

Analiza danych - Python

Skład: Robert Baca, Arkadiusz Bodziony, Wiktor Ciskał, Mikołaj Śnieżko

Praca wykonywana w okresie II kwartału roku 2024

1 Wprowadzenie

W ramach tego projektu zajmujemy się analizą danych szkolnych, w tym ocen i frekwencji uczniów z różnych regionów. Wykonamy różne analizy statystyczne i przedstawimy je za pomocą wykresów.

2 Przygotowanie danych

- Struktura pliku danych: `address`, `absences`, `Mjob`, `Fjob`, `math_grade`
- Kolumny `address` zawierają wartości U (miejski) i R (wiejski).
- Kolumna `absences` zawiera liczbę nieobecności ucznia.
- Kolumny `Mjob` i `Fjob` oznaczają zawody matki i ojca.
- Kolumna `math_grade` oznacza końcową ocenę z matematyki w skali 0-20.

3 Średnie oceny

```
1 def srednia_ocen(adr, oceny):
2     licznik_w = 0
3     licznik_m = 0
4     sum_wiejski = 0
5     sum_miejski = 0
6
7     for i, j in zip(oceny, adr):
8         if j == 'r':
9             sum_wiejski += i
10            licznik_w += 1
11        elif j == 'u':
12            sum_miejski += i
13            licznik_m += 1
14
15        srednia_wiejski = sum_wiejski / licznik_w if licznik_w != 0
16        else 0
17        srednia_miejski = sum_miejski / licznik_m if licznik_m != 0
18        else 0
19
20    return round(srednia_miejski, 2), round(srednia_wiejski, 2)
```

Średnia ocen dla miasta wynosi 10.69, a dla wsi 9.47. Z tego wynika, że uczniowie z miasta mają statystycznie lepsze stopnie w szkole.

4 Mediana ocen

```
1 def mediana_ocen(adr, oceny):
2     wiejski = [i for i, j in zip(oceny, adr) if j == 'r']
3     miejski = [i for i, j in zip(oceny, adr) if j == 'u']
4
5     return median(oceny), median(miejski), median(wiejski)
```

Mediana wszystkich ocen wynosi 11.0, dla miasta 11.0, dla wsi 11.0. Zatem mediana ocen jest identyczna dla obu kategorii.

5 Odchylenie standardowe

```
1 def sigma(x):
2     return sqrt((sum([(i-mean(x))**2 for i in x])/len(x)))

1 def odchyl_std(adr, oceny):
2     wiejski = [i for i, j in zip(oceny, adr) if j == 'r']
3     miejski = [i for i, j in zip(oceny, adr) if j == 'u']
4
5     return sigma(oceny), sigma(miejski), sigma(wiejski)
```

Odchylenie standardowe wszystkich ocen wynosi 4.58, dla miasta 11.0, dla wsi 10.0. Te pomiary mogą nas naprowadzić na tezę, że oceny uczniów z miasta osiągają większe amplitudy.

6 Moda

```
1 def obliczanie_mody(adr, oceny):
2     wiejski = [i for i, j in zip(oceny, adr) if j == 'r']
3     miejski = [i for i, j in zip(oceny, adr) if j == 'u']
4
5     return moda(oceny), moda(miejski), moda(wiejski)

1 def moda(liczby):
2     najczesciej = {}
3     for i in liczby:
4         if i in najczesciej:
5             najczesciej[i] += 1
6         else:
7             najczesciej[i] = 1
8
9     max_licznik = 0
10    najczesciej_liczba = 0
11
12    for wartosc, licznik in najczesciej.items():
13        if licznik > max_licznik:
14            max_licznik = licznik
15            najczesciej_liczba = wartosc
16
17    return najczesciej_liczba
```

Moda dla wszystkich ocen wynosi 10, dla miasta 10, dla wsi 10. Podobnie jak w przypadku mediany, moda jest równa w obu grupach.

7 Regresja liniowa i współczynnik R^2

```
1 def korelacja_regresja_r2(nobec, oceny):
2     return korelacja(nobec, oceny), reg liniowa(nobec, oceny), r2
       (nobec, oceny)
```

```
1 def korelacja(x, y):
2     return sum([(i-mean(x))*(j-mean(y)) for i, j in zip(x, y)]) /
3     sqrt(sum([(i-mean(x))**2 for i in x]) * sum([(i-mean(y))**2
       for i in y]))
```

```
1 def reg liniowa(x, y):
2     a1 = sum([(i-mean(x))*(j-mean(y)) for i, j in zip(x, y)]) /
3     sum([(i-mean(x))**2 for i in x])
4     a0 = mean(y-(a1*mean(x)))
5     return f'{round(a0, 2)}x + {round(a1, 2)}'
```

```
1 def r2(x, y):
2     return 1-(sum([(i-mean(y))**2 for i in y]) / sum([(j-i)**2
3     for i, j in zip(x, y)]))
```

Współczynnik korelacji między nieobecnością a oceną wynosi 0.03, regresja liniowa $10.3x + 0.02$, a współczynnik R^2 jest równy 0.08. Korelacja osiąga umiarkowaną wartość, więc oceny niekoniecznie są powiązane z frekwencją, jednak nie można uznać tego powiązania za nieistniejące.

8 Przewidywanie punktów

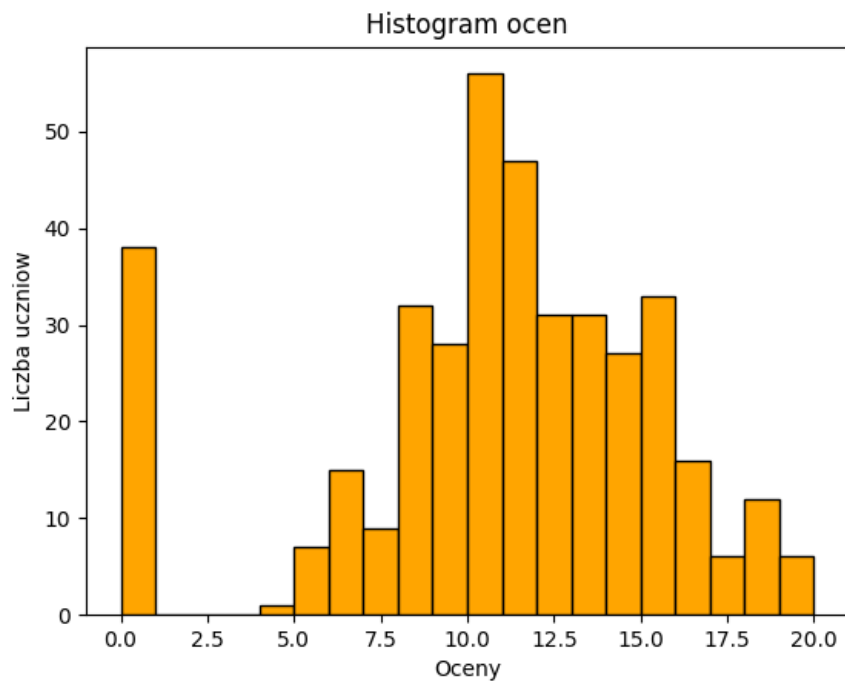
```
1 def punkty_predykcja(nieobecnosci, grupy):
2     liczba_zajec = 75
3     punkty = []
4     for liczba_nieobecnosci, grupa in zip(nieobecnosci, grupy):
5         liczba_obecnosci = (liczba_zajec - liczba_nieobecnosci)
6         punkty_obecnosci = liczba_obecnosci * 5
7         if liczba_nieobecnosci == 0:
8             punkty_obecnosci += 35
9         if grupa == 'r':
10             punkty_obecnosci += 2 * liczba_obecnosci
11         punkty.append(punkty_obecnosci)
12
13     return punkty
```

Punkty za obecności zostały przedstawione na wykresie w sekcji 'Wizualizacje'. Wypisywanie całej listy, zajęłoby sporo przestrzeni, a nie są to dane niezbędne. Zamieszczone jednak zostaną pierwsze 10 wyników: [345, 497, 325, 365, 355, 325, 410, 345, 410, 410].

9 Wizualizacje

9.1 Histogram ocen

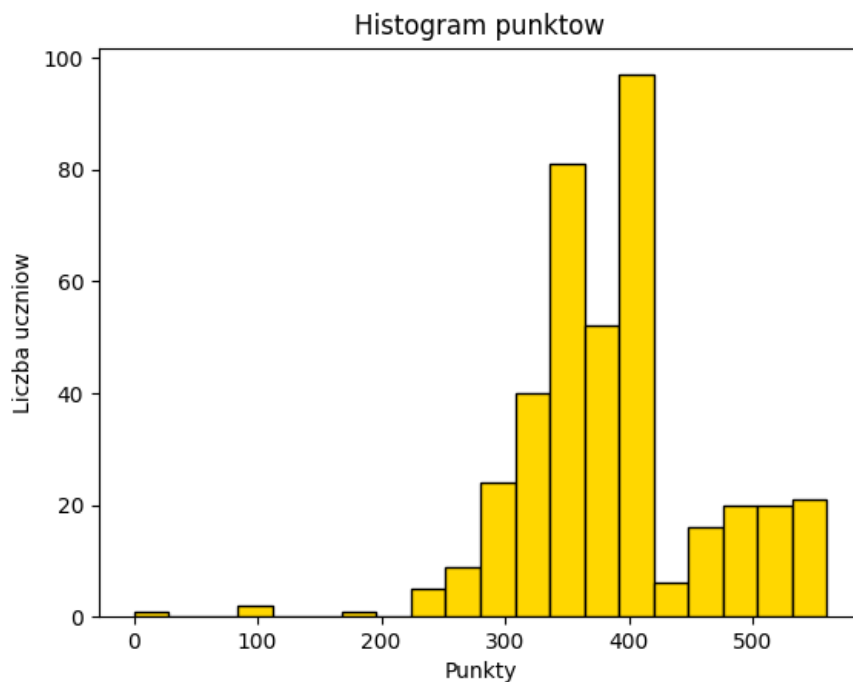
```
1 def histogram_ocen(oceny):  
2     plt.hist(oceny, bins=20, edgecolor='black', color='orange')  
3     plt.xlabel('Oceny')  
4     plt.ylabel('Liczba uczniów')  
5     plt.title('Histogram ocen')  
6     plt.show()
```



Krótki wniosek: Zdecydowana większość uczniów na koniec roku otrzymała oceny w przedziale [10-12]. Są to wyniki graniczące z warunkami zaliczenia, więc stwierdzić można, że przeważającą grupą uczniów są uczniowie przeciętni. Jednak warto zauważyć, jaka dysproporcja istnieje między ocenami 0, a [17-20]. Prawdopodobnie wynika to z faktu, że oceny 0 wystawiane mogły być za nieoddane projekty, natomiast najwyższe stopnie, otrzymali jedynie wybitni uczniowie. Tezę o ocenach zerowych potwierdzić może brak uczniów w przedziale [1-4].

9.2 Histogram punktów

```
1 def histogram_punktow(oceny):  
2     plt.hist(oceny, bins=20, edgecolor='black', color='gold')  
3     plt.xlabel('Punkty')  
4     plt.ylabel('Liczba uczniów')  
5     plt.title('Histogram punktów')  
6     plt.show()
```



Krótki wniosek: Z wykresu odczytać można, że przeważającym przedziałem jest [330-425]. Wynika to z danych otrzymanych ze szkoły oraz wzoru na wyliczenie przewidywanych punktów. Wzór: liczba obecności*5 + dodatkowe 2 za każdą obecność dla uczniów ze wsi. Większość uczniów opuściło jedynie kilka zajęć, co można odczytać z wykresu. Nieliczni nie opuścili żadnych zajęć, a osób, które nie pojawiły się wcale było bardzo mało.

9.3 Udział grup w klasie

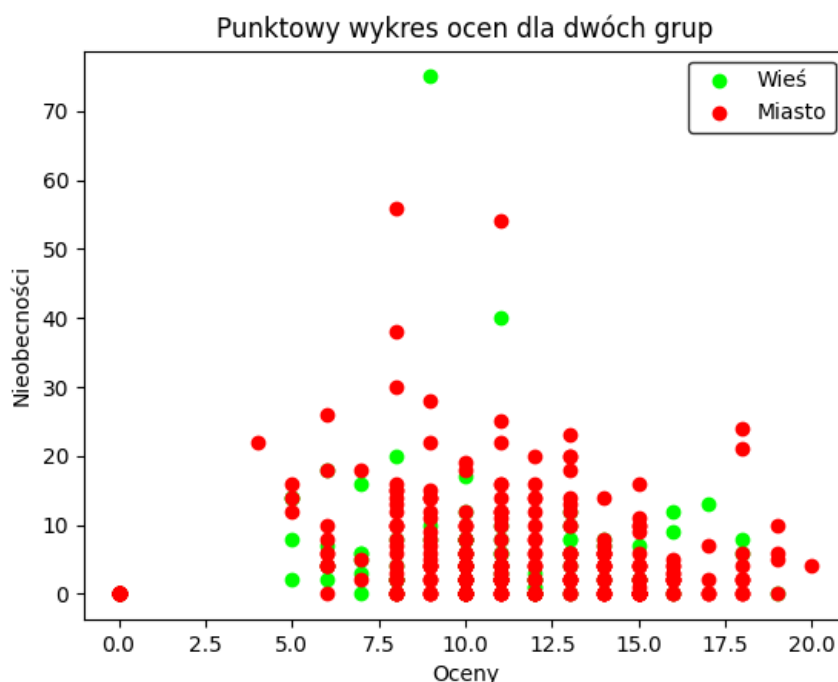
```
1 def udzial_grup(oceny, adr):
2     wiejski = [i for i, j in zip(oceny, adr) if j == 'r']
3     miejski = [i for i, j in zip(oceny, adr) if j == 'u']
4
5     plt.pie([len(wiejski)/len(oceny), len(miejski)/len(oceny)],
6             labels=
7             ['obszar wiejski', 'obszar miejski'], autopct='%0.1f%%',
8             colors=['lime', 'red'])
9     plt.legend(['obszar wiejski', 'obszar miejski'], edgecolor='
10                black', loc='upper left')
```



Krótki wniosek: W przypadku tego wykresu, ciężko o jakiegokolwiek wnioski poza tymi oczywistymi. Ponad 3/4 uczniów z badanej szkoły zamieszkuje tereny miejskie.

9.4 Wykres punktowy ocen uczniów z obu grup

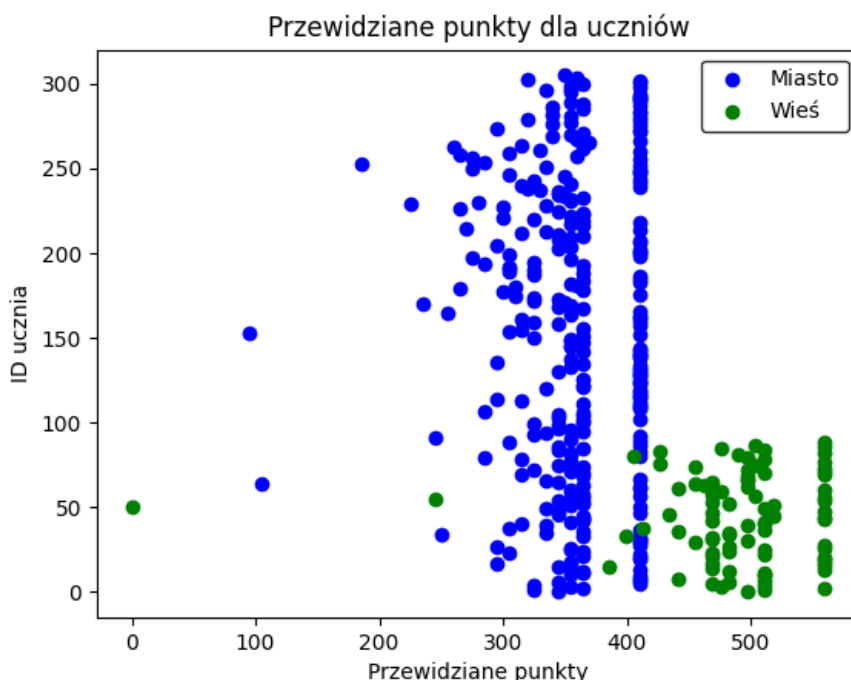
```
1 def punktowy(oceny, adr, nobec):
2     o_wiejski = [i for i, j in zip(oceny, adr) if j == 'r']
3     o_miejski = [i for i, j in zip(oceny, adr) if j == 'u']
4
5     n_wiejski = [i for i, j in zip(nobec, adr) if j == 'r']
6     n_miejski = [i for i, j in zip(nobec, adr) if j == 'u']
7
8     plt.scatter(o_wiejski, n_wiejski, color='lime')
9     plt.scatter(o_miejski, n_miejski, color='red')
10
11     plt.title('Punktowy wykres ocen dla dwóch grup')
12     plt.ylabel('Nieobecności')
13     plt.xlabel('Oceny')
14     plt.legend(['Wieś', 'Miasto'], edgecolor='black')
15     plt.show()
```



Krótki wniosek: Ten wykres może być nieco niezrozumiały na pierwszy rzut oka, dlatego spieszmy z pomocą. Zobrazowane tutaj są oceny w skali liczby nieobecności. Widać, że oceny nie są w dużym stopniu skorelowane z ocenami, dane są dość chaotyczne i niepowiązane. Możemy stwierdzić, że pewna korelacja istnieje, jednak zdecydowanie nie jest ona duża. Zdecydowana większość uczniów, była nieobecnych [0-20] razy i osiągają oni oceny [7-13]. Zauważyć możemy, że stosunkowo bardziej odstającymi w obie strony są uczniowie ze wsi, którzy pomimo niewielu nieobecności osiągają stosunkowo niskie i stosunkowo wysokie stopnie. Jednak, prawdopodobnie wynika to z próbki, która jest znacznie mniejsza od studentów z miasta.

9.5 Wykres przewidywanych punktów uczniów z obu grup

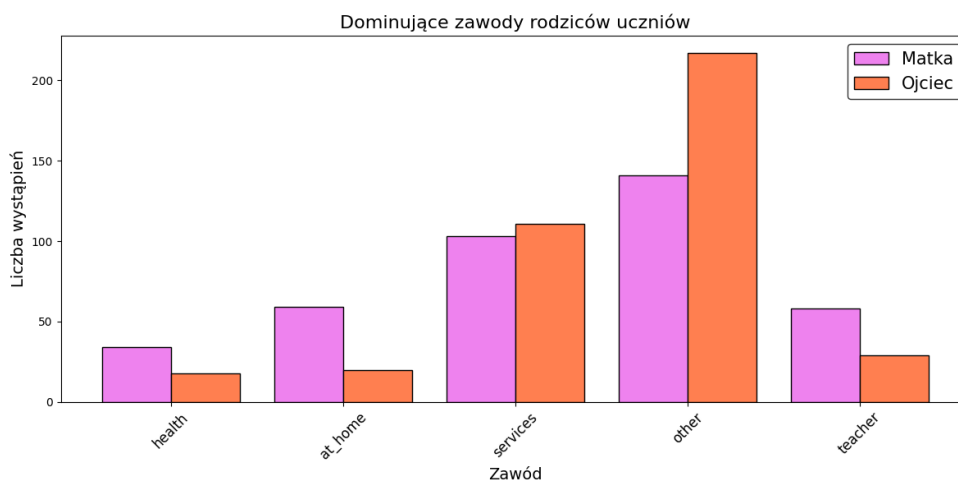
```
1 def wykres_punktowy_predykcji(predykcja, grupa):
2     punkty_wsi = [predykcja[i] for i in range(len(predykcja)) if
3                   grupa[i] == 'r']
4     punkty_miasta = [predykcja[i] for i in range(len(predykcja))
5                      if grupa[i] == 'u']
6
7     plt.scatter(punkty_miasta, [i for i in range(len(
8         punkty_miasta))], color='blue',
9                 label='Miasto')
10    plt.scatter(punkty_wsi, [i for i in range(len(punkty_wsi))],
11               color='green',
12               label='Wieś')
13
14    plt.xlabel('Przewidziane punkty')
15    plt.ylabel('ID ucznia')
16    plt.title('Przewidziane punkty dla uczniów')
17    plt.legend(edgecolor='black')
18    plt.show()
```



Krótki wniosek: Ten wykres również może wydawać się niezrozumiały, jednak ponownie - spokojnie, zaraz będzie on klarowniejszy. Patrząc na legendę, widzimy jakie punkty, co oznaczają. Jest to wykres przedstawiający wszystkich uczniów oraz ich punkty za frekwencję. Widać, że najwięcej uczniów osiągnie ok. 410 punktów. Bliżej maksymalnej wartości są uczniowie ze wsi, co jest dość oczywiste, biorąc pod uwagę nasz wzór na przewidywane punkty. Nie licząc pojedynczych przypadków, widzimy, że uczniowie ze wsi, pojawią się w placówce podobną liczbę razy, jednakże trzeba wziąć poprawkę na to, że mają znacznie cięższy dojazd.

9.6 Zawody wśród rodziców uczniów

```
1 def zawody_rodzicow(m_zawod, o_zawod):
2     job_counts_mother = {}
3     job_counts_father = {}
4
5     for job in m_zawod:
6         if job not in job_counts_mother:
7             job_counts_mother[job] = 0
8             job_counts_mother[job] += 1
9
10    for job in o_zawod:
11        if job not in job_counts_father:
12            job_counts_father[job] = 0
13            job_counts_father[job] += 1
14
15    jobs = list(set(m_zawod + o_zawod))
16    counts_mother = [job_counts_mother.get(job, 0) for job in
17                     jobs]
18    counts_father = [job_counts_father.get(job, 0) for job in
19                     jobs]
20
21    x = range(len(jobs))
22    width = 0.4
23
24    plt.figure(figsize=(12, 6))
25    plt.bar(x, counts_mother, width=width, label='Matka', align='
26            center', color='violet', edgecolor='black')
27    plt.bar([p + width for p in x], counts_father, width=width,
28            label='Ojciec', align='center', color='coral', edgecolor=
29            'black')
30
31    plt.xlabel('Zawód', fontsize=14)
32    plt.ylabel('Liczba wystąpień', fontsize=14)
33    plt.title('Dominujące zawody rodziców uczniów', fontsize=16)
34    plt.xticks([p + width / 2 for p in x], jobs, rotation=45,
35               fontsize=12)
36    plt.legend(fontsize=15, edgecolor='black')
37    plt.tight_layout()
38    plt.show()
39
40    return job_counts_mother, job_counts_father
```



Krótki wniosek: Na powyższym wykresie, widzimy liczbę wystąpień każdego z podanych zawodów, obu rodziców uczniów badanej placówki. Widzimy, że przeważającym zawodem (szczególnie wśród ojców) jest 'other', jest dość zrozumiałe, ze względu na to, że do tej grupy wchodzi wszystkie niewymienione w pliku specjalizacje, naturalnym następstwem tego faktu, jest że to matki dominują w pozostałych zawodach (poza usługami).

10 Wnioski

Podsumowanie wyników analizy i obserwacji.

10.1 Oceny i frekwencja:

- Wydaje się, że uczniowie zamieszkujący obszary miejskie odnoszą lepsze sukcesy szkolne w porównaniu do swoich kolegów z obszarów wiejskich. Statystyczna analiza ujawniła, że średnie oceny dla uczniów miejskich wynoszą 10.69, podczas gdy dla uczniów wiejskich jest to 9.47.
- Istnieje również pewna korelacja między frekwencją a wynikami w nauce, choć nie jest ona jednoznaczna. Współczynnik korelacji wynoszący 0.03 sugeruje, że większa frekwencja może mieć korzystny wpływ na wyniki edukacyjne.

10.2 Rozkład ocen i punktów:

- Analiza rozkładu ocen wykazała, że większość uczniów otrzymuje oceny w przedziale [10-12], co sugeruje, że są to przeważnie oceny średnie. Interesującą obserwacją jest również znaczna różnica między liczbą ocen zerowych a ocenami w górnym zakresie skali.
- Jeśli chodzi o punkty za obecność, większość uczniów zdaje się osiągać około 410 punktów. Warto zauważyć, że uczniowie z obszarów wiejskich zdają się osiągać nieco więcej punktów niż uczniowie z obszarów miejskich.

10.3 Zawody rodziców uczniów:

- Dominującym zawodem zarówno matek, jak i ojców uczniów jest kategoria "other". Jednakże, gdy rozważymy specjalizacje matki i ojca osobno, obserwujemy pewne różnice. Przykładowo, matki częściej pracują w zawodach związanych z edukacją i opieką zdrowotną, podczas gdy ojcowie częściej pracują w innych branżach.