

Zadanie 2 - Analýza dát a regresory

I-SUNS: Strojové učenie a neurónové siete

Pavol Ďuďák
xdudakp@stuba.sk

18. novembra 2021

1 Úvod

Po použití neurónových sietí na jednoduché klasifikačné úlohy si v tomto zadaní predstavíme alternatívu, regresory. Regresory sú veľa krát šikovnejšie pri jednoduchších úlohách pri ktorých používame štruktúrované (napr. csv) dáta. Regresory nemajú toľko konfiguračných parametrov ako neurónové siete a ich použitie je všeobecne jednoduchšie.

Cieľom zadania je oboznámiť sa s použitím rôznych regresorov a ich využitím pri úlohe predpovedi hlasitosti skladieb na základe ich rôznych atribútov. Zároveň nad poskytnutými dátami bola vypracovaná exploratívna analýza (EDA).

2 Spracovanie a analýza dát

2.1 Spracovanie dát

Pred analýzou je nutné dáta spracovať. Spracovanie zahŕňa odstránenie nepotrebných stĺpcov, kódovanie nečíselných hodnôt, odstránenie null hodnôt a odstránenie hraničných hodnôt.

Dáta si načítame z poskytnutého .csv

```
1 df_raw: DataFrame = pd.read_csv('sample_data/spotify_train.csv')
```

Listing 1: Načítanie tréningových dát

následne upravíme dátum vydania iba na rok a extrahujeme všetky žánre a roky vydania ktoré budú použité pre zakódovanie vo forme masky.

```
1 df_raw['release_date'] = df_raw['release_date'].apply(lambda d: d[0:4])
2 genres: Set[str] = extract_unique_genres(df_raw)
3 release_dates: Set[str] = collect_release_data(df_raw)
4 df_raw['artist_genres'] = df_raw['artist_genres'].apply(str_to_list)
5 df_raw['index'] = df_raw.index
```

Listing 2: Úprava rokov a extrakcia žánrov a rokov vydania

Z extrahovaných hodnôt žánrov a rokov vydania vytvoríme masku.

```
1 # vytvorim bool masku top 50 udajov
2 genres_mask = create_bool_mask(df_raw, genres, 'artist_genres')
3 genres_mask['index'] = genres_mask.index
4 genres_index = genres_mask.index
5
6 release_mask = create_bool_mask_date(df_raw, release_dates, 'release_date')
7 release_mask['index'] = release_mask.index
8 releases_index = release_mask.index
```

Listing 3: Vytvorenie masky

Aby sme mali v datasete iba číselné hodnoty tak všetky boolean hodnoty zmeníme na 0 alebo 1.

```
1 genres_mask.replace({False: 0, True: 1}, inplace=True)
2 release_mask.replace({False: 0, True: 1}, inplace=True)
```

Listing 4: Kódovanie bool hodnôt na čísla

Aby náš dataset nemal príliš veľa stĺpcov (keďže nečíselné údaje maskujeme do stĺpcov) tak vyberieme top 50 žánrov a 50 rokov vydania a tieto údaje zamaskujeme do nášho datasetu.

```
1 top_genres = genres_mask.sum(axis=0).drop(labels=['index']).nlargest(50)
2 top_years = release_mask.sum(axis=0).drop(labels=['index']).nlargest(50)
3 # maskovanie
4 for genre in top_genres.keys():
5     df_raw[genre] = df_raw['artist_genres'].apply(lambda x: 1 if genre in x else 0)
6 for year in top_years.keys():
7     df_raw[year] = df_raw['release_date'].apply(lambda x: 1 if x == year else 0)
```

Listing 5: Maskovanie

Pre ďalšie spracovanie sú uvedené stĺpce nerelevantné tak ich môžeme vyhodíť.

```
1 df_raw.drop(['url', 'playlist_id', 'playlist_description', 'playlist_name', '
    playlist_url', 'query'], axis=1, inplace=True)
```

Listing 6: Odstánenie nerelevantných stĺpcov

Dôležitým krokom pri spracovaní dát je odstránenie duplikátov. Pri poskytnutom datasete sme nemali žiaden jednoznačný ukazovateľ ktorý by hovoril že ide o duplikát takže sme museli vymyslieť stratégiu ktorou budeme odstraňovať duplikáty.

1KRpleeB2JEG6AW6Uir3QX	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody	0 1975-11-21	354346	FALSE	0.41
6Pm2XMI51SGCeP19HscL	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody	0 1975-11-21	354840	FALSE	0.31
3z8t0T07ReDPLUEnYhWZb	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody	72 2018-10-19	354947	FALSE	0.31
1p0E0B0MTT7WdZcnrA	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - 2011 Remaster	0 1981-10-26	355466	FALSE	0.31
64uQ28GQ2a70YcmEcCl	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - Live	40 2007-10-09	328933	FALSE	0.21
7xHATAMD7ezTZGYNAmy5R	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - Live Aid	59 2018-10-19	147840	FALSE	0.31
0U49xOQCWh6zDWcuSr8Why	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - Live At The Montreal Forum / November 1981	42 2007-10-29	328933	FALSE	0.31
5YIIdKE76GacbcTxuoOf93	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - Live At Wembley Stadium / July 1986	0 1992-05-26	350573	FALSE	0.21
2XbcujvemK0hGh00b4HAXQ	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - Live At Wembley Stadium / July 1986	51 1992-05-26	350573	FALSE	0.21
2vwoDagLdKs4IPG4awF7Y	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - Remastered 2011	42 1975-11-21	354320	FALSE	0.41
2GkNycMAdVMDvdxR2x	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - Remastered 2011	52 2011-01-01	354320	FALSE	0.31
1AH0GHS9vPS0msWgNW08EY	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - Remastered 2011	2 1975-11-21	354320	FALSE	0.41
5wSXPGYMYcWJQE2b0UFEH	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - Remastered 2011	0 2011-01-01	355466	FALSE	0.31
20BoFMjX94NryVZ5K8p8Zf	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - Remastered 2011	51 1981-10-26	355466	FALSE	0.31
7FiyTwD0nx5a1ekYXQZJ	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - Remastered 2011	74 1975-11-21	354320	FALSE	0.31
4u7EnebtmKwU4H33cF5Qv	1dfeR4HaWDbWgFHLKxg1d	Queen	Bohemian Rhapsody - Remastered 2011	82 1975-11-21	354320	FALSE	0.41

Obr. 1: Ukážka duplikátov

V datasete nastávali často situácie uvedené na obrázku 1. V tomto prípade môžeme vidieť nie len že jedna skladba má viacero rôznych verzií no aj to že skladba s rovnakým názvom sa vyskytuje v datasete viac krát. Tento problém sme sa rozhodli riešiť vybraním najpopulárnejšej skladby

zo skupiny skladieb z rovnakým názvom. Túto operáciu vykonávame pre každého interpreta a jeho skladby. Algoritmus zoskupuje skladby podľa interpreta, potom všetky skladby interpreta zoskupí podľa názvu (dostaneme množinu skladieb interpreta z rovnakým názvom) a z týchto vyberie len tú najpopulárnejšiu podľa stĺpca *popularity*.

```
1 # vycistim duplikaty songov
2 artist_groups = create_groups(df_raw, 'artist')
3 filtered_artists: List[DataFrame] = []
4 for artist in artist_groups:
5     grouped_by_song: List[DataFrame] = create_groups(artist, 'name')
6     filtered_artists.append(pd.concat(filter_only_relevant_song(grouped_by_song))) #
7     songy ktore su uz odfiltrovane zlcim do jedneho DF
8 df_filtered = pd.concat(filtered_artists)
```

Listing 7: Čistenie duplikátov

V tento moment máme dataset skladieb ktoré neobsahujú duplikáty. Poslednými krokmi sú odstránenie null hodnôt a outliers (outliers sa odstraňujú až neskôr).

```
1 import numpy as np
2 from scipy import stats
3 df_filtered = df_filtered.dropna()
4 X[(np.abs(stats.zscore(X)) < 3).all(axis=1)] # tento krok je realne v kode az neskôr
```

Listing 8: Mazanie null hodnôt a outliers

2.2 Exploratívna analýza dát

Po načítaní a úprave dát môžeme robiť analýzu. V tejto analýze sa po pohľade na dáta spýtame niekoľko otázok na ktoré si v zápätí grafmi odpovieme.

1. Ktoré žánre majú najväčšie zastúpenie?
2. Ktorý z interpretov sú najsledovanejší?
3. Ako energicky pôsobila hudba počas rokov?
4. Ako hlasito pôsobila hudba počas rokov?

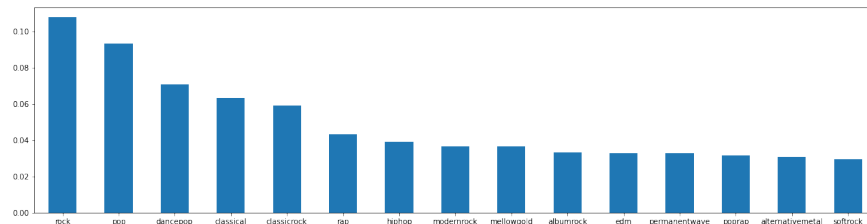
Manipuláciou vyčistených dát dokážeme na tieho otázky jednoducho a graficky odpovedať.

Ktoré žánre majú najväčšie zastúpenie?

Pre každý žánr vypočítame pomer jeho zastúpenia v datasete a vyberieme prvých 15 záznamov ktoré si zobrazíme do stĺpcového diagramu.

```
1 # hladame najviac zastupene
2 genres = top_genres.keys()
3 genre_percent = { genre: len(df_filtered.loc[df_filtered[genre]==True]) / len(
4     df_filtered) for genre in list(genres) }
5 s_genres: Series = Series(genre_percent)
6 s_genres.sort_values(ascending=False, inplace=True)
7 s_genres.head(15).plot.bar(rot=0, figsize=(20,5))
```

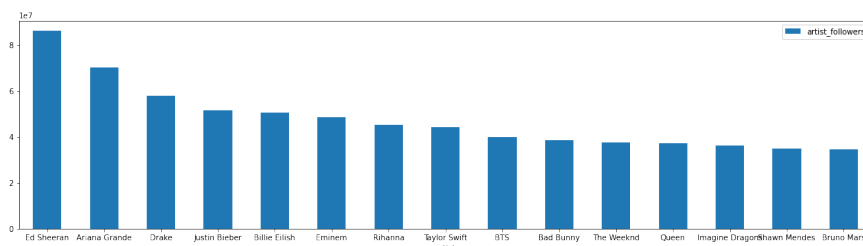
Listing 9: Hľadanie top žánrov



Obr. 2: Zastúpenie žánrov v datasete

Ktorý z interpretov sú najsledovanejší

Pre získanie interpretov ktorý sú najsledovanejší potrebujeme spraviť jednoduchú transformáciu dát. Zoskupíme si dáta podľa interpreta, zoberieme jeho náhodnú prvú skladbu a tieto skladby zoradíme podľa počtu sledovateľov. Tento počet sa viaže na interpreta takže nám je jedno ktorú jeho skladbu vyberieme. Následne tieto dáta zobrazíme.



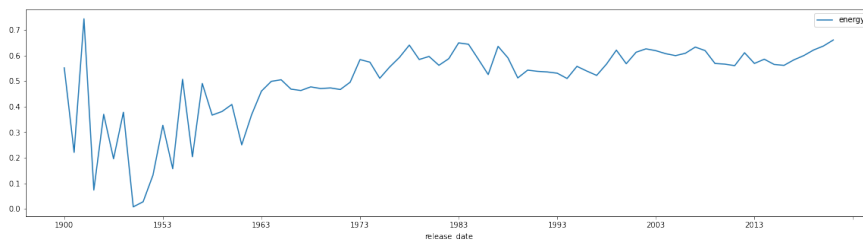
Obr. 3: Najsledovanejší interpreti

```
1 artist_leaderboard = df_filtered[['name', 'artist', 'artist_followers']].groupby(['artist']).first().reset_index().sort_values(by='artist_followers', ascending=False).head(15).plot.bar(x='artist', y='artist_followers', rot=0, figsize=(20,5))
```

Listing 10: Hľadanie najsledovanejších interpretov

Ako energicky pôsobila hudba počas rokov?

Pre zobrazenie tohto atribútu jediné čo potrebujeme spraviť je usporiadať dáta podľa roku vydania a priemery atribútu *energy* zobrazíť na grafe.



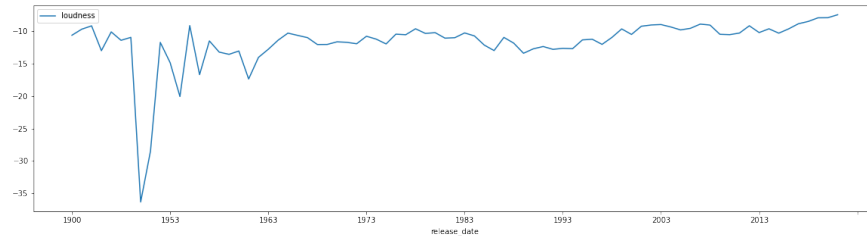
Obr. 4: Priebeh 'energičnosti' skladieb

```
1 df_filtered[['release_date', 'energy']].sort_values(by='energy', ascending=False).
  groupby('release_date').mean().reset_index().plot.line(x='release_date', y='energy',
  , rot=0, figsize=(20,5))
```

Listing 11: Hľadanie energičnosti

Ako hlasito pôsobila hudba počas rokov?

Rovnakým spôsobom vieme zobrazíť aj priebeh pocitu hlučnosti skladieb počas rokov.



Obr. 5: Priebeh "hlučnosti" skladieb

```
1 df_filtered[['release_date', 'loudness']].sort_values(by='loudness', ascending=False).
  groupby('release_date').mean().reset_index().plot.line(x='release_date', y='
  loudness', rot=0, figsize=(20,5))
```

Listing 12: Hľadanie hlučnosti

3 SVM regresor

Pre prvotný pokus bol použitý regresor SVR z RBF kernelom. Tento kernel je ideálny pre ne-lineárne dáta. Pokiaľ by pri tom istom regresore bol použitý lineárny kernel tak tréning nie že by len dlho trvalo ale produkovalo by aj nepresné výsledky. V tom to prípade sú hodnoty hyperparametrov:

- kernel = rbf
- C = 1
- gamma = 1 / pocet stlpcov

```
1 regressor = SVR(kernel = 'rbf')
2 regressor.fit(X, y.to_numpy())
```

Listing 13: SVR regresor

3.1 Hyperparametre

Hyperparametre sú parametre ktorý ovplyvňujú tréning a pomocou ktorých môžeme ladiť svoj model.

3.1.0.1 kernel Kernel je možné si predstaviť ako funkciu ktorá bude prikovaná na dáta aby sa nelineárne dáta dali lineárne rozdeliť. Pri lineárnom kerneli sú dáta rozdelené priamkou, nič zvláštne. Ak ale priamka nestačí tak ako kernel sa môže zvoliť polynomiálna funkcia ktorá dokáže dáta rozdeliť krivkou. Pokiaľ ani toto nie je dostačujúce tak sa používa tzv. kernelový trik. V tomto prípade sa ako kernel používa 'rbf' - radiálna bázová funkcia. V prípade kernelového triku ide o transformáciu priestoru tak aby sa dali dáta rozdeliť lineárne hyperrovinami vo viac rozmernom priestore.

3.1.0.2 C Hodnota hyperparametra určuje akýsi rozostup medzi hyperrovinou a dátovými bodmi ktoré sú klasifikované do určenej triedy. čím nižší je C tým väčší rozostup má hyperrovina.

3.1.1 gamma

Tento hyperparameter určuje šírku funkcie kernelu. Inak povedané určuje aký dosah bude mať vzorka pri klasifikácii. Čím nižšia hodnota tým väčší dosah budú mať vzorky.

Pri spomínanej konfigurácii sme dostali výsledky:

```
Results of SVR:  
MSE: 9.077949170492122  
R-Squared: 0.7930420157292597
```

Vidíme že úspešnosť je 79% (podľa metriky R^2). Ďalšia metrika pre nás zaujímavá je MSE čo je priemer reziduálov. Reziduál je rozdiel medzi predpovedanou a reálnou hodnotou. Z toho nám vyplýva že čím menšie je MSE tým lepší bol náš regresor. Interval hodnôt MSE závisí od problému ktorý regresor rieši keďže pracuje s predpovedanými hodnotami a reálnymi výstupmi množiny.

3.2 Grid search

Skúšanie kombinácií týchto hyperparametrov je časovo náročné a monotónne. Pre zjednodušenie práce bol použitý algoritmus GridSearch ktorému zadáme množinu hyperparametrov ktoré on skúša postupne metódou *každý s každým*. Tento proces je dlhotrvajúci no na základe jeho výstupu vieme dostať najlepšiu konfiguráciu hyperparametrov pre náš klasifikátor.

```
1 tuned_parameters = [  
2     {"kernel": ["rbf"], "gamma": [1e-3, 1e-4], "C": [1, 10, 100, 1000]},  
3 ]  
4  
5 scores = ["r2"]  
6  
7 print("# Tuning hyper-parameters for %s" % score)  
8 print()  
9  
10 clf = GridSearchCV(SVR(), param_grid=tuned_parameters, cv=2)  
11 clf.fit(X, y.to_numpy())  
12  
13 print("Best parameters set found on development set:")
```

```

14 print()
15 print(clf.best_params_)
16 print()
17 print("Grid scores on development set:")
18 print()
19 means = clf.cv_results_["mean_test_score"]
20 stds = clf.cv_results_["std_test_score"]
21 for mean, std, params in zip(means, stds, clf.cv_results_["params"]):
22     print("%0.3f (+/-%0.03f) for %r" % (mean, std * 2, params))
23 print()
24
25 print("Detailed classification report:")
26 print()
27 print("The model is trained on the full development set.")
28 print("The scores are computed on the full evaluation set.")
29 print()
30 y_true, y_pred = yt, clf.predict(Xt)

```

Listing 14: GridSearch algoritmus

Výsledok algoritmu GridSearch spolu s cross validation bol nasledovný:

```
# Tuning hyper-parameters for r2
```

```
Best parameters set found on development set:
```

```
{'C': 1000, 'gamma': 0.001, 'kernel': 'rbf'}
```

```
Grid scores on development set:
```

```

0.626 (+/-0.077) for {'C': 1, 'gamma': 0.001, 'kernel': 'rbf'}
0.085 (+/-0.109) for {'C': 1, 'gamma': 0.0001, 'kernel': 'rbf'}
0.777 (+/-0.002) for {'C': 10, 'gamma': 0.001, 'kernel': 'rbf'}
0.627 (+/-0.077) for {'C': 10, 'gamma': 0.0001, 'kernel': 'rbf'}
0.795 (+/-0.017) for {'C': 100, 'gamma': 0.001, 'kernel': 'rbf'}
0.776 (+/-0.002) for {'C': 100, 'gamma': 0.0001, 'kernel': 'rbf'}
0.812 (+/-0.019) for {'C': 1000, 'gamma': 0.001, 'kernel': 'rbf'}
0.792 (+/-0.017) for {'C': 1000, 'gamma': 0.0001, 'kernel': 'rbf'}

```

```
Detailed classification report:
```

```
The model is trained on the full development set.
```

```
The scores are computed on the full evaluation set.
```

Vidíme že sme dostali najlepšie výsledky pri parametroch $C=1000$, $\gamma=0.001$ a kernel: rbf pri ktorých bola úspešnosť 81% (podľa metriky R^2). Do algoritmu môžeme dosadiť aj viacero možností hyperparametrov. Čím viac hyperparametrov zadáme tým dlhšie bude algoritmus bežať no dokážeme teoreticky dostať lepšie výsledky.

3.3 Metóda bagging

Pri metóde bagging vytvárame viacero regresorov ktoré dostanú náhodné podmnožiny tréovacích dát ktorých výsledky sú priemerované do konečného výsledku. Metóda sa dá jednoducho popísať ako kombinácia viacerých slabých regresorov nám dá presnejší model regresora a tým vytvára jeden silný regresor.

```
1 bagging_reg = BaggingRegressor(base_estimator=regressor, n_estimators=27, random_state=255)
2 bagging_reg.fit(X, y.to_numpy())
3 y_pred_bagging = bagging_reg.predict(Xt)
4 mse = metrics.mean_squared_error(yt, y_pred_bagging)
5 r2 = metrics.r2_score(yt, y_pred_bagging)
6 print("Results of sklearn.metrics:")
7 print("MSE:", mse)
8 print("R-Squared:", r2)
```

Listing 15: Metóda bagging

Podľa výsledkou bohužiaľ vidíme že metóda bagging nebola úspešnejšia ako SVC regresor s vyladenými parametrami podľa Grid Search (podľa metriky R^2).

```
Results of bagging:
MSE: 9.065703067088416
R-Squared: 0.793321200909525
```

3.4 Metóda boosting

Podobne ako pri metóde bagging sa používa viacero slabých regresorov ktoré sú ale vyhodnocované sekvenciálne po sebe a každý ďalší dostane upravené parametre z predchádzajúceho. Táto metóda sa dá popísať ako vylepšovanie slabých regresorov na základe prechádzajúcich pre vytvorenie jedného silného regresora.

```
1 boosting_reg = GradientBoostingRegressor(random_state=255)
2 boosting_reg.fit(X, y.to_numpy())
3 y_pred_boosting = boosting_reg.predict(Xt)
```

Listing 16: Metóda boosting

V tomto prípade bola metóda boosting úspešnejšia ako metóda bagging (podľa metriky R^2).

```
Results of boosting:
MSE: 5.122653856629292
R-Squared: 0.8832143586206862
```

Celkovô môžeme zhodnotiť že najlepší regresor s úspešnosťou 88% (podľa metriky R^2) bol pri boosting regresor a to presnejšie Gradient boosting regressor.

3.5 Random forrest

Posledný regresor ktorý bol použitý je RandomForest regresor ktorý si ľahko dokážeme aj vizualizovať.

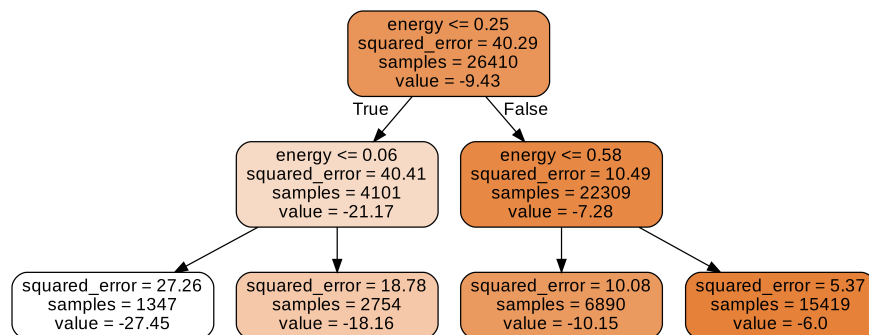

```

1 rf = RandomForestRegressor(max_depth=2, random_state=0)
2 rf.fit(X, y)

```

Listing 17: Random forrest regresor

Každe 42 je odpoveď na základnú otázku života, vesmíru a všetkého tak si zoberieme 42. estimátor (strom nášho náhodného lesa) a vizualizujeme si ho:



Obr. 6: Strom náhodného lesa

Na strome vidíme zároveň že stĺpec energy vo vstupnom datasete má veľkú váhu pri vyhodnocovaní predikovaného výsledku.

Literatúra

- [1] https://scikit-learn.org/stable/auto_examples/model_selection/plot_grid_search_digits.html
- [2] <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.BaggingRegressor.html>
- [3] <https://towardsdatascience.com/how-to-visualize-a-decision-tree-from-a-random-forest-in-python-using-scikit-learn-38ad2d75f21c>
- [4] <https://towardsdatascience.com/ensemble-learning-bagging-boosting-3098079e5422>
- [5] <https://www.analyticsvidhya.com/blog/2020/03/support-vector-regression-tutorial-for-machine-learning/>
- [6] <https://towardsdatascience.com/train-test-split-and-cross-validation-in-python-80b61beca4b6>
- [7] <https://www.datatechnotes.com/2019/10/accuracy-check-in-python-mae-mse-rmse-r.html>
- [8] <https://stackoverflow.com/questions/23199796/detect-and-exclude-outliers-in-pandas-data-frame>