

Politechnika Rzeszowska

Katedra Informatyki i Automatyki

Sprawozdanie z przedmiotu
Analiza danych w językach R i Python

TEMAT PROJEKTU

Prowadzący:

Krzysztof Kołodziej, L01,
169562@stud.prz.edu.pl
Jakub Duda, L01,
169532@stud.prz.edu.pl

Sprawozdanie z przedmiotu Analiza danych w językach R i Python, KiIA PRz

Tutaj znajduje się informacja o tym w jaki sposób należy dostarczyć kolejne etapy sprawozdania prowadzącym.

Kolejne etapy sprawozdania należy dostarczyć prowadzącemu w podanych niżej terminach za pomocą strony **sprawozdania.kia.prz.edu.pl**. Należy tam w pierwszej kolejności założyć konto (jeżeli jeszcze takowego nie ma) i dołączyć do odpowiedniej grupy projektowej z przedmiotu. Każdy kolejny etap sprawozdania należy dopisać do pliku zawierającego poprzedni, już wysłany, tak, aby ostatnie sprawozdanie było kompletne. Jeśli kolejny etap wymaga poprawek w treści poprzedniego są one dozwolone. Przed każdorazowym wysłaniem sprawozdania należy odświeżyć pole spisu treści!

Każde sprawozdanie ma być wysłane w załączniku maila z tematem:

ADRP<numer_rozdziału>_L<numer_grupy>_<I1>.<Nazwisko1>_<I2>.<Nazwisko2>

na przykład:

ADRP_R1_L01_J.Kowalski_M.Nowak

Ostatnie sprawozdanie powinien zawierać też załącznik z implementacją projektu, plikami danych opisującymi problemy rozważane w sprawozdaniu oraz plikami wyników jeśli takie pliki istnieją.

Każdy wykonany projekt trzeba obronić podczas osobistej rozmowy grupy projektowej z prowadzącym w wyznaczonym terminie. Na tą rozmowę należy być przygotowanym do zaprezentowania działającego projektu. Każdy członek grupy zostanie oceniony osobno na podstawie jego wkładu w wykonanie projektu.

Terminy nadsyłania kolejnych sprawozdania:

Rozdział sprawozdania	Termin
Rozdział numer 1.	28 X
Rozdział numer 2.	25 XI
Rozdział numer 3. etapów	23 XII

Pozostałe rozdziały należy wypełnić podczas pracy nad zasadniczymi trzema wymienionymi wyżej. Niewywiązanie się z powyższych terminów będzie skutkowało drastycznym obniżeniem oceny z przedmiotu. Termin oddawania i obrony projektów zostanie ustalona ze starostą roku.

Spis treści

1. Opis problematyki klasyfikacji danych	4
1.1 Wstęp.....	4
1.2 Etapy procesu klasyfikacji.....	4
1.3 Metody stosowane do poprawności wiarygodności i dokładności budowanych modeli	7
2. Charakterystyka danych.....	9
2.1 Źródło i typ danych	9
2.2 Informacje charakteryzujące dane	9
2.3 Przygotowanie danych do budowy modelu klasyfikacji.....	20
2.4 Normalizacja danych.....	22
3. Opis eksperymentów i opracowanie uzyskanych wyników	23
Podsumowanie	24
Bibliografia	25

1. Opis problematyki klasyfikacji danych

1.1 Wstep

Klasyfikacja danych jest kluczowym elementem eksploracji danych, czyli procesu automatycznego odkrywania nieznanych wcześniej zależności i wzorców w zbiorach danych. Proces ten pozwala na przekształcenie surowych danych w wartościowe informacje, które mogą wspierać podejmowanie decyzji w różnych dziedzinach. Informacje uzyskane dzięki klasyfikacji mogą być wykorzystywane m.in. do przewidywania zachowań czy generowania raportów. [2]

Klasyfikację można podzielić na: [2]

- **Klasyfikacja dwuklasowa:** Przypisuje obiekt do jednej z dwóch dostępnych klas (np. "ssak" lub "nie ssak").
- **Klasyfikacja wieloklasowa:** Przyporządkowuje obiekt do jednej z wielu klas (np. rozpoznanie gatunku zwierzęcia).

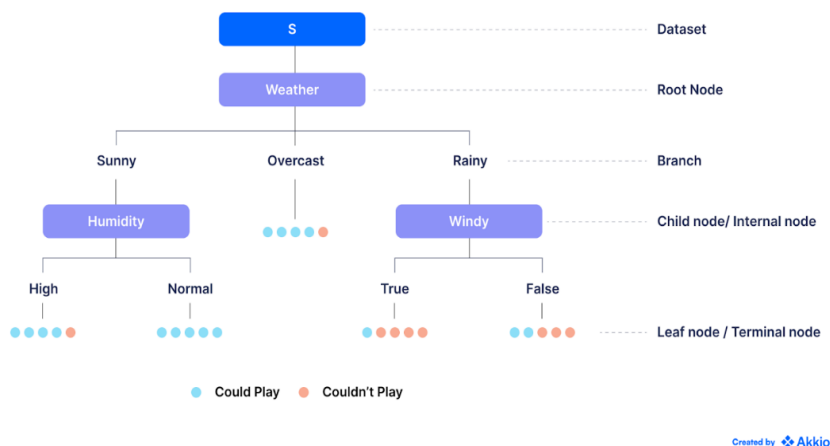
1.2 Etapy procesu klasyfikacji

Proces klasyfikacji można podzielić na dwa główne etapy. Pierwszy to **etap uczenia się**, w którym algorytm analizuje dane treningowe, aby nauczyć się wzorców i reguł, które później będą wykorzystywane do przypisywania nowych danych do odpowiednich klas. W tym etapie model buduje swoją wiedzę na podstawie dostarczonych informacji. Drugi to **etap predykcji**, w którym model, opierając się na wyuczonych wzorcach, dokonuje klasyfikacji nowych danych, przypisując je do odpowiednich kategorii. [2]

Klasyfikację można zaliczyć do metod **nadzorowanych**, ponieważ algorytm uczy się na wcześniej oznaczonych danych. Oznacza to, że każdy obiekt w zbiorze treningowym ma już przypisaną klasę, a zadaniem modelu jest nauczyć się, jak na podstawie cech obiektu przypisać go do odpowiedniej kategorii. W fazie uczenia algorytm otrzymuje dane wejściowe wraz z prawidłowymi wynikami (czyli z etykietami klas), zakładając, że są one prawidłowe, co pozwala mu wyciągać wnioski dotyczące przyszłych klasyfikacji. [2]

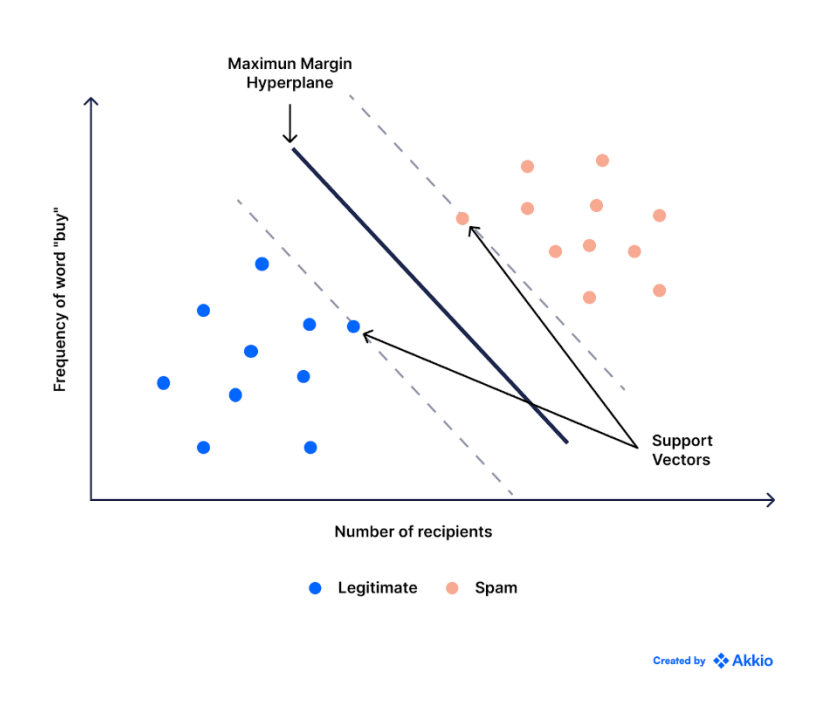
Przykładami algorytmów stosowanych w klasyfikacji nadzorowanej są:

- **Drzewa decyzyjne:** Bardzo powszechny sposób przedstawiania i wizualizacji możliwych wyników decyzji lub działania na podstawie prawdopodobieństw. Każde drzewo jest hierarchiczną reprezentacją przestrzeni wyników, w której każdy węzeł reprezentuje działanie lub wybór, a liście reprezentują stany wyniku. Są łatwe w zaimplementowaniu i szczególnie przydatne w sytuacjach, gdy istotna jest przejrzystość procesu decyzyjnego oraz łatwość interpretacji wyników [1]



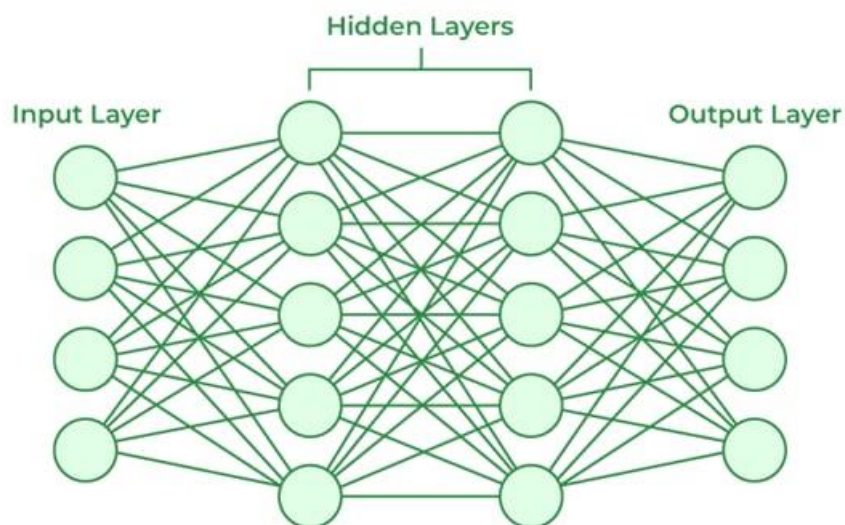
Rysunek 1 - algorytm drzewa decyzyjnego [1]

- **Maszyna wektorów nośnych (SVM):** Algorytm, który tworzy granicę decyzyjną między klasami, tak aby maksymalizować odległość między najbliższymi przykładami różnych klas. Jest szczególnie dobry w oddzielaniu podobnych danych. Jest ceniony za skuteczność w sytuacjach i stosowany w sytuacjach, gdy dane są trudne do rozdzielenia, ponieważ tworzy hiperplan maksymalnie oddzielający różne klasy, co zmniejsza ryzyko błędnej klasyfikacji. [1]



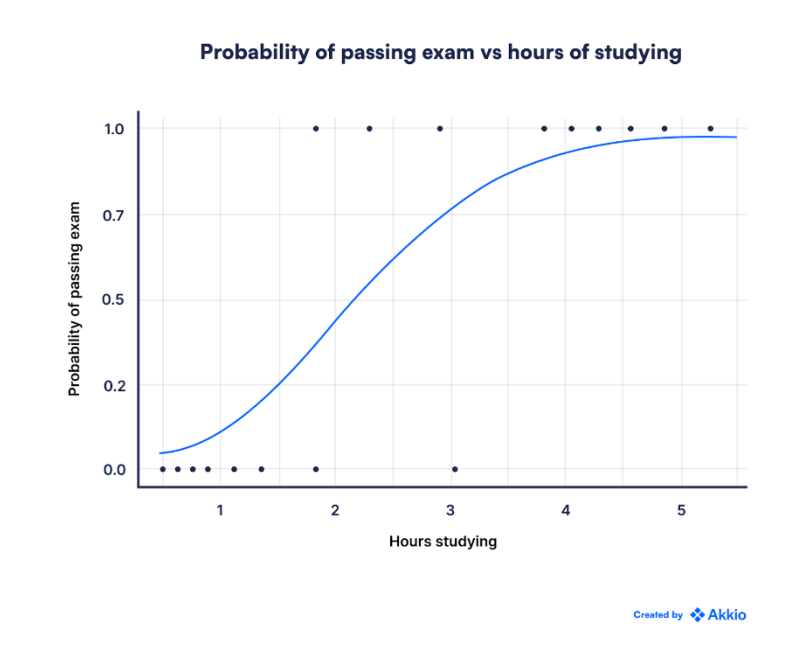
Rysunek 1 - algorytm SVM [1]

- **Sztuczne sieci neuronowe:** sama nazwa wskazuje, emulują biologiczne sieci neuronowe w komputerach. Sieci neuronowe działają poprzez trenowanie zestawu parametrów na danych. Parametry te są następnie używane do określania wyników modelu. Sieci mogą być jednokierunkowe lub rekurencyjne w których występuje sprzężenie zwrotne. Sztuczne sieci neuronowe mają szerokie zastosowanie w wielu dziedzinach, zwłaszcza tam, gdzie występuje duża złożoność danych lub trudność w wykrywaniu wzorców. Dzięki swojej zdolności do przetwarzania dużych ilości danych i "uczenia się" z nich, sieci neuronowe znajdują zastosowanie w takich obszarach, jak: przetwarzanie języka naturalnego, rozpoznawanie obrazów i wideo itp. [1]



Rysunek 2 - sztuczne sieci neuronowe [3]

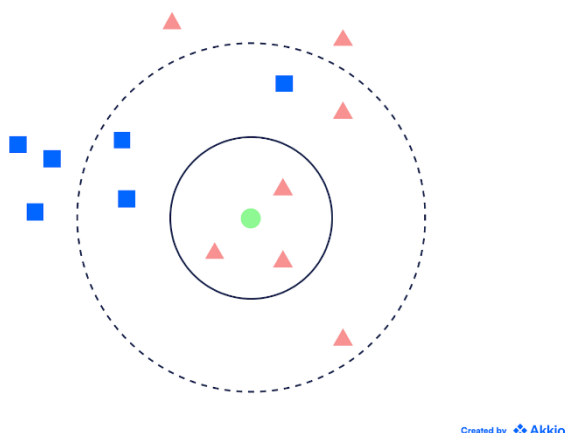
- **Regresja logistyczna:** W odróżnieniu od regresji liniowej, która prognozuje wartości ciągłe, regresja logistyczna stosuje funkcję sigmoidalną (logistyczną), aby wyniki mieściły się w zakresie od 0 do 1, czyli odzwierciedlały prawdopodobieństwo wystąpienia danej klasy. Jest często stosowana w klasyfikacji binarnej, np. do diagnozy chorób (zdrowy vs chory), klasyfikacji wiadomości (spam vs nie-spam) czy przewidywania wyników (sukces vs porażka). [1]



Rysunek 3 - regresja logistyczna [1]

- **K-Najbliższych sąsiadów (KNN):** Jego zasada działania opiera się na założeniu, że obiekty o podobnych cechach będą należeć do podobnych klas lub mieć zbliżone wartości. W klasyfikacji KNN, aby przewidzieć klasę nowego punktu, algorytm wybiera k najbliższych sąsiadów tego punktu (na podstawie metryki odległości, np. euklidesowej) i przypisuje mu klasę, która występuje najczęściej w tej grupie sąsiadów. Algorytm KNN znajduje szerokie zastosowanie w różnych dziedzinach, szczególnie tam, gdzie ważne jest

dopasowanie nowych danych do wzorców na podstawie istniejących danych historycznych np. systemy rekomendacji, rozpoznawanie obrazów i wideo itd. [1]



Rysunek 4 - algorytm K-Najbliższych sąsiadów [1]

- **Naiwny klasyfikator Bayesa:** Algorytm probabilistyczny, który zakłada niezależność cech. Klasyfikatory Naiwnego Bayesa są często używane w klasyfikacji tekstów, ponieważ łatwo jest obliczyć prawdopodobieństwo na podstawie częstości, a tekst zwykle ma dużą liczbę cech (np. pojedyncze tokeny w słowach). [1]

Metody te wykorzystują wiedzę z fazy uczenia, aby skutecznie dokonywać predykcji w fazie wyznaczania wyniku, przypisując nowe dane do właściwych klas.

1.3 Metody stosowane do poprawności wiarygodności i dokładności budowanych modeli

Istnieje kilka kluczowych metod stosowanych w uczeniu maszynowym do poprawy wiarygodności i dokładności modeli. Te techniki pozwalają na ocenę skuteczności modelu i minimalizację błędów związanych z przeuczeniem lub niedouczeniem. Oto kilka z nich:

- **Walidacja krzyżowa z pozostawieniem jednej wartości (LOOCV):** W tej metodzie przeprowadzamy trening na całym zestawie danych, ale pozostawiamy tylko jeden punkt danych z dostępnego zestawu danych, a następnie iterujemy dla każdego punktu danych. W LOOCV model jest trenowany na $N-1$ próbki i testowane na jednej pominiętej próbce, powtarzając ten proces dla każdego punktu danych w zestawie danych. Wykorzystujemy wszystkie punkty danych, co sprawia, że jest ona mało obciążona błędem. Testujemy jednocześnie tylko jeden punkt danych. Jeśli punkt danych jest wartością odstającą, może to prowadzić do większej zmienności. Metoda ta zajmuje również dużo czasu wykonania, ponieważ iteruje się przez „liczbę punktów danych” razy. [4]
- **Walidacja krzyżowa K-Fold:** Polega na dzieleniu zbioru danych na k podzbiorów (znanych jako folds), a następnie przeprowadzeniu treningu na wszystkich podzbiórach, pozostawiając jeden $(k-1)$ podzbiór do oceny wytrenowanego modelu. Jej zaletą jest możliwość uzyskania bardziej wiarygodnych i stabilnych wyników oceny modelu, zwłaszcza na ograniczonym zbiorze danych [4]
- **Stratyfikowana walidacja krzyżowa:** Jest to technika stosowana w uczeniu maszynowym, aby zapewnić, że każdy etap procesu walidacji krzyżowej zachowuje ten sam rozkład klas, co cały zestaw danych. Jest to szczególnie ważne w przypadku nie zrównoważonych zestawów danych, w których pewne klasy mogą być niedoreprezentowane. W tej metodzie:
 1. Zbiór danych podzielony jest na k części, przy zachowaniu proporcji klas w każdej części.
 2. Podczas każdej iteracji jeden z elementów jest używany do testowania, a pozostałe do trenowania.

3. Proces powtarza się k razy, przy czym każdy podzbiór służy jako zbiór testowy dokładnie raz.

Stratyfikowana walidacja krzyżowa jest niezbędna w przypadku problemów klasyfikacji, w których zachowanie równowagi rozkładu klas ma kluczowe znaczenie dla prawidłowego uogólniania modelu na niewidziane wcześniej dane. [4]

- **Walidacja holdout:** W tej metodzie przeprowadzamy trening na 50% danego zestawu danych, a pozostałe 50% jest używane do celów testowych. To prosty i szybki sposób na ocenę modelu. Główną wadą tej metody jest to, że przeprowadzamy trening na 50% zestawu danych, może się zdarzyć, że pozostałe 50% danych zawiera pewne ważne informacje, które pomijamy podczas trenowania naszego modelu. [4]

2. Charakterystyka danych

2.1 Źródło i typ danych

Zbiór danych **Red Wine Quality Dataset** (Cortez et al., 2009) pochodzi z UCI Machine Learning Repository i jest używany do analiz w kontekście oceny jakości win czerwonych. Zawiera dane na temat fizykochemicznych właściwości win oraz ich ocenę jakościową [5].

Kolumny w zbiorze danych [5]:

1. **Kwasowość stała** (pH) w winie (liczba zmiennoprzecinkowa, np. 7.4)
2. **Kwasowość lotna** (pH) w winie (liczba zmiennoprzecinkowa, np. 0.7)
3. **Kwas cytrynowy** (w gramach na litr) (liczba zmiennoprzecinkowa, np. 0.0)
4. **Cukier resztkowy** (w gramach na litr) (liczba zmiennoprzecinkowa, np. 1.9)
5. **Zawartość chlorków** (w gramach na litr) (liczba zmiennoprzecinkowa, np. 0.076)
6. **Wolny dwutlenek siarki** (w miligramach na litr) (liczba zmiennoprzecinkowa, np. 11)
7. **Całkowity dwutlenek siarki** (w miligramach na litr) (liczba zmiennoprzecinkowa, np. 34)
8. **Gęstość** wina (liczba zmiennoprzecinkowa, np. 0.9978)
9. **pH** wina (liczba zmiennoprzecinkowa, np. 3.51)
10. **Zawartość siarczanów** (w gramach na litr) (liczba zmiennoprzecinkowa, np. 0.56)
11. **Zawartość alkoholu** (w procentach) (liczba zmiennoprzecinkowa, np. 9.4)
12. **Quality** – ocena jakości wina (skala od 0 do 10, liczba całkowita)

2.2 Informacje charakteryzujące dane

Ten zbiór danych zawiera 1599 rekordów oraz 12 kolumn. Wszystkie kolumny poza kolumną Quality zawierają informacje typu zmiennoprzecinkowego. Kolumna Quality zawiera dane typu całkowitego.

Za pomocą następującego programu wyznaczamy wartości minimalne, maksymalne, średnią, odchylenie standardowe oraz kwartyle [6].

```
import pandas as pd

#załadowanie danych
df = pd.read_csv('winequality-red.csv')
#generowanie podstawowych statystyk opisowych
statistics = df.describe()
# Zaokrąglenie statystyk do trzech miejsc po przecinku
statistics = statistics.round(3)
# Wyświetlenie wyników
print(statistics)
```

Sprawozdanie z przedmiotu Analiza danych w językach R i Python, KiIA PRz

Index	Kwasowość stała	Kwasowość lotna	Kwas cytrynowy	Cukier resztkowy	Zawartość chlorków	Wolny dwutlenek siarki
Liczba	1599	1599	1599	1599	1599	1599
Max	15.9	1.58	1	15.5	0.611	72
75%	9.2	0.64	0.42	2.6	0.09	21
Średnia	8.32	0.528	0.271	2.539	0.087	15.875
50%	7.9	0.52	0.26	2.2	0.079	14
25%	7.1	0.39	0.09	1.9	0.07	7
Min	4.6	0.12	0	0.9	0.012	1
std	1.741	0.179	0.195	1.41	0.047	10.46

Rysunek 5 Informacje charakteryzujące dane

Index	Całkowity dwutlenek siarki ▼	Gęstość	pH	Siarczany	Alkohol	Ocena Jakości
Liczba	1599	1599	1599	1599	1599	1599
max	289	1.004	4.01	2	14.9	8
75%	62	0.998	3.4	0.73	11.1	6
Średnia	46.468	0.997	3.311	0.658	10.423	5.636
50%	38	0.997	3.31	0.62	10.2	6
std	32.895	0.002	0.154	0.17	1.066	0.808
25%	22	0.996	3.21	0.55	9.5	5
min	6	0.99	2.74	0.33	8.4	3

Rysunek 6 Informacje charakteryzujące dane cd.

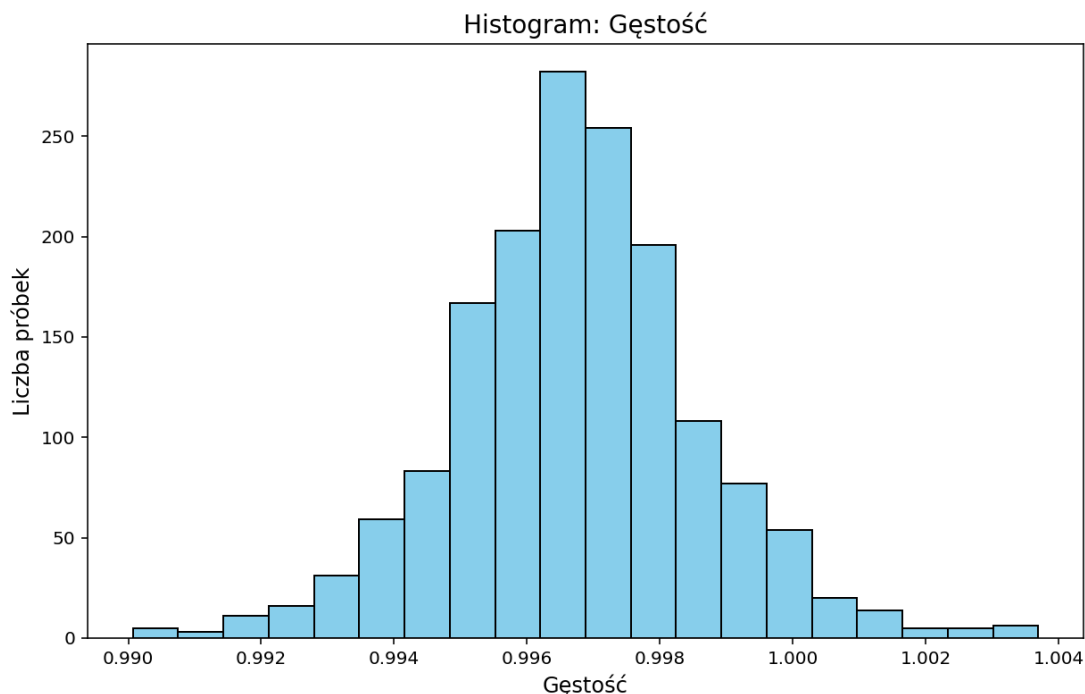
Za pomocą następującego programu wyznaczymy rozkład wartości każdej z danych [6,7].

```
import pandas as pd
import matplotlib.pyplot as plt

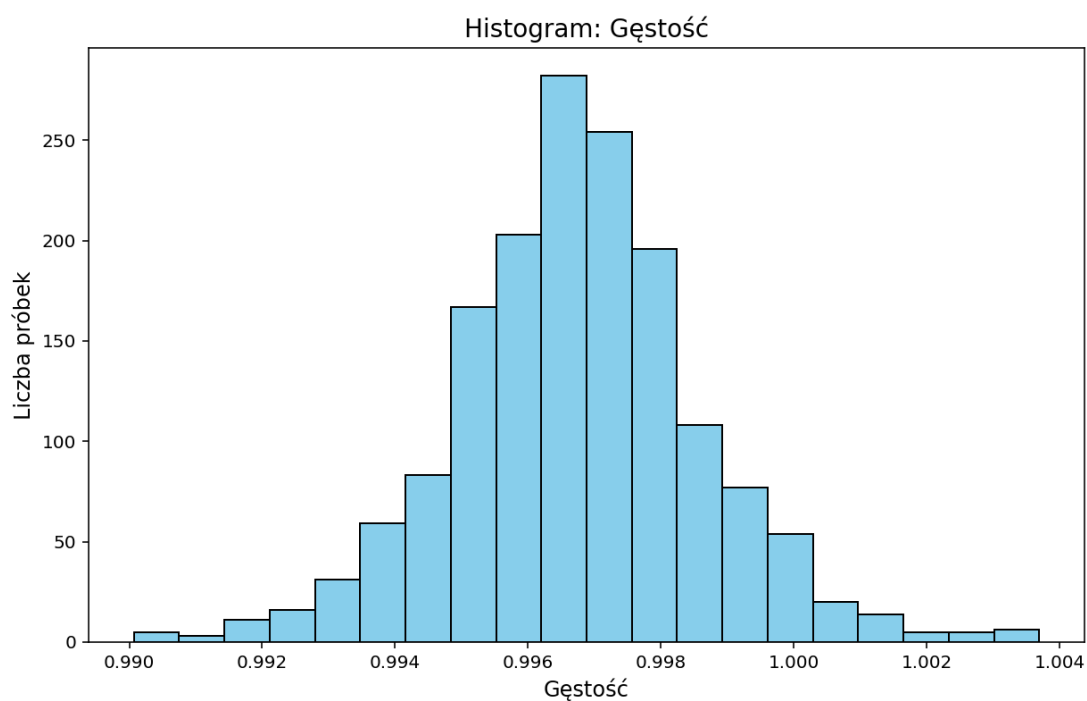
# Załadowanie danych
df = pd.read_csv('winequality-red.csv')
# Mapowanie nazw kolumn na polski
column_names_pl = {
    'fixed acidity': 'Kwasowość stała',
    'volatile acidity': 'Kwasowość lotna',
    'citric acid': 'Kwas cytrynowy',
    'residual sugar': 'Cukier resztkowy',
    'chlorides': 'Chlorki',
    'free sulfur dioxide': 'Wolny dwutlenek siarki',
    'total sulfur dioxide': 'Całkowity dwutlenek siarki',
```

```
'density': 'Gęstość',  
'pH': 'pH',  
'sulphates': 'Siarka',  
'alcohol': 'Alkohol',  
'quality': 'Jakość'  
}  
  
# Generowanie histogramów dla każdej kolumny w zbiorze danych  
for column in df.columns:  
    # Utworzenie wykresu  
    plt.figure(figsize=(10, 6))  
    plt.hist(df[column], bins=20, color='skyblue', edgecolor='black')  
    # Ustawienie tytułu wykresu i etykiet  
    # Użycie polskiego tłumaczenia lub oryginalnej nazwy  
    column_name_pl = column_names_pl.get(column, column)  
    plt.title(f'Histogram: {column_name_pl}', fontsize=14)  
    plt.xlabel(column_name_pl, fontsize=12)  
    plt.ylabel('Liczba próbek', fontsize=12)  
    # Zapisanie wykresu do pliku  
    filename = f'{column}_histogram.png'  
    plt.savefig(filename) # Zapisanie wykresu w formacie PNG  
    plt.close() # Zamknięcie wykresu po zapisaniu, aby przygotować nowy  
  
print("Wszystkie wykresy zostały zapisane.")
```

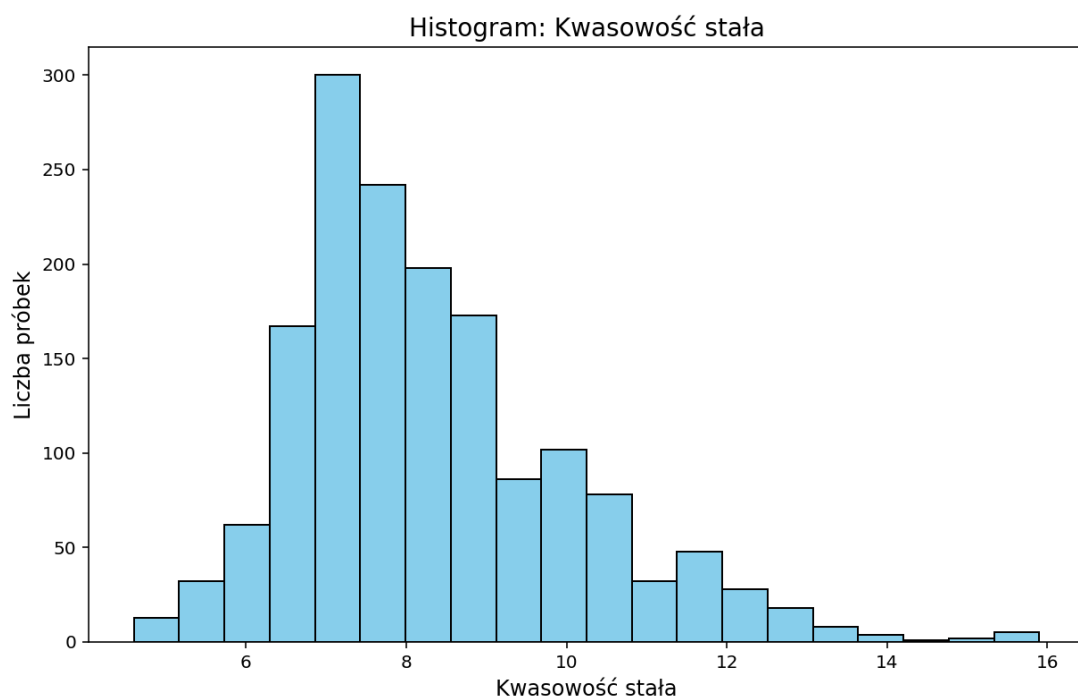
Na następnych zdjęciach zostaną przedstawione histogramy rozkładu wartości poszczególnych cech danych.



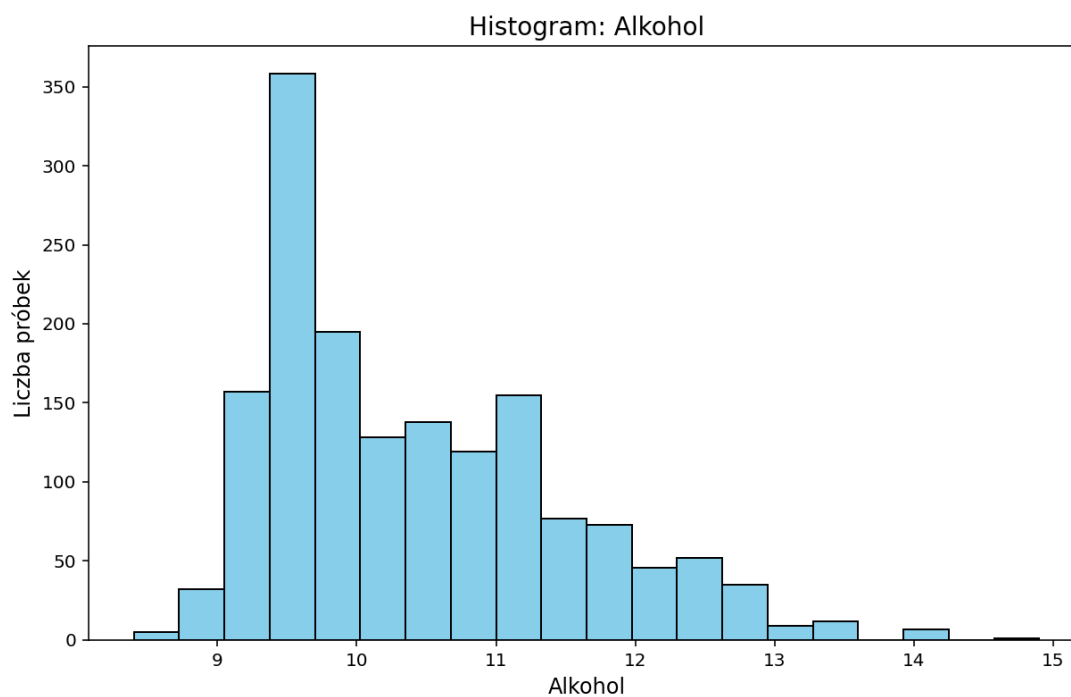
Rysunek 7 Histogram rozkładu danych gęstości win



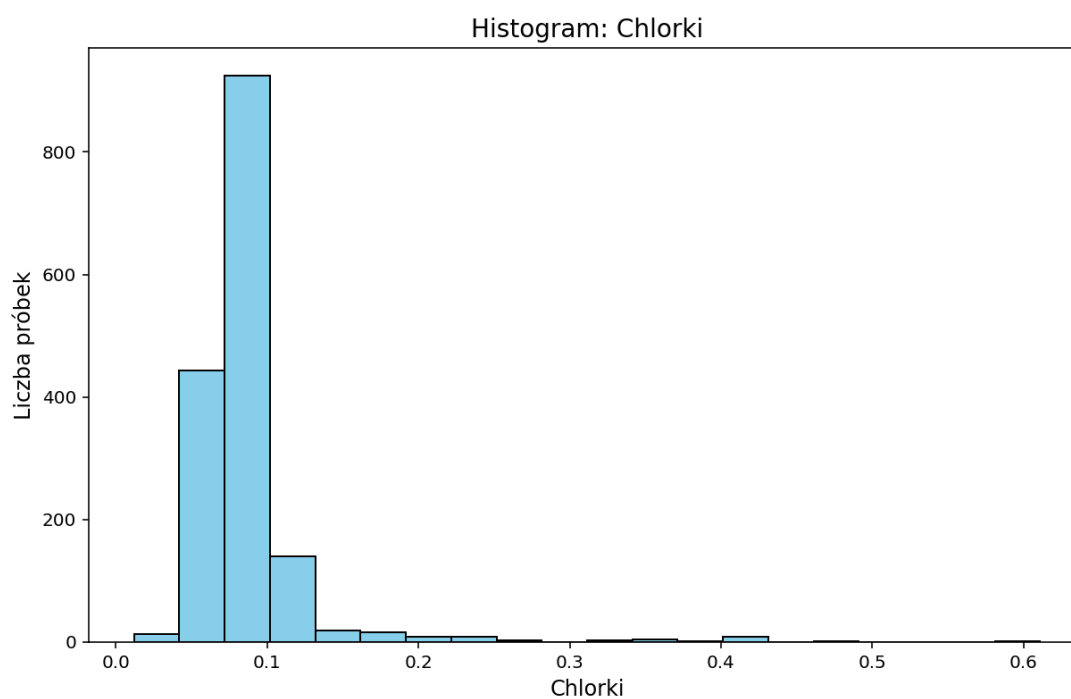
Rysunek 8 Histogram rozkładu danych kwasowości ulotnej



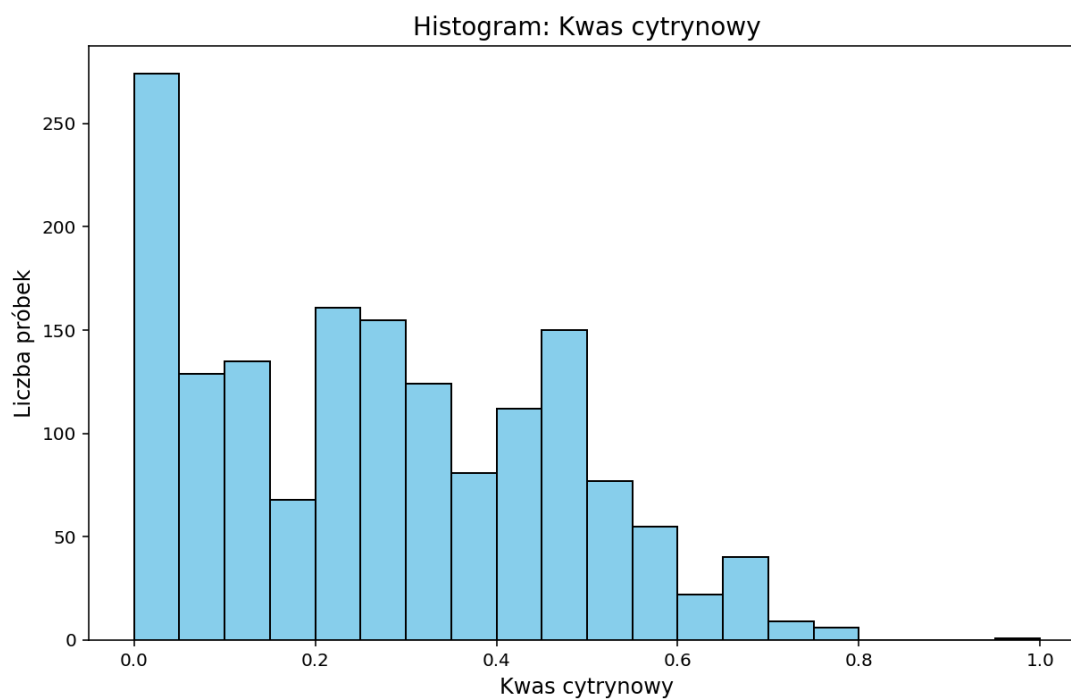
Rysunek 9 Histogram rozkładu danych kwasowości stałej



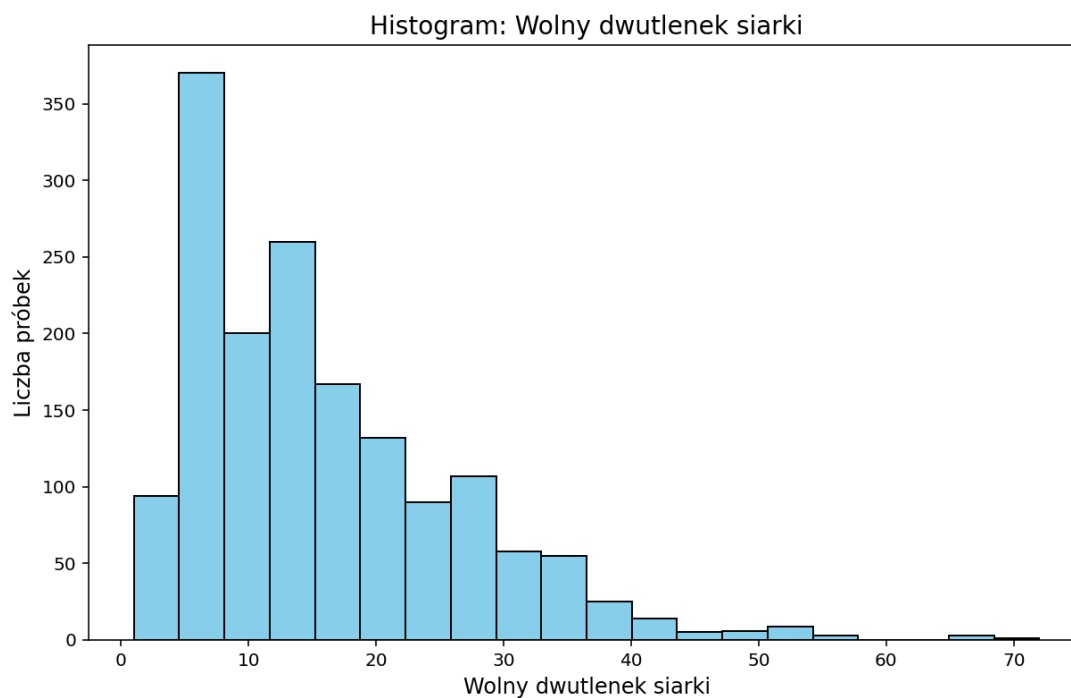
Rysunek 10 Histogram rozkładu zawartości alkoholu w winie



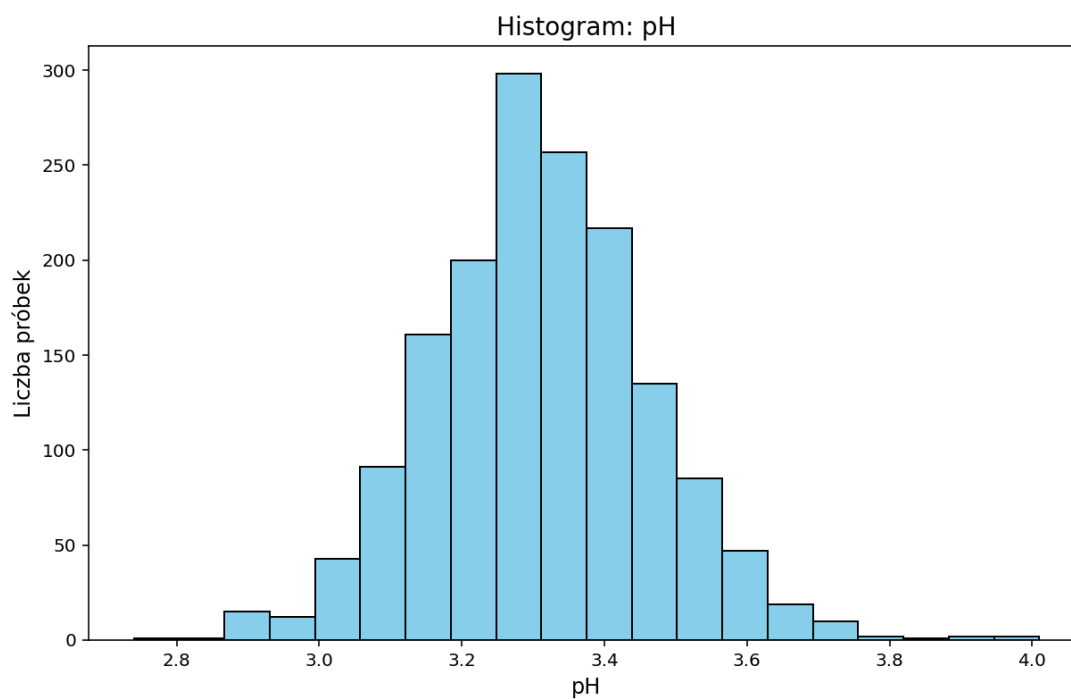
Rysunek 11 Histogram rozkładu zawartości chlorków



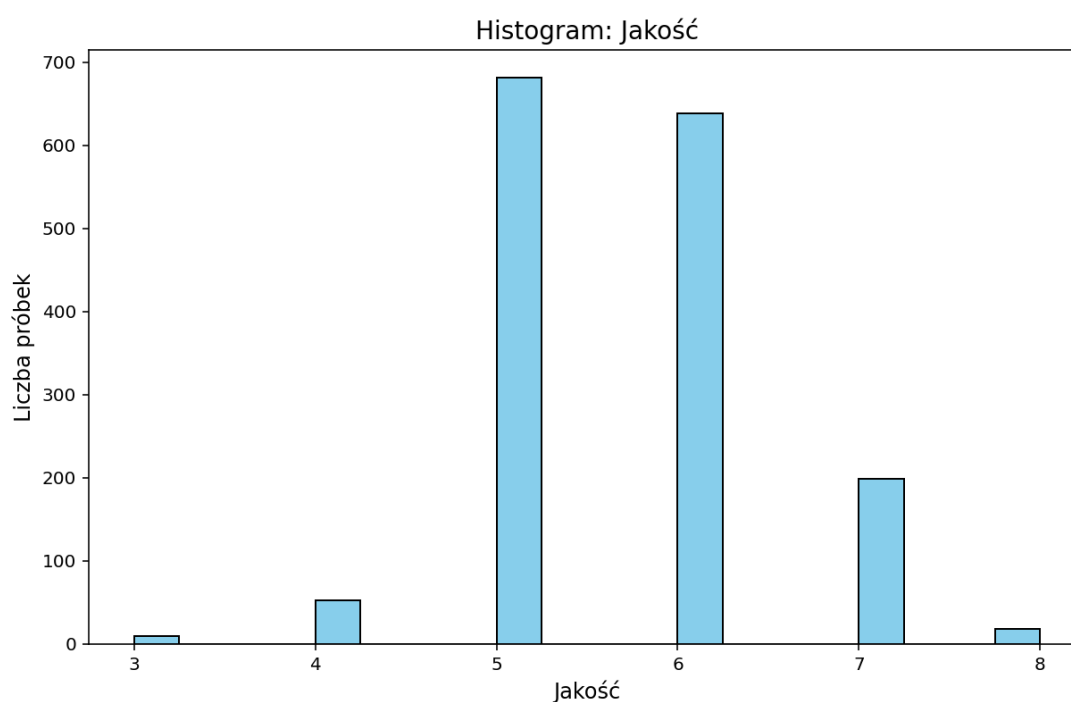
Rysunek 12 Histogram rozkładu zawartości kwasu cytrynowego



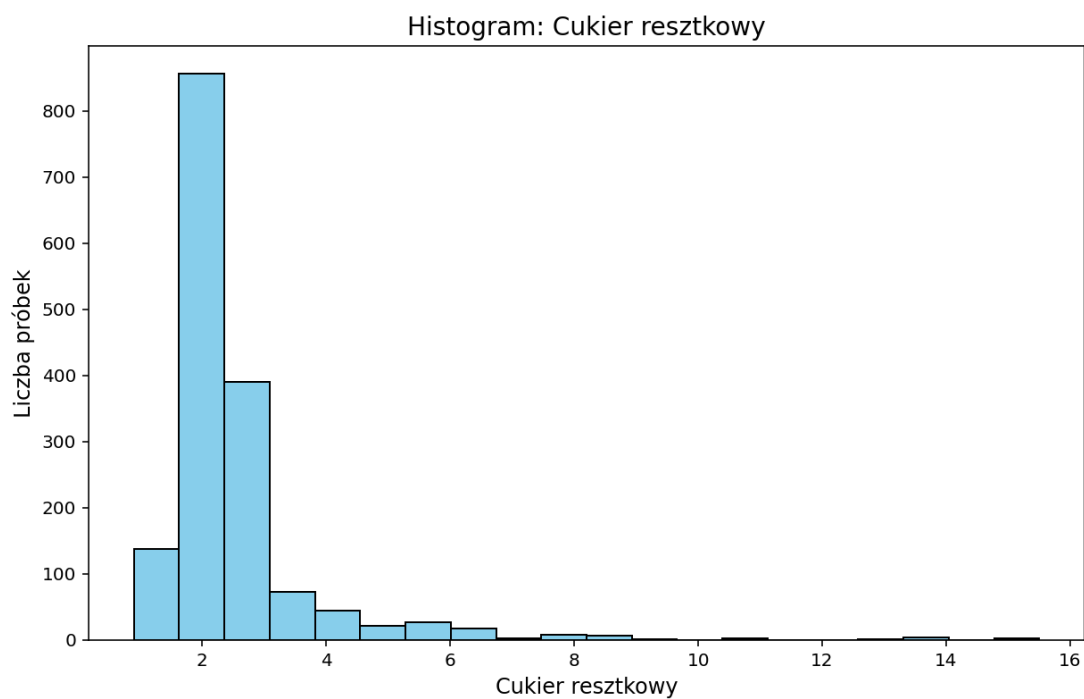
Rysunek 13 Histogram rozkładu zawartości wolnego dwutlenku siarki



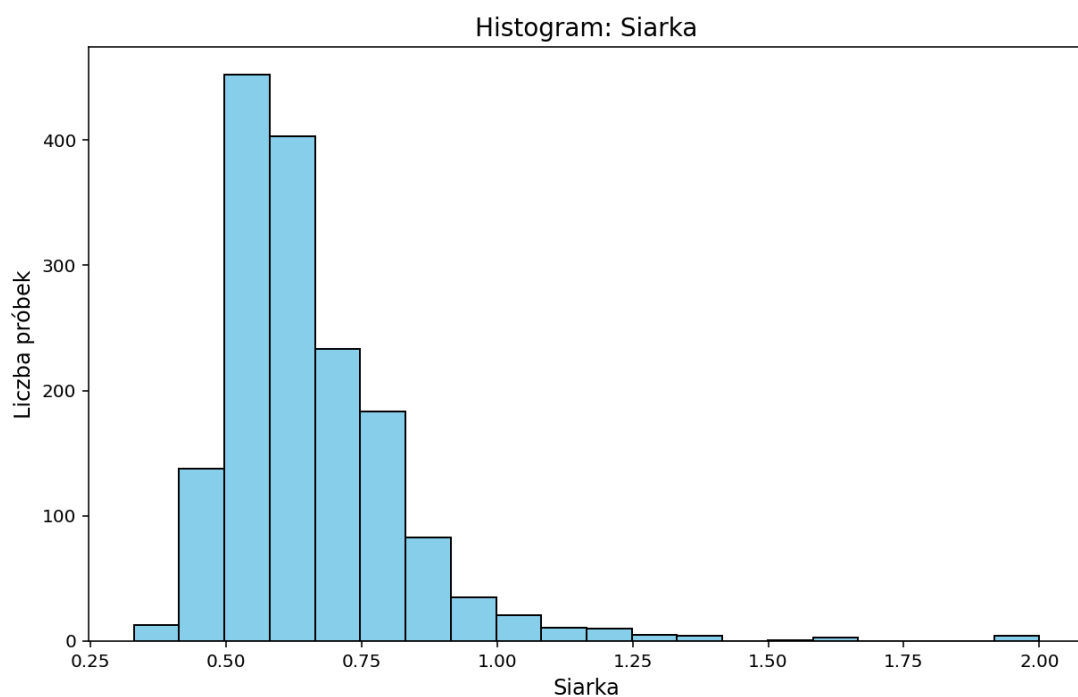
Rysunek 14 Histogram poziomu pH w winie



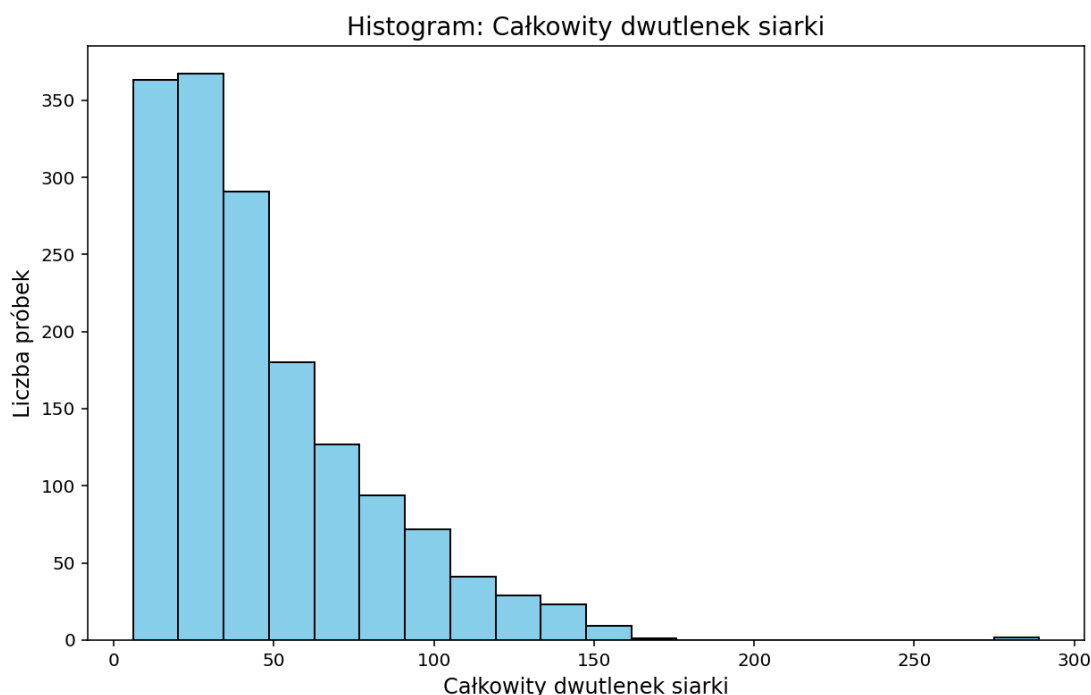
Rysunek 15 Histogram jakości wina



Rysunek 16 Histogram zawartości cukru reszkowego



Rysunek 17 Histogram rozkładu zawartości siarczanów w winie



Rysunek 18 Histogram rozkładu całkowitej zawartości dwutlenku siarki w winie

- **Kwasowość stała:** Średnia wynosi 8.32, przy czym rozkład jest zróżnicowany (odchylenie standardowe 1.741). Wartości wahają się od 4.6 do 15.9, co sugeruje szeroką zmienność tej cechy.
- **Kwasowość lotna:** Średnia to 0.528, z mniejszą zmiennością (std: 0.179). Maksymalne wartości sięgają 1.58, co może być uważane za odstające w stosunku do typowych wartości (między 0.39 a 0.64 w kwartylach).
- **Kwas cytrynowy:** Wartość średnia to 0.271, a odchylenie standardowe 0.195. Warto zauważyć, że minimalna wartość wynosi 0, co może oznaczać próbki wina, w których kwas cytrynowy nie występuje.
- **Cukier resztkowy:** Średnia wynosi 2.539, jednak maksymalna wartość aż 15.5 wskazuje na istnienie próbek o bardzo słodkich winach. Odchylenie standardowe 1.41 sugeruje znaczną różnorodność.
- **Chlorki:** Średnia wynosi 0.087, z niewielką zmiennością (std: 0.047). Maksymalna wartość (0.611) może być odstająca.
- **Wolny dwutlenek siarki:** Średnia to 15.875, przy dużym rozrzucie (std: 10.46). Wartości maksymalne (72.0) mogą być odstające w stosunku do typowych wartości w kwartylach (7.0–21.0).
- **Całkowity dwutlenek siarki:** Średnia wynosi 46.468, ale odchylenie standardowe aż 32.895 sugeruje, że niektóre próbki mogą zawierać bardzo wysokie ilości tej substancji (maksymalnie 289.0).
- **Gęstość:** Średnia wynosi 0.997, a zmienność jest niewielka (std: 0.002). Większość próbek mieści się w zakresie typowym dla win (0.996–0.998 w kwartylach).
- **pH:** Średnia wartość 3.311 wskazuje na lekko kwaśne pH, co jest typowe dla win. Rozkład jest stosunkowo wąski (std: 0.154).
- **Siarka:** Średnia wynosi 0.658, z maksymalnymi wartościami sięgającymi 2.0. Wyższe wartości mogą być odstającymi.
- **Alkohol:** Średnia to 10.423%, przy czym rozkład jest umiarkowany (std: 1.066). Minimalne i maksymalne

Sprawozdanie z przedmiotu Analiza danych w językach R i Python, KiIA PRz

wartości (8.4–14.9) pokazują różnorodność w zawartości alkoholu.

- **Jakość:** Średnia wynosi 5.636, z odchyleniem standardowym 0.808. Wartości wahają się od 3 do 8, co sugeruje skupienie wokół średniej z niewielką liczbą bardzo niskich i bardzo wysokich ocen.

Dodatkowo dla następujących kolumn: **total sulfur dioxide** oraz osobno: **residual sugar, chlorides, fixed acidity, volatile acidity, citric acid, sulphates** wykonamy wykres boxowy aby lepiej zobrazować te dane. Wykorzystamy do tego następujący program[6,7]:

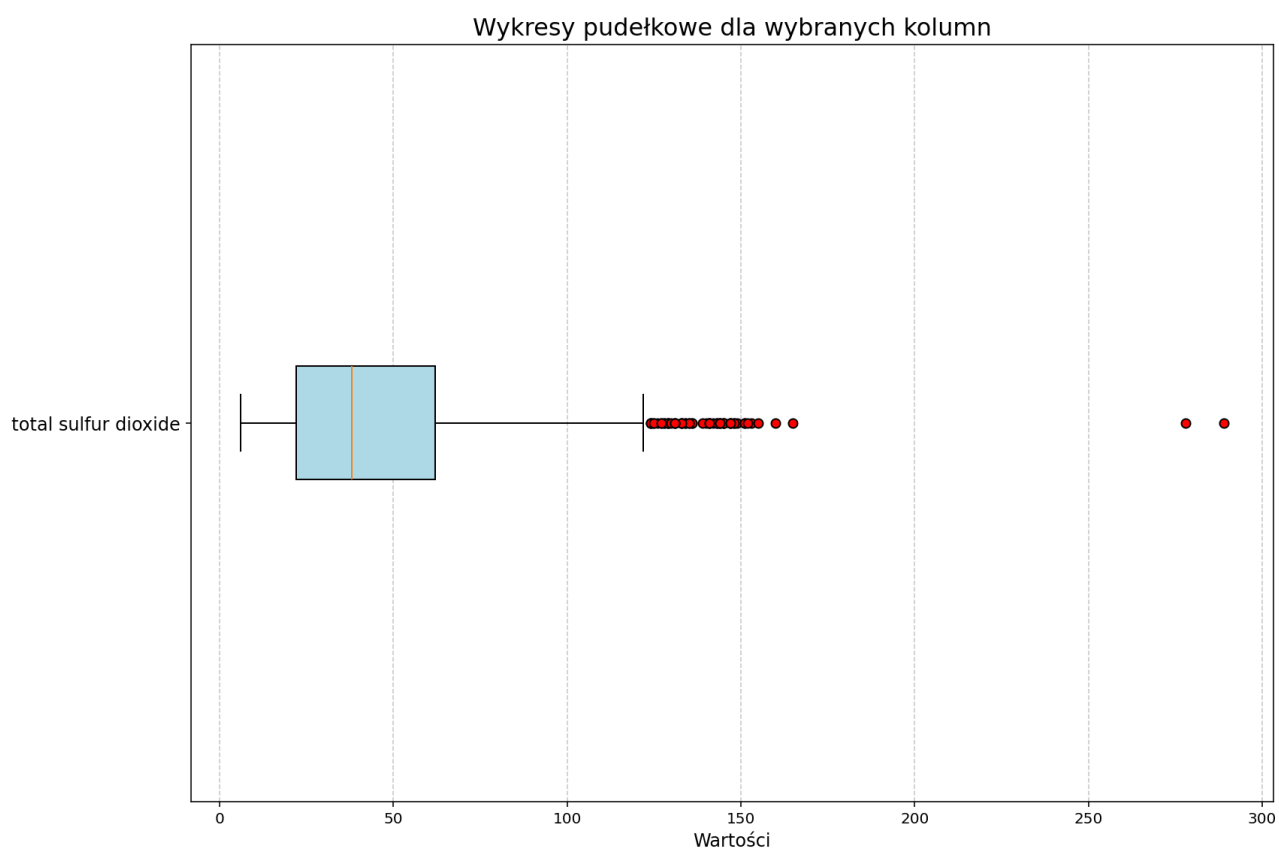
```
import pandas as pd
import matplotlib.pyplot as plt

# Załadowanie danych
df = pd.read_csv('winequality-red.csv')

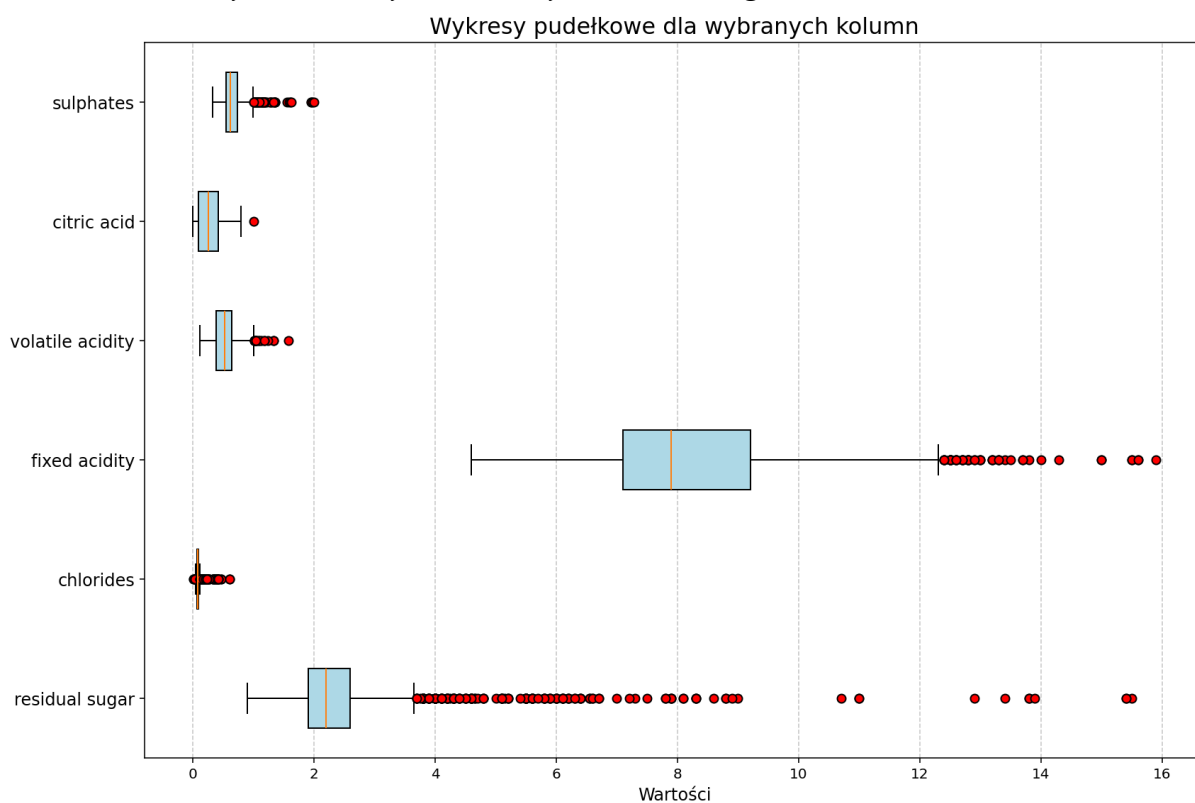
# Wybrane kolumny do analizy
columns_to_plot = [
    'total sulfur dioxide',
    'residual sugar',
    'chlorides',
    'fixed acidity',
    'volatile acidity',
    'citric acid',
    'sulphates'
]

# Tworzenie jednego wykresu boxplot
plt.figure(figsize=(12, 8))
plt.boxplot([df[column] for column in columns_to_plot],
            vert=False, # Poziome wykresy pudełkowe
            patch_artist=True,
            boxprops=dict(facecolor='lightblue', color='black'),
            whiskerprops=dict(color='black'),
            capprops=dict(color='black'),
            flierprops=dict(markerfacecolor='red', marker='o'))
plt.xticks(range(1, len(columns_to_plot) + 1), columns_to_plot, fontsize=12) # Etykiety osi Y
plt.title('Wykresy pudełkowe dla wybranych kolumn', fontsize=16)
plt.xlabel('Wartości', fontsize=12)
plt.grid(axis='x', linestyle='--', alpha=0.7)
plt.tight_layout()

# Zapis wykresu do pliku
plt.savefig('boxplot_all_columns.png')
plt.show()
```



Rysunek 19 Wykres boxowy dla całkowitego dwutlenku siarki



Rysunek 20 Wykresy boxowe dla reszty wybranych danych

2.3 Przygotowanie danych do budowy modelu klasyfikacji

Aby móc skutecznie zbudować model klasyfikacji, dane należy wyczyścić, znormalizować i usunąć wartości puste.

Dzięki użyciu następującego kodu w języku Python, ładujemy zbiór danych a **następnie sprawdzamy, czy występują wartości puste** [6].

```
import pandas as pd

df = pd.read_csv('winequality-red.csv')

print(df.isnull().values.any())
```

Otrzymaliśmy *False*, więc ten zbiór danych nie zawiera wartości pustych. Następnym krokiem będzie **usunięcie duplikatów**. Możemy to wykonać następującym programem [6].

```
import pandas as pd

# Załadowanie danych
df = pd.read_csv('winequality-red.csv')

# Sprawdzenie, czy dane zawierają duplikaty
if df.duplicated().any():
    print("Znaleziono duplikaty w danych.")

    # Wyświetlenie liczby duplikatów
    print(f"Liczba duplikatów: {df.duplicated().sum()}")

    # Usunięcie duplikatów
    df = df.drop_duplicates()
    print("Duplikaty zostały usunięte.")
else:
    print("Brak duplikatów w danych.")

# Zapisanie oczyszczonego zbioru danych do nowego pliku
df.to_csv('cleaned_winequality-red.csv', index=False)
print("Oczyszczone dane zapisano w pliku 'cleaned_winequality-red.csv'.")
```

Po uruchomieniu programu, okazało się, że znalezione zostały 240 duplikaty. Zostały one usunięte a oczyszczony zbiór danych został zapisany do nowego pliku.

Sprawozdanie z przedmiotu Analiza danych w językach R i Python, KiIA PRz

Kolejnym krokiem będzie usunięcie danych odstających. Najwięcej danych odstających zauważyliśmy w kolumnach: **total sulfur dioxide**, **residual sugar**, **chlorides**, **fixed acidity**, **volatile acidity**, **citric acid**, **sulphates**. Użyjemy do tego rozstępu międzykwartylowego (IQR). Wykorzystamy do tego następujący program [6].

```
import pandas as pd

# Załadowanie danych
df = pd.read_csv('cleaned_winequality-red.csv')

# Funkcja do usuwania odstających wartości na podstawie IQR
def remove_outliers_iqr(df, column):
    q1 = df[column].quantile(0.25) # 1. kwartyl
    q3 = df[column].quantile(0.75) # 3. kwartyl
    iqr = q3 - q1 # rozstęp międzykwartylowy
    lower_bound = q1 - 1.5 * iqr # dolna granica
    upper_bound = q3 + 1.5 * iqr # górna granica
    return df[(df[column] >= lower_bound) & (df[column] <= upper_bound)]

# Usuwanie wartości odstających dla wybranych kolumn
columns_to_check = [
    'total sulfur dioxide',
    'residual sugar',
    'chlorides',
    'fixed acidity',
    'volatile acidity',
    'citric acid',
    'sulphates'
]

for column in columns_to_check:
    df = remove_outliers_iqr(df, column)

# Sprawdzenie efektów
print(df.describe())
# Zapisanie przetworzonych danych do nowego pliku CSV
df.to_csv('cleaned_data.csv', index=False)
```

2.4 Normalizacja danych

Użyjemy normalizacji Min-Max, która zmienia zakres wartości od 0 do 1. W przypadku tych danych jest to o tyle przydatne, że mamy różne jednostki danych, a potrzebujemy je znormalizować w celu pozytywnej dalszej analizy i budowy modeli. Program normalizacji metodą Min-Max [6].

```
import pandas as pd

# Załadowanie danych
df = pd.read_csv('cleaned_data.csv')

# Wybór tylko kolumn liczbowych
numeric_columns = df.select_dtypes(include=['float64', 'int64']).columns

# Normalizacja Min-Max za pomocą Pandas
df[numeric_columns] = df[numeric_columns].apply(lambda x: (x - x.min()) / (x.max() - x.min()))

# Sprawdzenie wyników
print(df.describe())

# Zapisanie znormalizowanych danych do nowego pliku CSV
df.to_csv('normalized_cleaned_data.csv', index=False)
```

3. Opis eksperymentów i opracowanie uzyskanych wyników

Ten rozdział powinien zawierać wyniki eksperymentów wykonanych dla pozyskanych danych za pomocą przynajmniej 3 wybranych metod klasyfikacji. Powinny zostać przedstawione wyniki budowy modeli klasyfikacji dla różnych parametrów wybranych algorytmów oraz porównanie uzyskanych wyników w formie tabelarycznej i graficznej. Kolejne prezentowane eksperymenty powinny prowadzić do poprawienia parametrów algorytmów użytych w następnych eksperymentach. Uzyskane wyniki oraz dyskusję ich poprawności, dokładności, czasu uzyskania oraz parametrów i ich wpływu na wyniki należy umieścić w tym rozdziale. Całość należy udokumentować fragmentami kodu w języku Python lub R.

Podsumowanie

Tutaj należy umieścić podsumowanie wykonanych prac, własne wnioski dotyczące rozwiązywanego problemu klasyfikacji, metod oraz ich parametrów, zastosowanych narzędzi.

Bibliografia

- [1] Akkio.com. (2024, 1 8). Pobrano z lokalizacji <https://www.akkio.com/post/5-types-of-machine-learning-classification-algorithms>
- [2] Fatyga Piotr, P. R. (2023, 11 14). Klasyfikacja danych - przegląd wybranych metod, .
- [3] geekforgeeks.org. (2024, 8 7). Pobrano z lokalizacji <https://www.geeksforgeeks.org/artificial-neural-networks-and-its-applications/>
- [4] geeksforggeeks.org. (2024, 8 7). Pobrano z lokalizacji <https://www.geeksforggeeks.org/cross-validation-machine-learning/>
- [5] <https://www.kaggle.com/datasets/uciml/red-wine-quality-cortez-et-al-2009>
- [6] https://pandas.pydata.org/docs/user_guide/index.html#user-guide
- [7] https://matplotlib.org/3.5.3/api/_as_gen/matplotlib.pyplot.html