulation_9am	1089	0.20307895225211126	1.5939521253574893	0.0	24.019999999999907
rain_duration_9am	1092	294.1080522756142	1598.0787786601481	0.0	17704.0
relative_humidity_9am	1095	34.24140205923536	25.472066802250055	6.0900000000001012	92.62000000000002
relative_humidity_3pm	1095	35.34472714825898	22.524079453587273	5.30000000000006855	92.25000000000003

```
# summary statistics for a column
df.select('air_pressure_9am').show()
```



Co-funded by the  
Erasmus+ Programme  
of the European Union





# Data Exploration in Spark IV

```
# drop non-available value
df2 = df.na.drop(subset=['air_pressure_9am'])
# compute correlation
df2.stat.corr('rain_accumulation_9am', 'rain_duration_9am')
```

