

2. Vizualizace dat

Zadání:

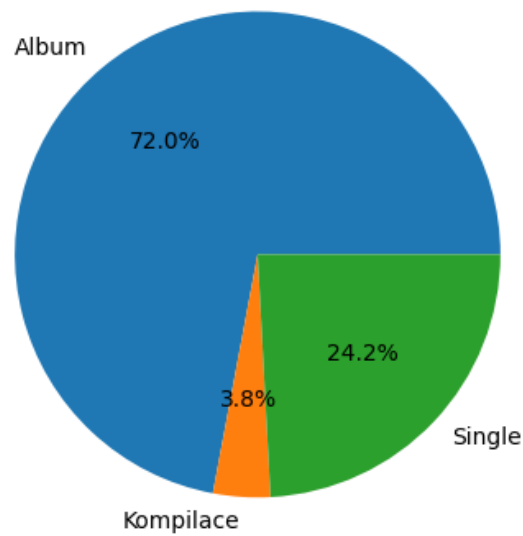
V jednom ze cvičení jste probírali práci s moduly pro vizualizaci dat. Mezi nejznámější moduly patří matplotlib (a jeho nadstavby jako seaborn), pillow, opencv, aj. Vyberte si nějakou zajímavou datovou sadu na webovém portále Kaggle a proveďte datovou analýzu datové sady. Využijte k tomu různé typy grafů a interpretujte je (minimálně alespoň 5 zajímavých grafů). Příklad interpretace: z datové sady pro počasí vyplynulo z liniového grafu, že v létě je vyšší rozptyl mezi minimální a maximální hodnotou teploty. Z jiného grafu vyplývá, že v létě je vyšší průměrná vlhkost vzduchu. Důvodem vyššího rozptylu může být absorpce záření vzduchem, který má v létě vyšší tepelnou kapacitu.

Řešení:

Jako zdroj dat jsem využil Kaggle a jejich obsáhlou databázi datasetů. Dataset který jsem se nakonec rozhodl využít „[Spotify and Youtube](#)“ dataset. Dataset obsahuje 20000 unikátních songů. Sloupců je dohromady 28 s informacemi které se týkají, počtu zhlédnutí a streamů, generované sloupce vypovídající o jejich vlastnostech a dalších. Zajímalo mě hned několik témat mezi nimi bylo i využití generovaných sloupců ke kategorizaci.

Na vizualizaci dat jsem využil matplotlib a pandas. Vytvořil jsem 5 grafů (2 další jsem bral jako nedokonalé a rozhodl jsem se je vynechat). Program začíná načtením dat, přičemž následují jednotlivé funkce, které definují jednotlivé grafy. Za každou funkcí se zavolají. To mi sloužilo hlavně k rychlému debuggingu, abych pouze ukazoval grafy, na kterých zrovna pracuji.

Rozdělení Alb do kategorií



Z grafu lze vyčíst, že většina songů v datasetu je zařazená v albech poté v menším počtu následují singly a ve velice malé míře jsou to jen kompilace.

```
10 def TypySongu(DataSet):
11     ...
12     ...
13     Rozdělení songů podle typu:
14     Album,
15     Single,
16     Kompilace
17     ...
18     grouped = DataSet.groupby(["Album_type"]).size()
19     grouped = grouped.rename(lambda x: x.title())
20     grouped = grouped.rename(
21         {
22             "Compilation": "Kompilace",
23             "Single": "Single"
24         }
25     )
26     Indexes = grouped.index
27     Indexes.name = ""
28     ...
29     fig, ax = plt.subplots()
30     grouped.plot().pie(x=grouped, labels=Indexes, autopct='%1.1f%%')
31     ax.set_title("Rozdělení Alb do typů")
32     ...
33     plt.show()
```

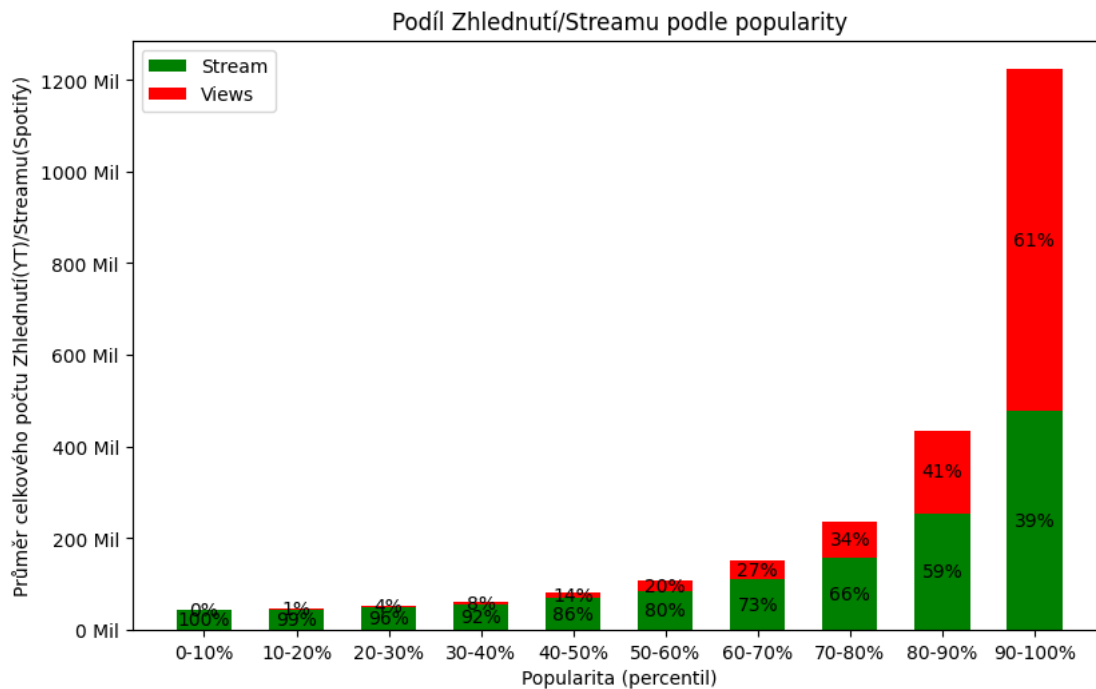


Z grafu lze vyčíst pár kategorií do kterých jsem songy zařadil. Kvůli nejednotnosti, jestli song je určité kategorie, jsem za pomoci této funkce zjišťoval přes prahy, jestli je daný song určitého typu, jelikož sloupce obsahují hodnoty „jistoty“ která je v rozmezí od 0-1.

Mnoho songů je nezařazených, ale určitě pokud by se pohrálo správně s prahy, tak by jich mohlo být o dost méně a však by to bylo za cenu jistoty, že song je určitého typu.

```
def zpracujDataSet(x):
    """
    Zde jde o nastavení hranic jak moc chceme být jistý, že určité songy jsou určitého charakteru
    vytáhnu z popisu Datasetu:
    Speechiness: detects the presence of spoken words in a track.
        The more exclusively speech-like the recording (e.g. talk show, audio book, poetry), the closer to 1.0 the attribute value. Values above 0.66
        describe tracks that are probably made entirely of spoken words. Values between 0.33 and 0.66 describe tracks that may contain both music and speech,
        either in sections or layered, including such cases as rap music. Values below 0.33 most likely represent music and other non-speech-like tracks.
    Acousticness: a confidence measure from 0.0 to 1.0 of whether the track is acoustic. 1.0 represents high confidence the track is acoustic.
    Instrumentalness: predicts whether a track contains no vocals. "Ooh" and "aah" sounds are treated as instrumental in this context.
        Rap or spoken word tracks are clearly "vocal". The closer the instrumentalness value is to 1.0,
        the greater likelihood the track contains no vocal content. Values above 0.5 are intended to represent instrumental tracks,
        but confidence is higher as the value approaches 1.0.
    Liveness: detects the presence of an audience in the recording. Higher liveness values represent an increased probability that the track was performed live.
        A value above 0.8 provides strong likelihood that the track is live.
    """
    x["Speechiness"] = x["Speechiness"] > 0.1 and x["Speechiness"] < 0.66
    x["Acousticness"] = x["Acousticness"] > 0.75
    x["Instrumentalness"] = x["Instrumentalness"] > 0.5
    x["Liveness"] = x["Liveness"] > 0.8
    return x

66 TypovyDataSet = DataSet(["#", "Speechiness", "Acousticness", "Instrumentalness", "Liveness"])
67 TypovyDataSet = TypovyDataSet.apply(zpracujDataSet, axis=1)
68
69 SongCount = TypovyDataSet["#"].count()
70 TypovyDataSet=TypovyDataSet[TypovyDataSet==True].count(axis=0)
71 TypovyDataSet["Nezařazené"] = SongCount-sum(TypovyDataSet)
72 TypovyDataSet=TypovyDataSet.rename(
73     {
74         "#": "Počet Songů",
75         "Speechiness": "Se zpěvem",
76         "Acousticness": "Akustické",
77         "Instrumentalness": "Instrumentální",
78         "Liveness": "Naživo"
79     }
80 )
81
82 bar_colors = ['black', 'green', 'red', 'yellow', 'blue']
83 helpMe = TypovyDataSet.plot.pie(colors=bar_colors)
84 helpMe.set_title("Rozdělení do kategorií")
85 plt.show()
```

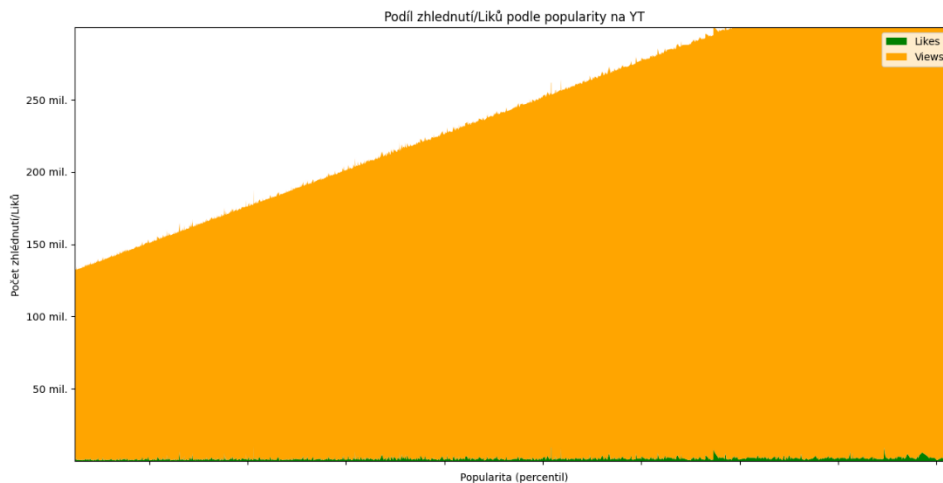
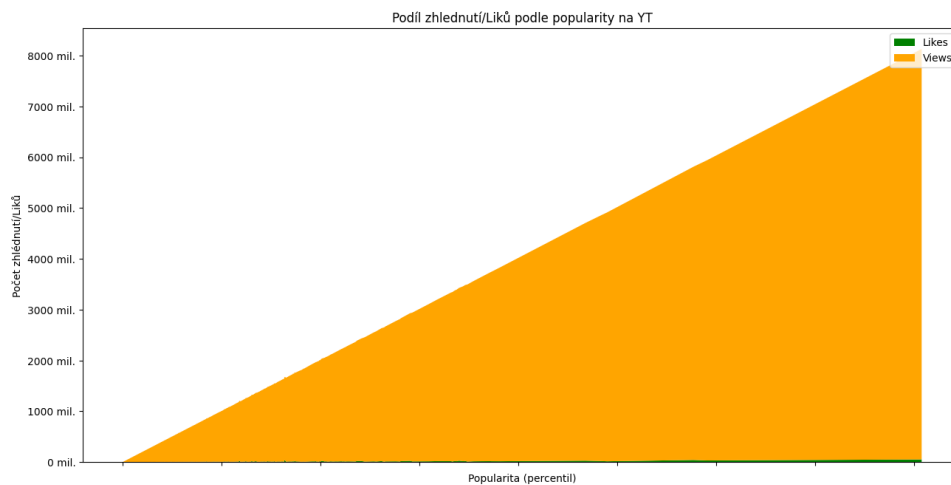


Graf, který ukazuje průměrný zhlédnutí a počet streamů s viditelným podílem Spotify vs YouTube. Lze vidět, že méně známé songy se více poslouchají na Spotify mezitím co ty známější jsou více a více spouštěny na YouTube.

```

89 def ZhlédnutíStreamy(DataSet):
90     views_duration = pd.DataFrame(zip(DataSet["Stream"],DataSet["Views"]), columns=['Stream','Views'])
91     views_duration["Stream"] = views_duration["Stream"].apply(lambda x : x /1000000)
92     views_duration["Views"] = views_duration["Views"].apply(lambda x : x /1000000)
93     views_duration = views_duration.sort_values(by="Views")
94     views_duration = views_duration.reset_index(drop=True)
95     views_duration = views_duration.groupby(pd.qcut(views_duration.index, 10)).mean()
96     views_duration = views_duration.reset_index(drop=True)
97     label1 = "Stream"
98     label2 = "Views"
99     fig,ax1 = plt.subplots()
100
101
102     views_duration["Stream"] = views_duration["Stream"].apply(lambda x : round(x,3))
103     views_duration["Views"] = views_duration["Views"].apply(lambda x : round(x,3))
104     views_duration["ViewsRatio"] = views_duration["Views"] / (views_duration["Views"] + views_duration["Stream"]) * 100
105     views_duration["StreamRatio"] = views_duration["Stream"] / (views_duration["Views"] + views_duration["Stream"]) * 100
106
107     ax1.yaxis.set_major_formatter(ticker.FuncFormatter(lambda y,pos: f"{y:.0f} Mil"))
108     for index, row in views_duration.iterrows():
109         p = plt.bar(str((index)*10)+"-"+str((index+1)*10)+"%",row["Stream"],width=0.6,label=label1, color="green")
110         ax1.bar_label(p,fmt=f"{row.StreamRatio:.0f}%", label_type='center')
111         p = plt.bar(str((index)*10)+"-"+str((index+1)*10)+"%",row["Views"],width=0.6,label=label2,bottom=row["Stream"],color="red")
112         ax1.bar_label(p,fmt=f"{row.ViewsRatio:.0f}%", label_type='center')
113
114     ax1.set_xlabel("Popularita (percentil)")
115     ax1.set_ylabel("Průměr celkového počtu Zhlédnutí(YT)/Streamu(Spotify)")
116     ax1.set_title("Podíl Zhlédnutí/Streamu podle popularity")
117     ax1.legend(["Stream","Views"])
118     print(views_duration)
119     plt.show()

```

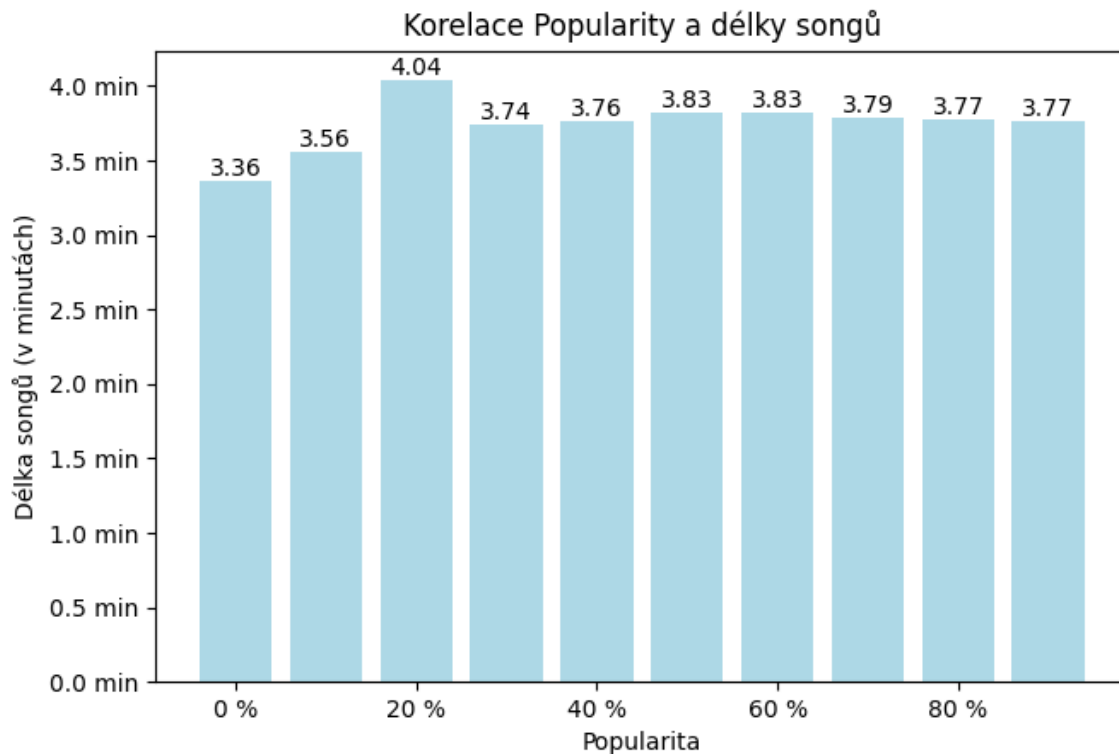


Graf, na kterém jsem chtěl ukázat podíl zhlédnutí a lajků které song na YouTube dostal. Jde vidět, jak absurdně malý podíl to je (podíl je lajků na: 0.6993400313173325 %). I přes malou čitelnost grafu jsem se ho rozhodl ponechat, aby byla vidět rozdíl mezi tím kolik lidí video zhlédlo a kolik jich jej lajknulo

```

125 def SongsAndLikes(DataSet):
126     """
127     Youtube songy popularita a počet liků
128     """
129     views_likes = pd.DataFrame(zip(DataSet["Views"],DataSet["Likes"]), columns=['Views','Likes'])
130     views_likes = views_likes.dropna()
131     views_likes = views_likes.sort_values(by="Views")
132     views_likes = views_likes.reset_index(drop=True)
133
134     label1 = "Views"
135     label2 = "Likes"
136     fig,ax1 = plt.subplots()
137     print(f"Průměrný podíl je: {(views_likes['Likes'].sum() / (views_likes['Views'].sum() + views_likes['Likes'].sum()))*100}%")
138
139     views_likes["LikesRatio"] = views_likes["Likes"] / views_likes["Views"] * 100
140     views_likes["ViewsRatio"] = 100 - views_likes["LikesRatio"]
141
142     ax1.xaxis.set_major_formatter(ticker.FuncFormatter(lambda x,pos: ""))
143     ax1.yaxis.set_major_formatter(ticker.FuncFormatter(lambda y,pos: f"{y/1000000:.0f} mil."))
144     ax1.stackplot(views_likes["Views"], (views_likes["Likes"],views_likes["Views"]),colors=("green","orange"))
145     ax1.set_xlabel("Popularita (percentil)")
146     ax1.set_ylabel("Počet zhlédnutí/Liků")
147     ax1.set_title("Podíl zhlédnutí/Liků podle popularity na YT")
148     ax1.legend(["Likes","Views"])
149     plt.show()

```

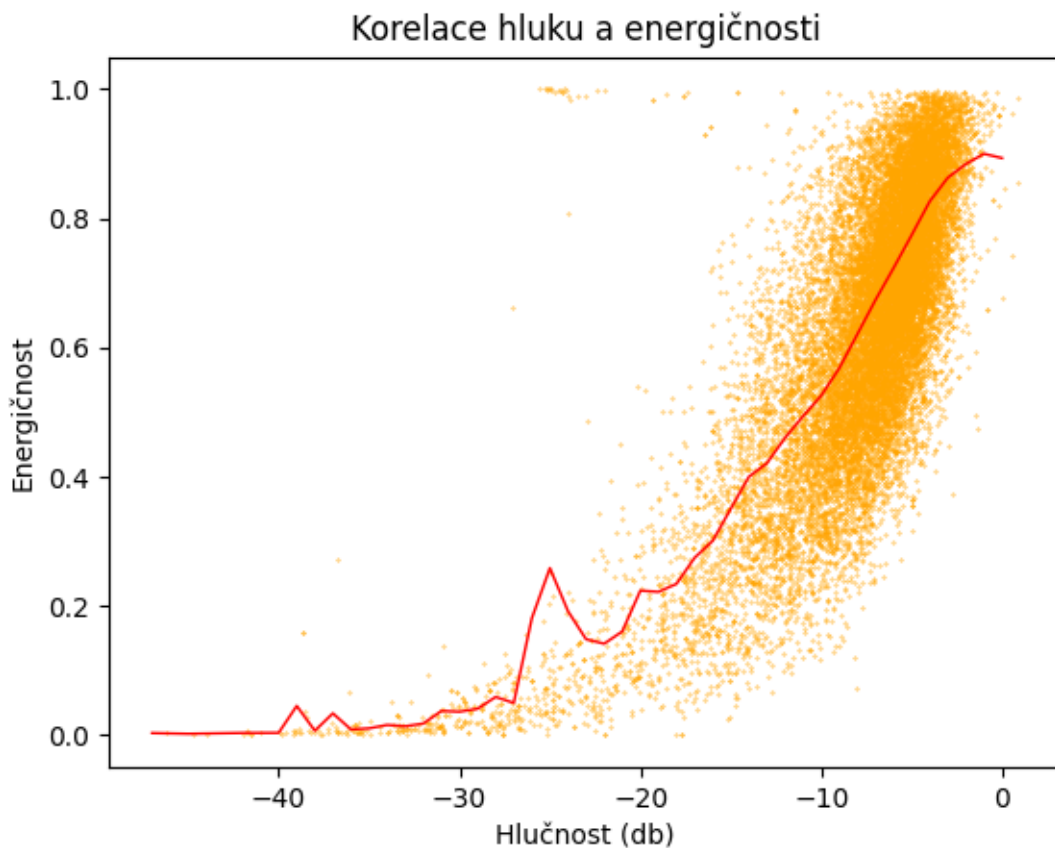


Tímto grafem jsem chtěl zjistit, jestli je nějaká korelace mezi délkou songů a popularity. Lze vidět, že na délce songu zas tolik nezáleží, ale lze vidět mírné výkyvy u méně známějších písničků

```

151 def PopularitaADelkaSongu(DataSet):
152     duration = pd.DataFrame(zip(DataSet["Duration_ms"],DataSet["Stream"],DataSet["Views"]), columns=['Duration_ms','Stream','Views'])
153
154     duration = duration.apply(lambda x : x /1000 / 60)
155     duration["Popularity"] = duration.apply(lambda x : x["Stream"] + x["Views"] ,axis=1)
156     duration = duration.sort_values(by="Popularity")
157     duration = duration.reset_index(drop=True)
158     duration = duration.groupby(pd.qcut(duration["Popularity"], 10)).mean()
159     duration = duration.reset_index(drop=True)
160     fig,ax1 = plt.subplots()
161
162     ax1.yaxis.set_major_formatter(ticker.FuncFormatter(lambda y,pos: f"{y:.1f} min"))
163     ax1.xaxis.set_major_formatter(ticker.FuncFormatter(lambda x,pos: f"{x*10:.0f} %"))
164
165
166     ax1.set_title("Průměrná délka písniček v po 10 percentilech")
167
168     bars = ax1.bar(duration.index,duration["Duration_ms"],color="lightblue")
169     ax1.bar_label(bars,fmt="%.2f")
170     ax1.set_title("Korelace Popularity a délky songů")
171     ax1.set_xlabel("Popularita")
172     ax1.set_ylabel("Délka songů (v minutách)")
173     plt.show()

```



U tohoto grafu mě zajímalo, jestli hlučnost autoři použili ke generování měřítka energičnosti. Lze vidět nějaká korelace mezi těmito dvěma vlastnostmi, přesto lze vidět spoustu výjimek jako například obrovský nárůst a pokles v půlce grafu

```

180 def HlukAEnergeticnost(DataSet):
181     duration = pd.DataFrame(zip(DataSet["Loudness"],DataSet["Energy"]), columns=['Loudness','Energy'])
182     duration = duration.sort_values(by="Loudness")
183     duration = duration.reset_index(drop=True)
184
185     fig,ax1 = plt.subplots()
186     ax1.scatter(duration["Loudness"],duration["Energy"],s=0.1,color="orange")
187
188     line = duration
189     line = line.dropna()
190     line["Loudness"] = line["Loudness"].apply(lambda x : math.floor(x))
191
192     line = line.groupby("Loudness").mean()
193     ax1.set_title("Korelace hluku a energičnosti")
194     ax1.set_xlabel("Hlučnost (db)")
195     ax1.set_ylabel("Energičnost")
196
197     ax1.plot(line.index,line["Energy"], label='Data', linewidth=1, color="red")
198     plt.show()

```