

Analýza atributů Pokémonů

Julie Vondráčková

29. 6. 2024

1 Úvod

Jako svou statistickou práci jsem si vybrala analyzovat data o Pokémonech z tohoto odkazu:

[Kaggle datasheet](#).

Má analýza zkoumá souvislost mezi typy Pokémonů a jejich výškou, váhou a útokem. Používáme data z datasetu Pokémonů dostupného na Kaggle.

K vizualizaci a výpočtům jsem použila Python 3.12 a tyto knihovny:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from scipy.stats import f_oneway, shapiro
```

2 Popisná tabulka dat

Popisná tabulka výšky, váhy a útoku všech 781 Pokémonů v datasetu:

Statistika	Výška (m)	Váha (kg)	Útok
Průměr	1.16	61.38	77.69
Std	1.08	109.35	32.22
Min	0.10	0.10	5.00
25%	0.60	9.00	55.00
50%	1.00	27.30	75.00
75%	1.50	64.80	100.00
Max	14.50	999.90	185.00

Data v tabulce jsem získala pomocí tohoto kódu:

```
1 relevantni_sloupce = ['type1', 'height_m', 'weight_kg', 'attack']
2 pokemon_data_cleaned = pokemon_data.dropna(subset=['height_m', 'weight_kg'])
3
4 popisna_statistika = pokemon_data_cleaned[relevantni_sloupce].describe()
5
6 count = popisna_statistika.loc['count']
7 mean = popisna_statistika.loc['mean']
8 std = popisna_statistika.loc['std']
9 min_val = popisna_statistika.loc['min']
10 q25 = popisna_statistika.loc['25%']
11 q50 = popisna_statistika.loc['50%']
12 q75 = popisna_statistika.loc['75%']
13 max_val = popisna_statistika.loc['max']
```

3 Vizualizace

3.1 Rozdělení jednotlivých vlastností

Nejprve jsem si vykreslila histogramy výšky (v metrech), váhy (v kilogramech) a útoku.

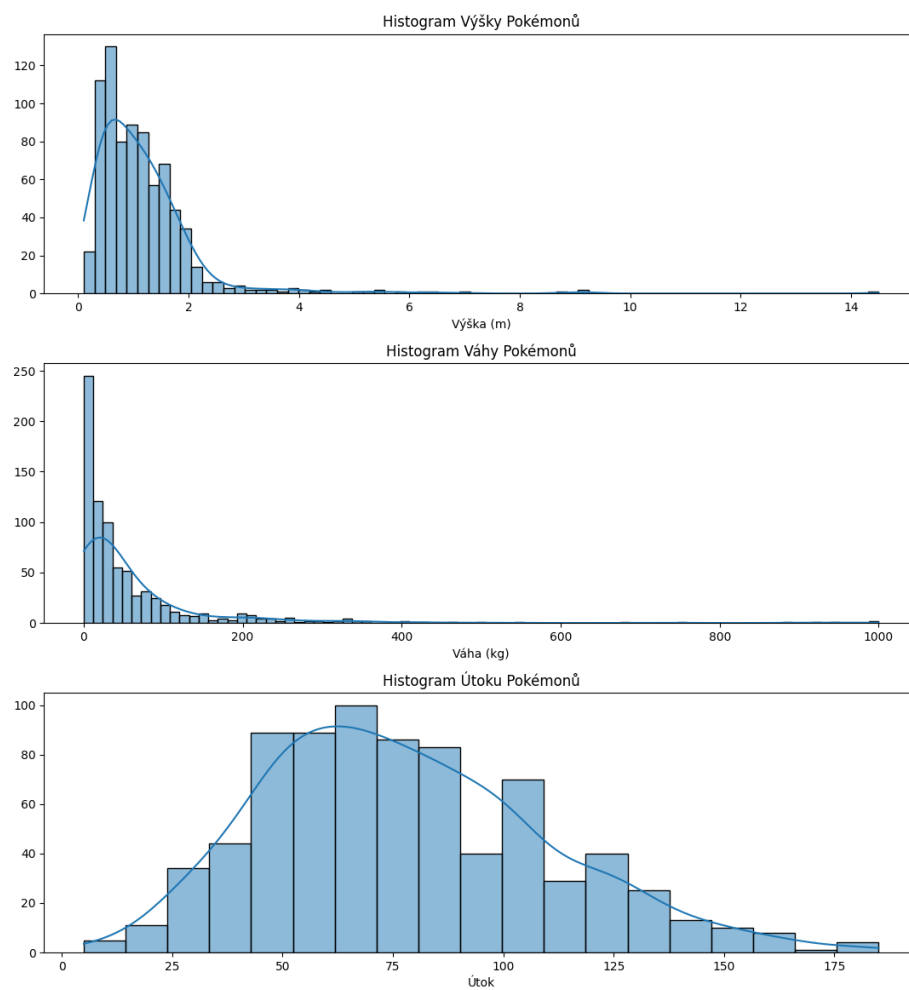


Figure 1: Histogramy výšky, váhy a útoku Pokémonů

Popis: Na histogramu výšky je patrné, že většina Pokémonů má výšku mezi 0.5 a 2 metry. Váha Pokémonů je velmi rozmanitá, většina z nich váží do 100 kg, ale jsou zde i extrémní případy vážící až 999.9 kg. Histogram útoku ukazuje, že útoky jsou poměrně rovnoměrně rozloženy, s průměrem kolem 77.

3.2 Distribuce jednotlivých vlastností podle typu

Vykreslila jsem boxploty pro výšku, váhu a útok podle typu Pokémona.

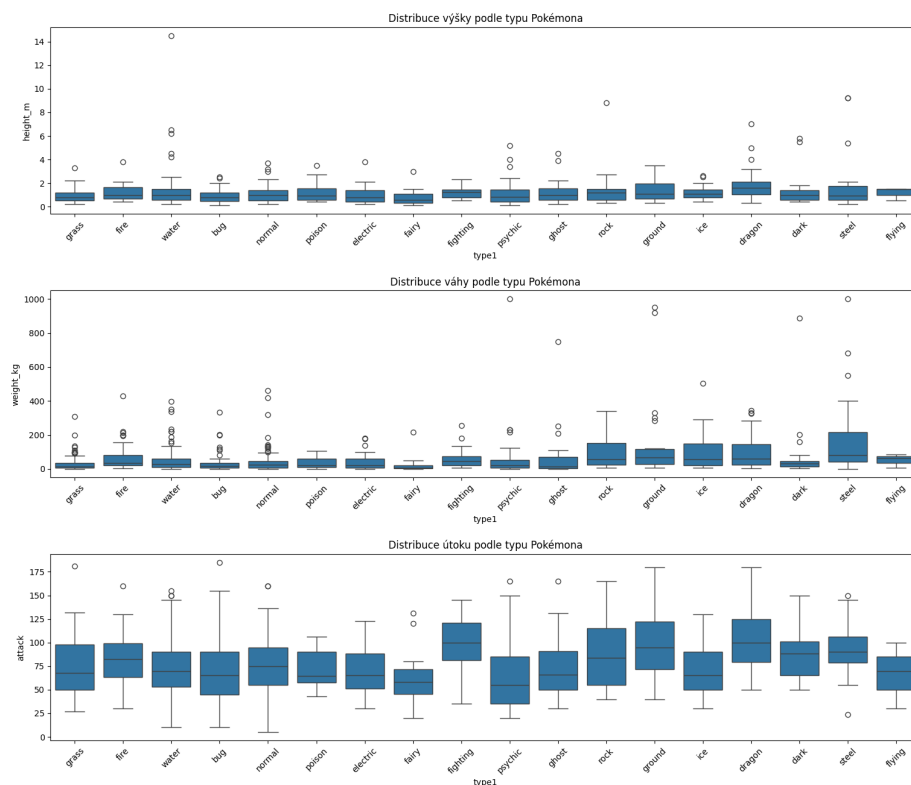


Figure 2: Distribuce výšky, váhy a útoku podle typu Pokémona

Popis: Boxploty ukazují rozložení výšky, váhy a útoku Pokémonů podle jejich typu. Je vidět, že některé typy (např. Dragon, Steel) mají větší rozptýl v hodnotách váhy a útoku. Výška Pokémonů je relativně konzistentní napříč různými typy.

3.3 Scatter Plots

Scatter ploty byly vytvořeny pro vizualizaci vztahů mezi výškou a váhou, výškou a útokem a váhou a útokem, barevně odlišený podle typu Pokémona.

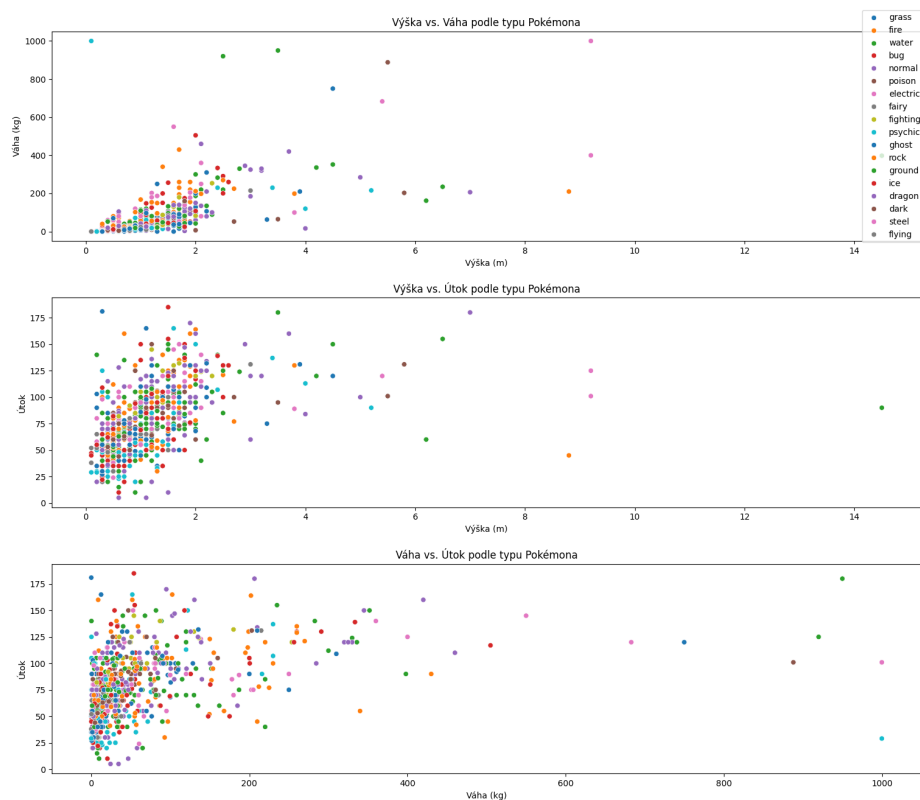


Figure 3: Scatter plots zobrazující vztahy mezi výškou, váhou a útokem Pokémonů

Popis: Scatter ploty ukazují, že existuje pozitivní korelace mezi výškou a váhou Pokémonů, což je očekávané. Vztah mezi výškou a útokem a mezi váhou a útokem není tak výrazný, ale stále lze vidět, že těžší a vyšší Pokémoni mají tendenci mít vyšší útok.

4 Statistická analýza

4.1 ANOVA (ANalysis Of VAriance) Testy

ANOVA testy byly provedeny pro určení, zda existují statisticky významné rozdíly ve výšce, váze a útoku mezi různými typy Pokémonů. Výsledky jsou následující:

Atribut	F-Statistika	p-hodnota
Výška	2.50	7.07e-04
Váha	5.20	5.94e-11
Útok	5.49	9.76e-12

Popis: P-hodnoty ukazují, že existují významné rozdíly ve výšce, váze a útoku mezi různými typy Pokémonů. To znamená, že typ Pokémona má vliv na jeho fyzické a bojové vlastnosti. Konkrétně, typy jako Dragon a Steel mají výrazné rozdíly ve srovnání s ostatními typy.

4.2 Test normality

Pro testování normality distribuce výšky, váhy a útoku byl použit Shapiro-Wilk test. Výsledky jsou následující:

Atribut	W-statistika	p-hodnota
Výška	0.76	1.25e-31
Váha	0.29	0.00e+00
Útok	0.97	1.43e-13

Popis: Výsledky Shapiro-Wilk testu naznačují, že distribuce výšky, váhy a útoku Pokémonů se významně liší od normální distribuce. To je patrné zejména u váhy, kde je p-hodnota extrémně nízká, což potvrzuje, že váhy Pokémonů nejsou normálně rozloženy.

Shapiro–Wilkův test

Shapiro–Wilkův test je statistický test normality, který byl představen Samuel Shapirovem a Martinem Wilkem v roce 1965. Jeho cílem je zjistit, zda daný soubor dat pochází z normálního rozdělení.

Definice

Testová statistika W je definována jako:

$$W = \frac{\left(\sum_{i=1}^n a_i x_{(i)}\right)^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

kde:

- $x_{(i)}$ jsou uspořádané vzestupně (od nejmenšího k největšímu) hodnoty ve vzorku,
- a_i jsou konstanty závislé na očekávaných hodnotách a kovariancích uspořádaných statistických hodnot normálního rozdělení,
- x_i jsou původní hodnoty ve vzorku,
- \bar{x} je průměr hodnot ve vzorku,
- n je počet pozorování ve vzorku.

Interpretace výsledků

Hodnota testové statistiky W se pohybuje mezi 0 a 1. Pokud je hodnota W nízká, znamená to, že data pravděpodobně nepocházejí z normálního rozdělení. Kritické hodnoty pro různé hladiny významnosti jsou tabulkově dostupné, a podle těchto hodnot se určuje, zda nulovou hypotézu (data pocházejí z normálního rozdělení) zamítneme nebo ne.

Použití

Shapiro–Wilkův test se často používá v oblastech, kde je třeba ověřit předpoklad normality před provedením dalších statistických analýz, jako jsou regresní analýzy nebo analýzy rozptylu (ANOVA).

5 Závěr

Na základě provedené analýzy mohu potvrdit, že existují statisticky významné rozdíly ve výšce, váze a útoku mezi různými typy Pokémonů. Tyto rozdíly byly identifikovány pomocí ANOVA testů a potvrzeny Shapiro-Wilk testy normality. Vizualizace v podobě histogramů, boxplotů a scatter plotů poskytují jasný pohled na rozložení a vztahy mezi těmito atributy.

6 Použitý kód

Níže uvádím použitý kód v jazyce Python pro načtení dat, výpočet popisné statistiky, vytvoření vizualizací a provedení statistických testů.

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from scipy.stats import f_oneway, shapiro
5
6 # Načtení datasetu
7 cesta_k_souboru = "path_to_your_pokemon.csv"
8 pokemon_data = pd.read_csv(cesta_k_souboru)
9
10 # Příprava dat
11 relevantni_sloupce = ['type1', 'height_m', 'weight_kg', 'attack']
12 pokemon_data_cleaned = pokemon_data.dropna(subset=['height_m', 'weight_kg'])
13
14 # Popisná statistika
15 popisna_statistika = pokemon_data_cleaned[relevantni_sloupce].describe()
16
17 # Box Ploty
18 plt.figure(figsize=(18, 12))
19
20 # Výška podle typu
21 plt.subplot(3, 1, 1)
22 sns.boxplot(x='type1', y='height_m', data=pokemon_data_cleaned)
23 plt.title('Distribuce výšky podle typu Pokémona')
24 plt.xticks(rotation=45)
25
26 # Váha podle typu
27 plt.subplot(3, 1, 2)
28 sns.boxplot(x='type1', y='weight_kg', data=pokemon_data_cleaned)
29 plt.title('Distribuce váhy podle typu Pokémona')
30 plt.xticks(rotation=45)
31
32 # Útok podle typu
33 plt.subplot(3, 1, 3)
34 sns.boxplot(x='type1', y='attack', data=pokemon_data_cleaned)
35 plt.title('Distribuce útoku podle typu Pokémona')
36 plt.xticks(rotation=45)
37
38 plt.tight_layout()
39 plt.show()
40
41 # Scatter Plots with a single legend
42 fig, axes = plt.subplots(3, 1, figsize=(18, 12))
43
44 # Výška vs. Váha
45 scatter1 = sns.scatterplot(x='height_m', y='weight_kg', hue='type1',
46                             data=pokemon_data_cleaned, palette='tab10', ax=axes[0])
47 axes[0].set_title('Výška vs. Váha podle typu Pokémona')
48 axes[0].set_xlabel('Výška (m)')
49 axes[0].set_ylabel('Váha (kg)')
50 handles, labels = scatter1.get_legend_handles_labels()
```



```

50 scatter1.legend_.remove()
51
52 # Výška vs. Útok
53 scatter2 = sns.scatterplot(x='height_m', y='attack', hue='type1',
54                             data=pokemon_data_cleaned, palette='tab10', legend=False, ax=
55                                 axes[1])
56 axes[1].set_title('Výška vs. Útok podle typu Pokémona')
57 axes[1].set_xlabel('Výška (m)')
58 axes[1].set_ylabel('Útok')
59
60 # Váha vs. Útok
61 scatter3 = sns.scatterplot(x='weight_kg', y='attack', hue='type1',
62                             data=pokemon_data_cleaned, palette='tab10', legend=False, ax=
63                                 axes[2])
64 axes[2].set_title('Váha vs. Útok podle typu Pokémona')
65 axes[2].set_xlabel('Váha (kg)')
66 axes[2].set_ylabel('Útok')
67
68 # Add a single legend to the upper right corner of the figure
69 fig.legend(handles, labels, loc='upper right', bbox_to_anchor=(1,
70     1), ncol=1)
71
72 plt.tight_layout()
73 plt.show()
74
75 # Histogramy
76 plt.figure(figsize=(18, 12))
77
78 # Histogram Výšky
79 plt.subplot(3, 1, 1)
80 sns.histplot(pokemon_data_cleaned['height_m'], kde=True)
81 plt.title('Histogram Výšky Pokémonů')
82 plt.xlabel('Výška (m)')
83 plt.ylabel('Frekvence')
84
85 # Histogram Váhy
86 plt.subplot(3, 1, 2)
87 sns.histplot(pokemon_data_cleaned['weight_kg'], kde=True)
88 plt.title('Histogram Váhy Pokémonů')
89 plt.xlabel('Váha (kg)')
90 plt.ylabel('Frekvence')
91
92 # Histogram Útoku
93 plt.subplot(3, 1, 3)
94 sns.histplot(pokemon_data_cleaned['attack'], kde=True)
95 plt.title('Histogram Útoku Pokémonů')
96 plt.xlabel('Útok')
97 plt.ylabel('Frekvence')
98
99 plt.tight_layout()
100 plt.show()
101
102 # ANOVA Testy
103 anova_vyska = f_oneway(*[group['height_m'].dropna() for name, group
104     in pokemon_data_cleaned.groupby('type1')])
105 anova_vaha = f_oneway(*[group['weight_kg'].dropna() for name, group
106     in pokemon_data_cleaned.groupby('type1')])

```

```

100 anova_utok = f_oneway(*[group['attack'].dropna() for name, group in
    pokemon_data_cleaned.groupby('type1')])
101
102 anova_vysledky = {
103     'Atribut': ['Výška', 'Váha', 'Útok'],
104     'F-Statistika': [anova_vyska.statistic, anova_vaha.statistic,
    anova_utok.statistic],
105     'p-hodnota': [anova_vyska.pvalue, anova_vaha.pvalue, anova_utok
    .pvalue]
106 }
107
108 anova_vysledky_df = pd.DataFrame(anova_vysledky)
109 print("ANOVA Výsledky:")
110 print(anova_vysledky_df)
111
112 # Test normality (Shapiro-Wilk test)
113 shapiro_vyska = shapiro(pokemon_data_cleaned['height_m'])
114 shapiro_vaha = shapiro(pokemon_data_cleaned['weight_kg'])
115 shapiro_utok = shapiro(pokemon_data_cleaned['attack'])
116
117 normalita_vysledky = {
118     'Atribut': ['Výška', 'Váha', 'Útok'],
119     'W-statistika': [shapiro_vyska.statistic, shapiro_vaha.
    statistic, shapiro_utok.statistic],
120     'p-hodnota': [shapiro_vyska.pvalue, shapiro_vaha.pvalue,
    shapiro_utok.pvalue]
121 }
122
123 normalita_vysledky_df = pd.DataFrame(normalita_vysledky)
124 print("Výsledky testu normality (Shapiro-Wilk):")
125 print(normalita_vysledky_df)

```