

Charles Deliege - Arthur Clavier - DIA2

Python Project

News Popularity Prediction

01 - Modules utilisés

Numpy - Pandas - Plotly.express
Seaborn -Matplotlib

```
import numpy as np
import pandas as pd
import plotly.express as px
import seaborn as sns
import matplotlib.pyplot as plt
```

02 - Le Dataset

61 Colonnes

dtypes : float64(59), int64(1), object(1)

39644 Rows, 18.5 MB

	url	timedelta	n_tokens_title	n_tokens_content	n_unique_tokens	n_non_stop_words	n_non_stop_unique_tokens	num_hrefs	lang
0	http://mashable.com/2013/01/07/amazon-instant... http://mashable.com/2013/01/07/samsung-spons... http://mashable.com/2013/01/07/apple-40-billion... http://mashable.com/2013/01/07/astronaut-notre... http://mashable.com/2013/01/07/att-universe-apps/	731.0 731.0 731.0 731.0 731.0	12.0 9.0 9.0 9.0 13.0	219.0 255.0 211.0 531.0 1072.0	0.663594 0.604743 0.575130 0.503788 0.415646	1.0 1.0 1.0 1.0 1.0	0.815385 0.791946 0.663866 0.665635 0.540890	4 3 3 3 19	en en en en en
1									
2									
3									
4									

5 rows × 61 columns

Data Preprocessing

Let's separate values that contains outliers from the other with only few possible values (binary columns)

```
low_number_of_value=[]
for col in num.columns:
    if np.abs(df[col].skew()) > 1 and df[col].nunique()<5:
        low_number_of_value.append(col)

low_number_of_value
```

```
['data_channel_is_lifestyle',
'data_channel_is_entertainment',
'data_channel_is_bus',
'data_channel_is_socmed',
'data_channel_is_tech',
'data_channel_is_world',
'weekday_is_monday',
'weekday_is_tuesday']

truely_continuous_value=[]
plt.figure(figsize=(20,20))
for col in num.columns:
    if np.abs(df[col].skew()) > 1 and df[col].nunique()>15
        truly_continuous_value.append(col)

truely_continuous_value
```

```
['n_tokens_content',
'n_unique_tokens',
'n_non_stop_words',
'n_non_stop_unique_tokens',
'num_hrefs',
'num_self_hrefs',
'num_imgs',
'num_videos',
```

Select columns

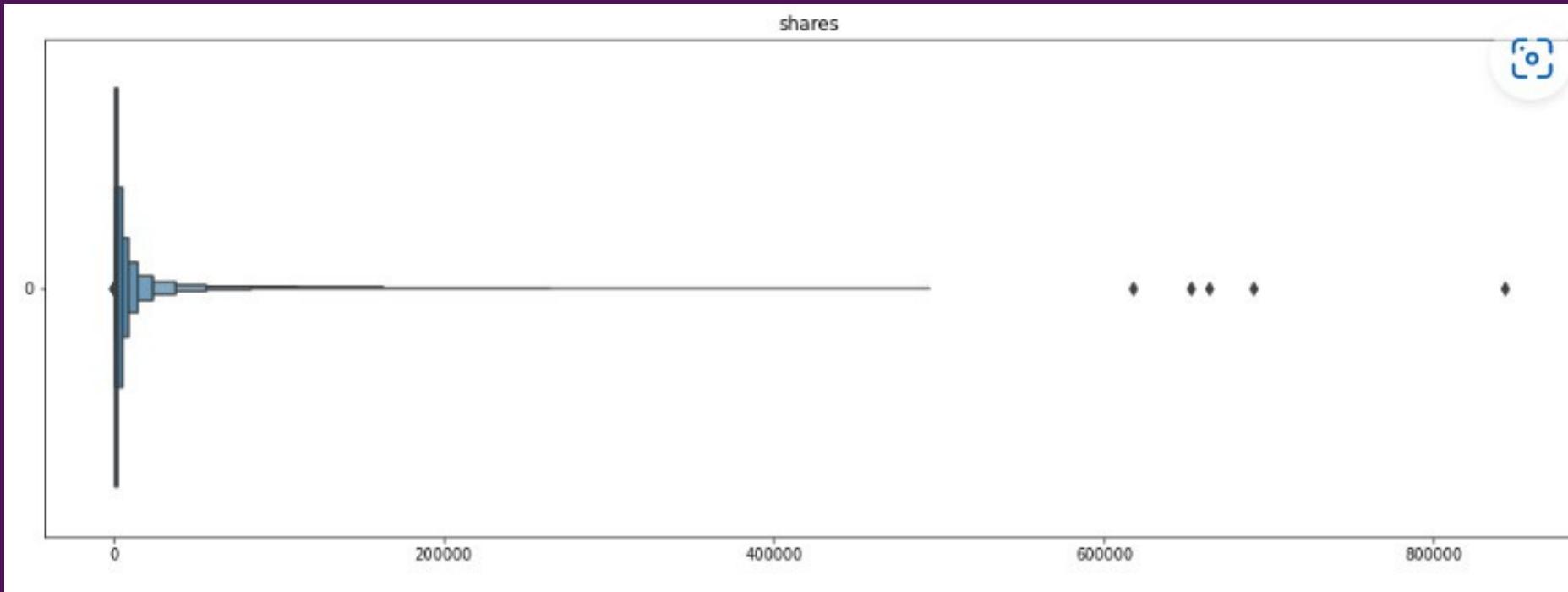
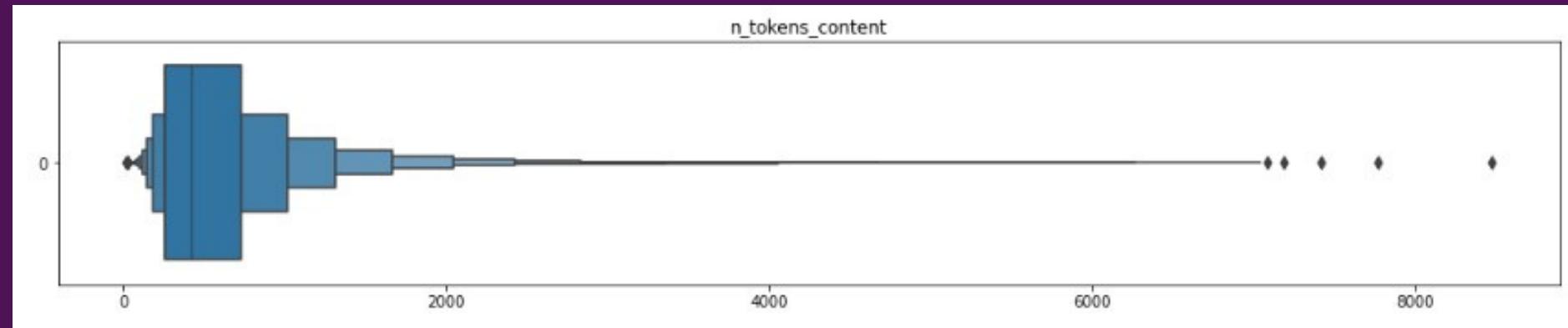
No null data

```
number_zero=df.isnull().sum().sum()
number_Nan=df.isna().sum().sum()
number_zero,number_Nan
```

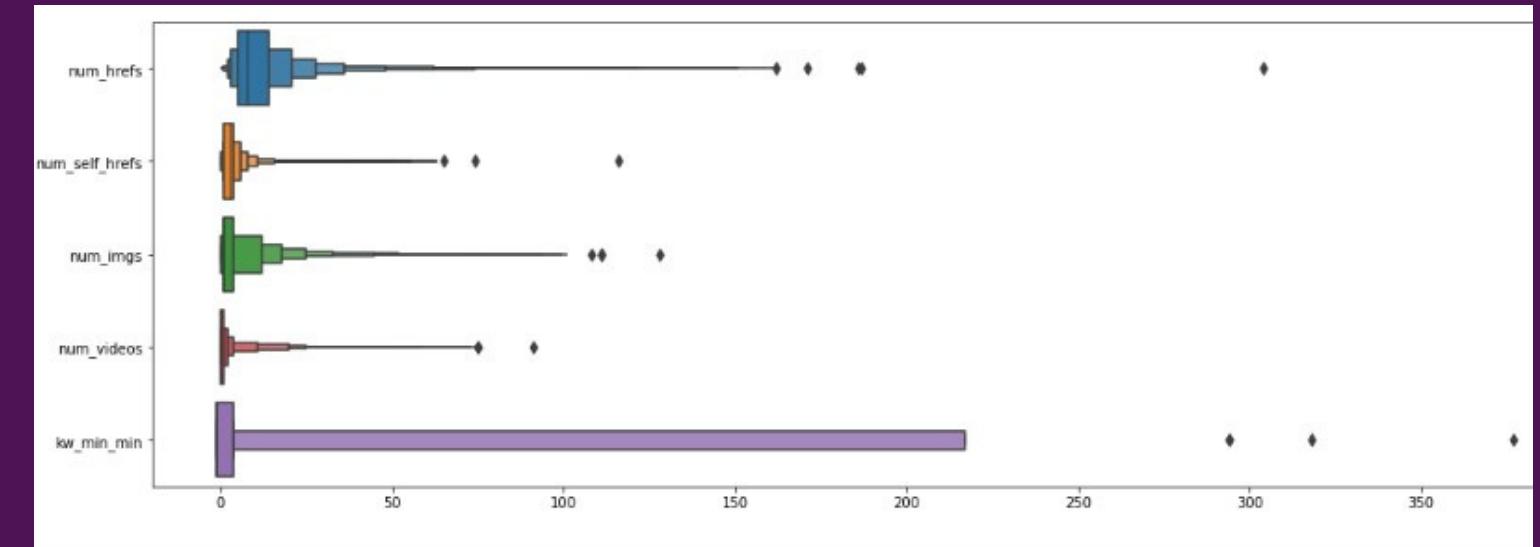
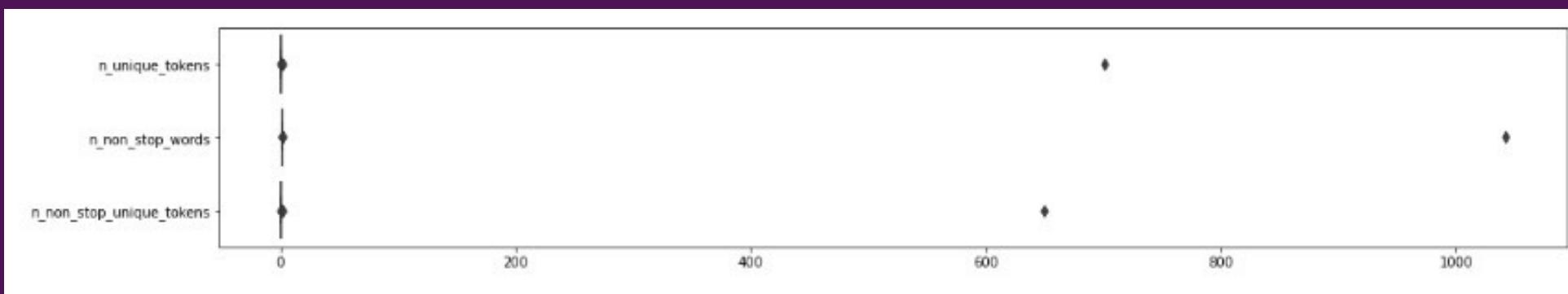
(0, 0)

Data Preprocessing

Outlier Check



We will use the shares to create our target variable



We could easily remove these outliers

Popularity variable

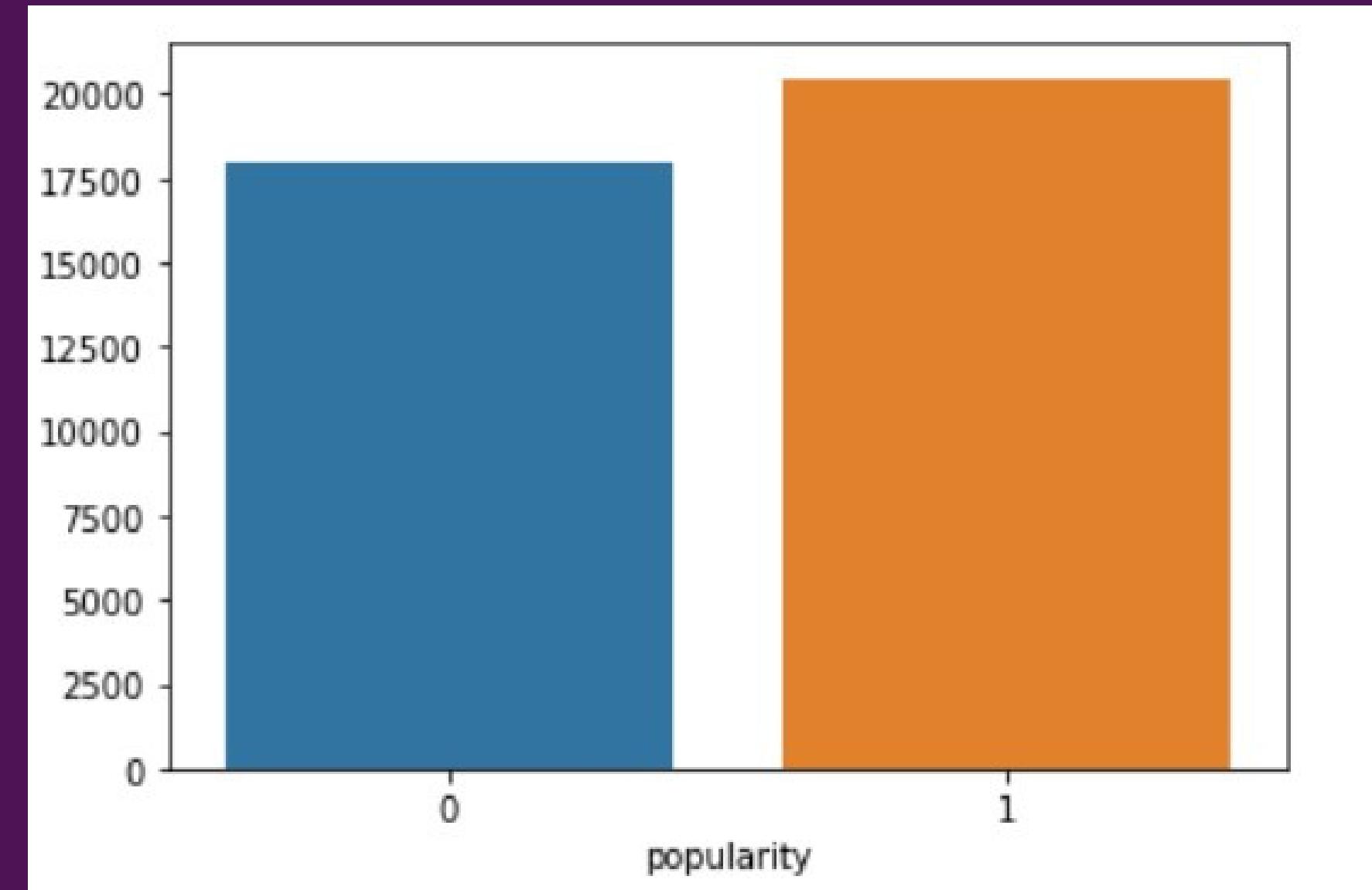
Creation of the target variable

```
df['shares'].describe()

count    38463.000000
mean     3355.360398
std      11585.968776
min      1.000000
25%     945.000000
50%    1400.000000
75%    2700.000000
max   843300.000000
Name: shares, dtype: float64
```

```
df['popularity']=df['shares'].apply(lambda x: 1 if x>=1400 else 0)
df['popularity'].value_counts()

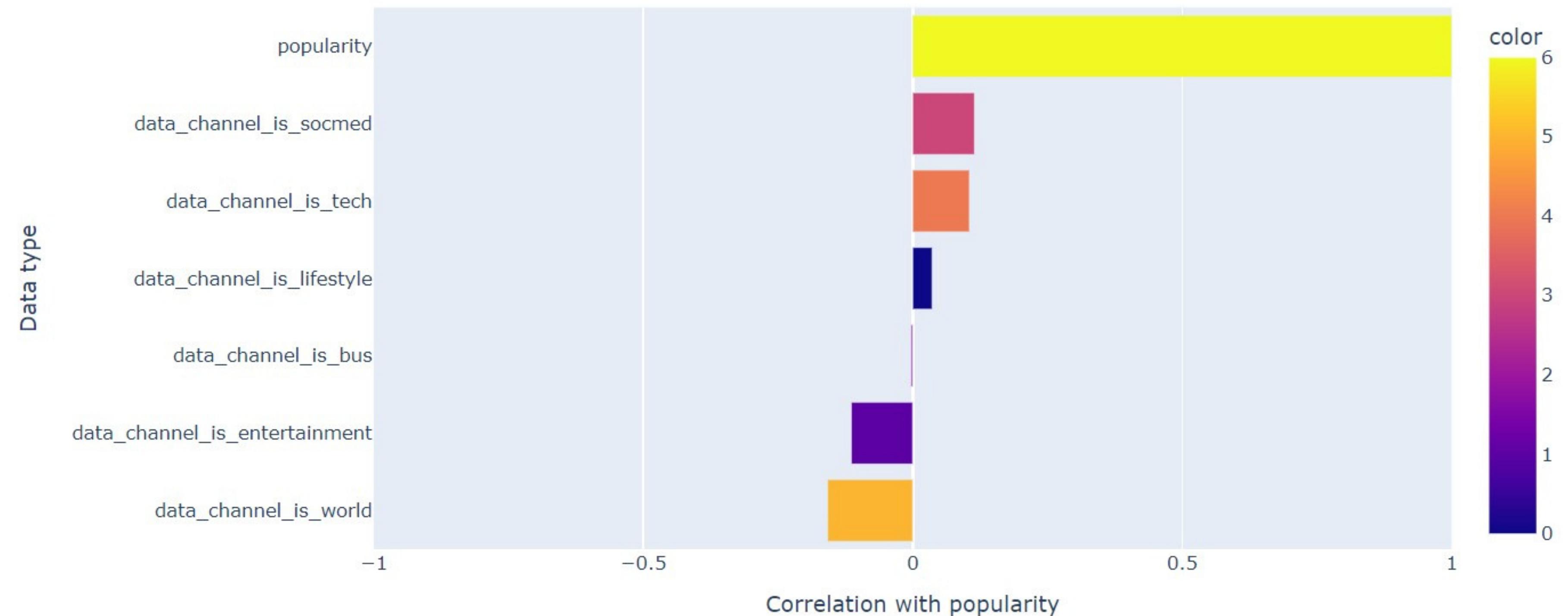
1    20464
0    17999
```



Correlations

Popularity and subject

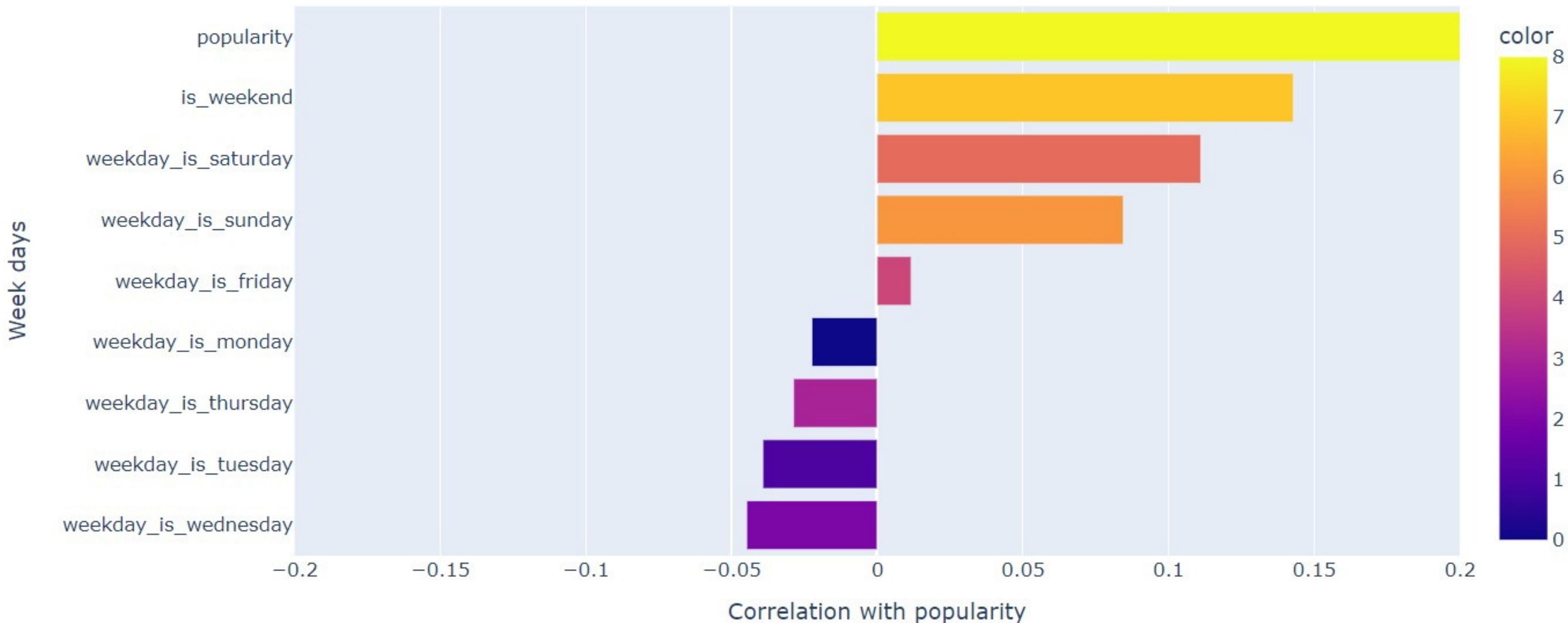
Correlation between the popularity and the subject of articles



Correlations

Correlation between week days and popularity

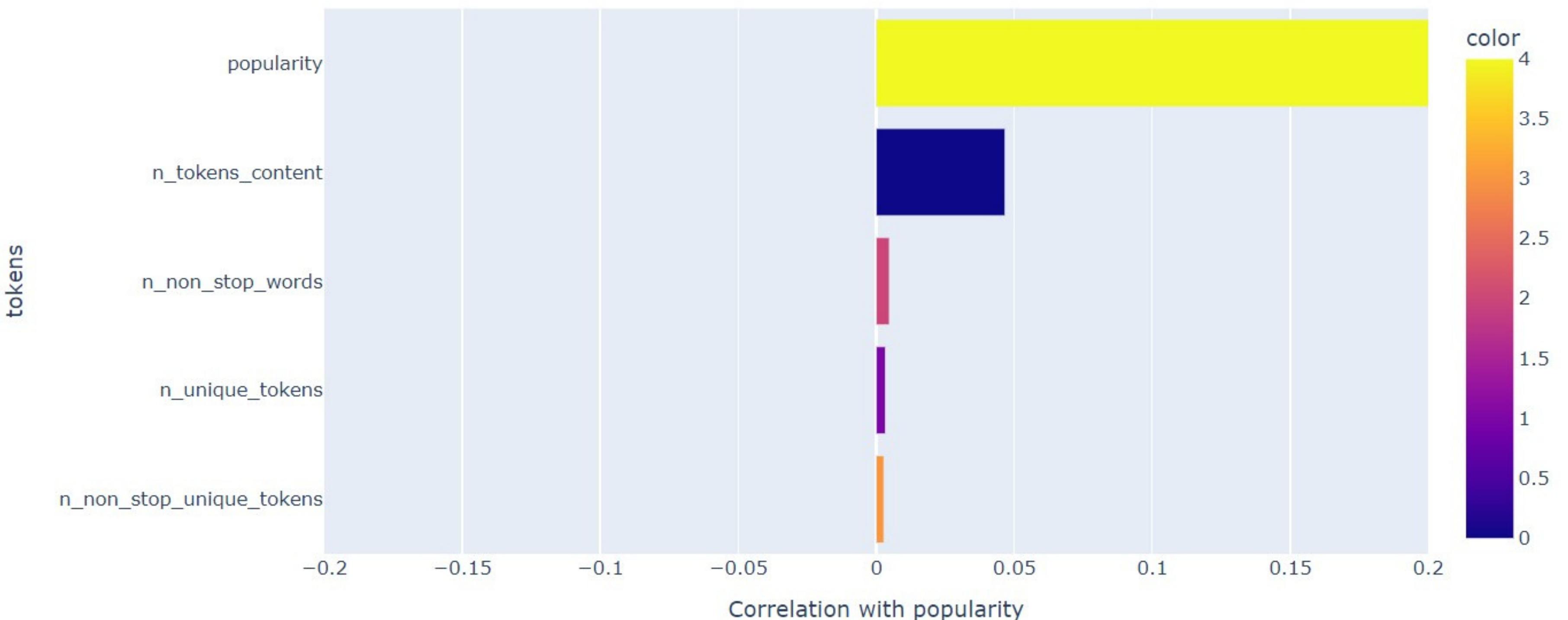
Popularity and weekdays



Correlations

Popularity and tokens

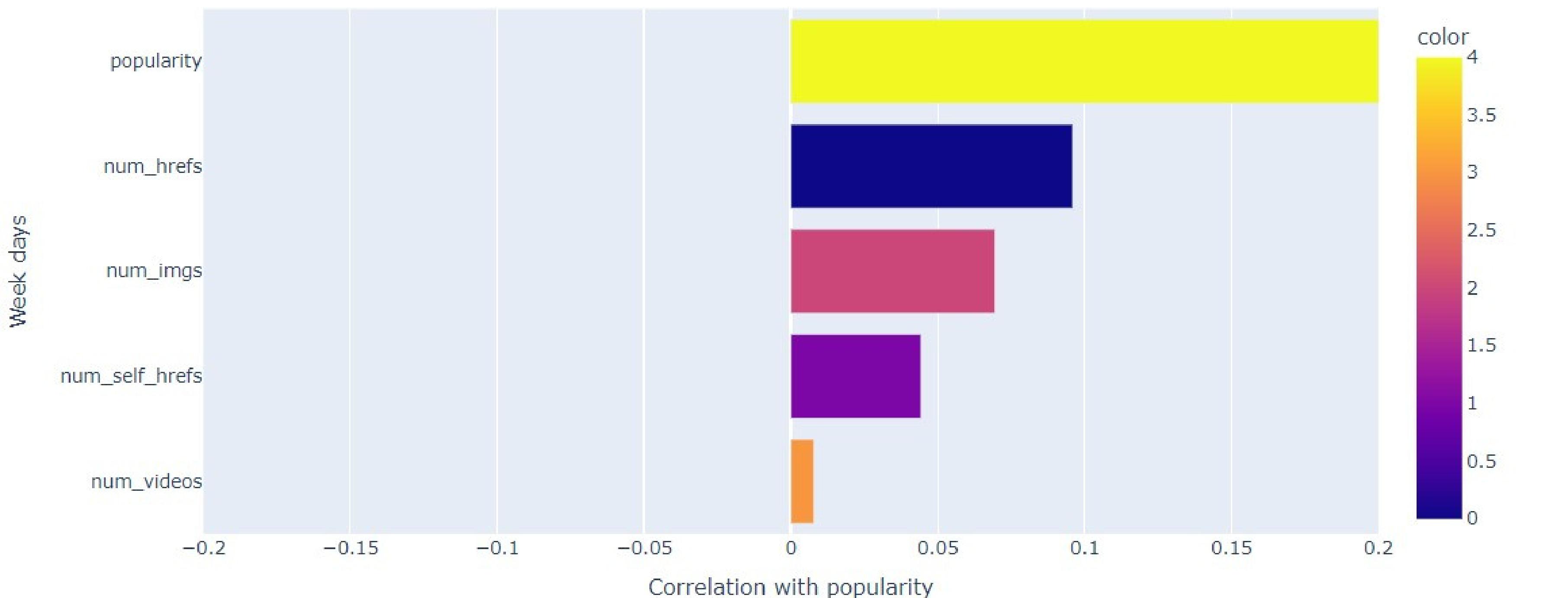
Correlation between tokens and popularity



Correlations

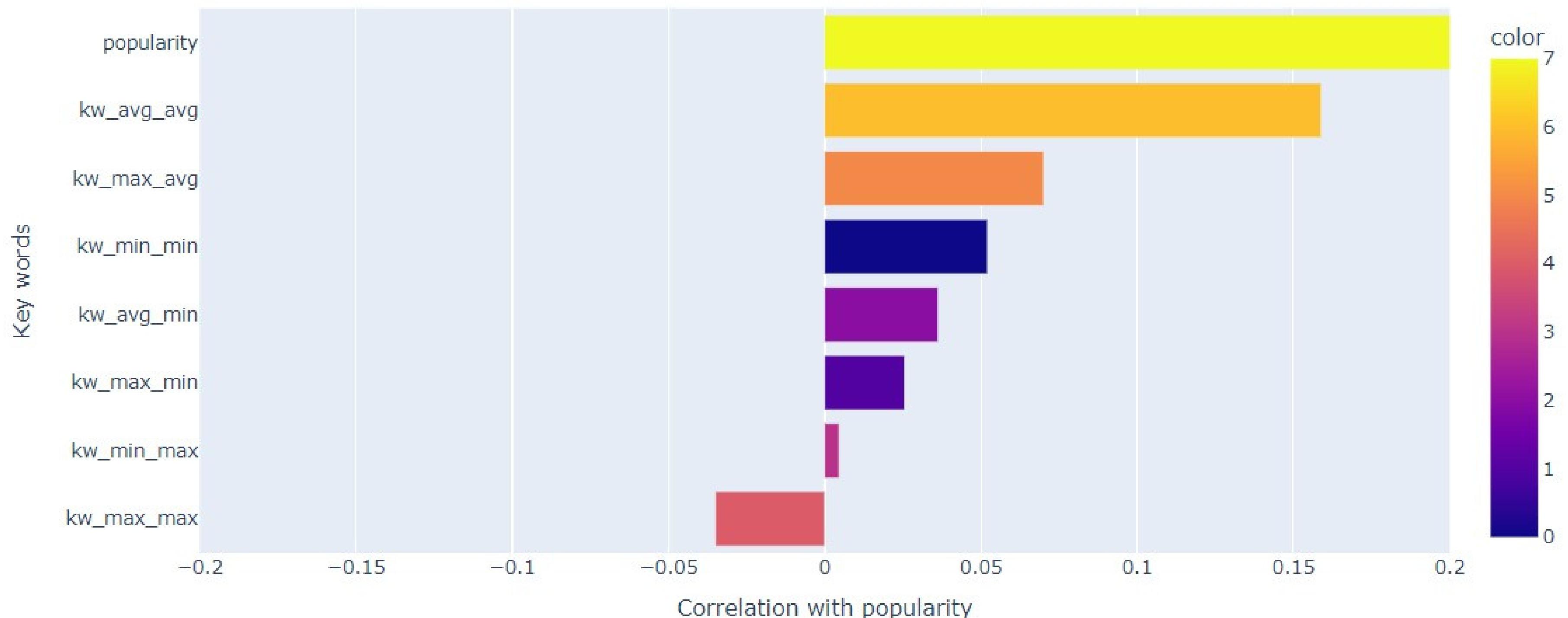
Popularity and number of attached objects

Correlation between number of objects attached and popularity



Correlations -

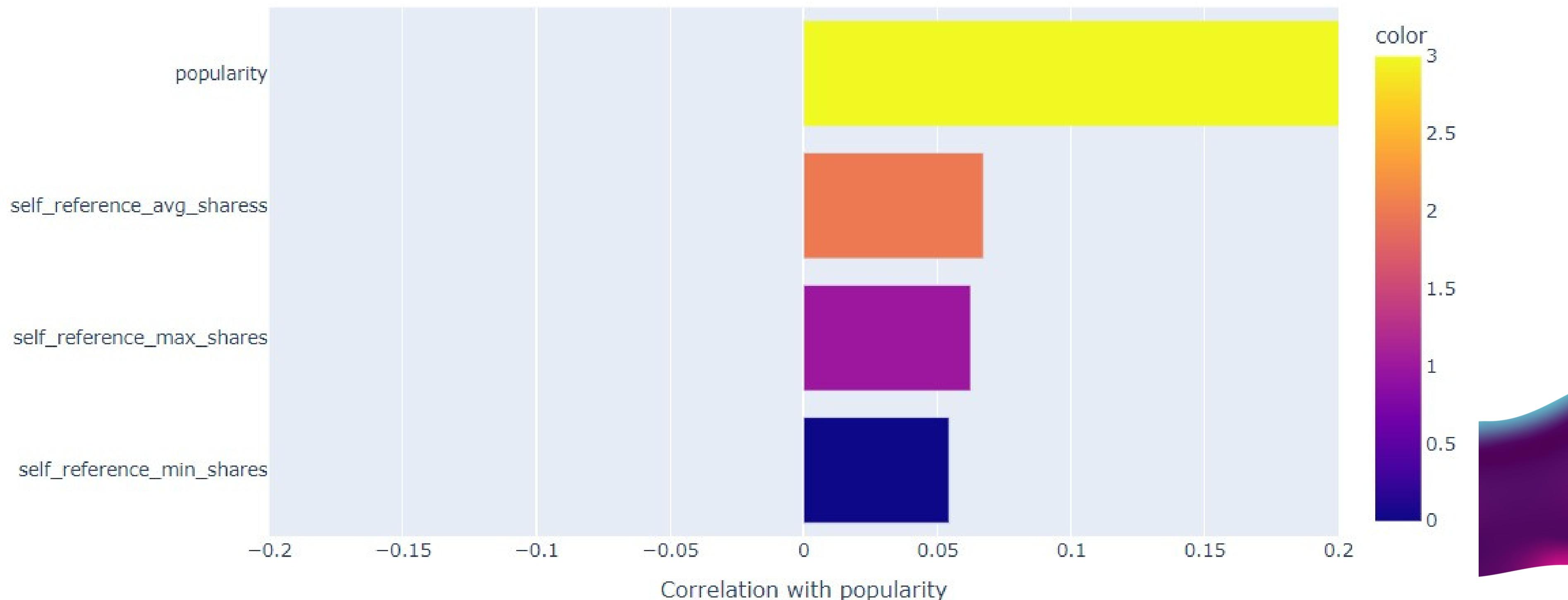
Popularity and key words



Correlations -

Correlation between share reference stats and popularity

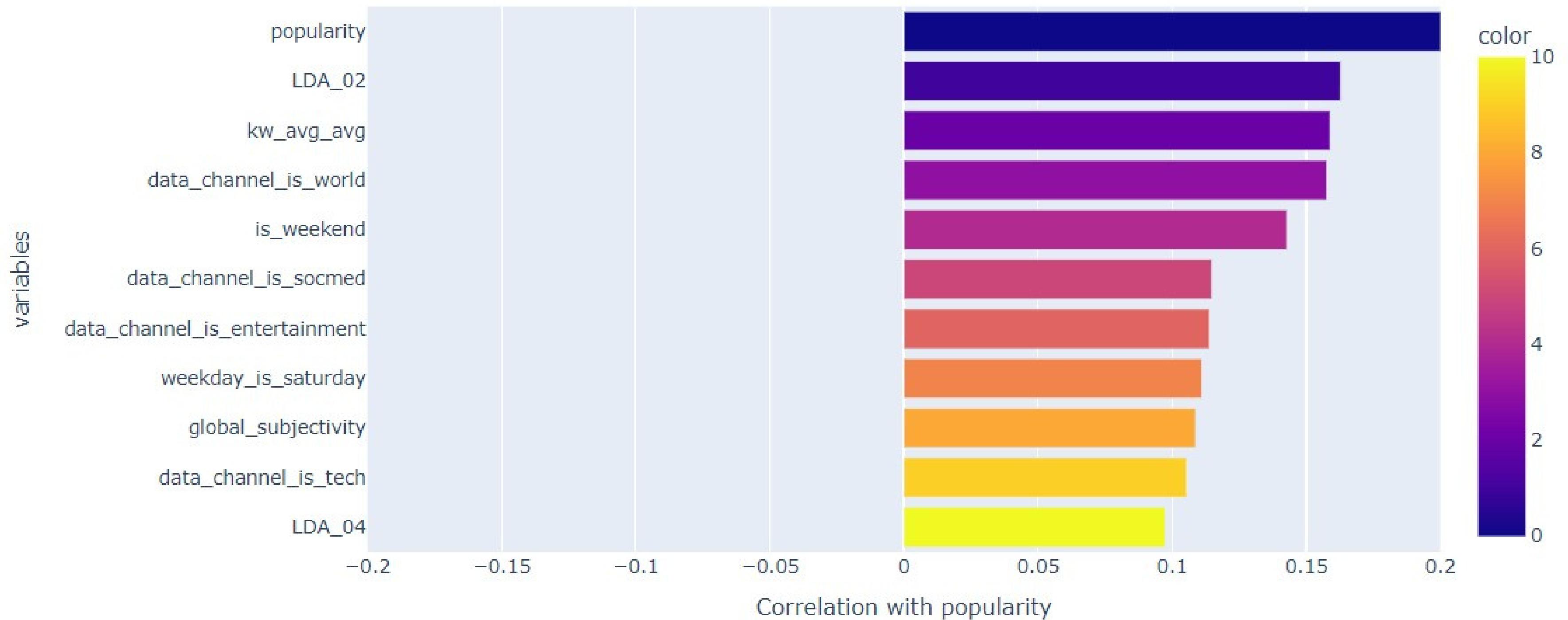
Popularity and share reference stats



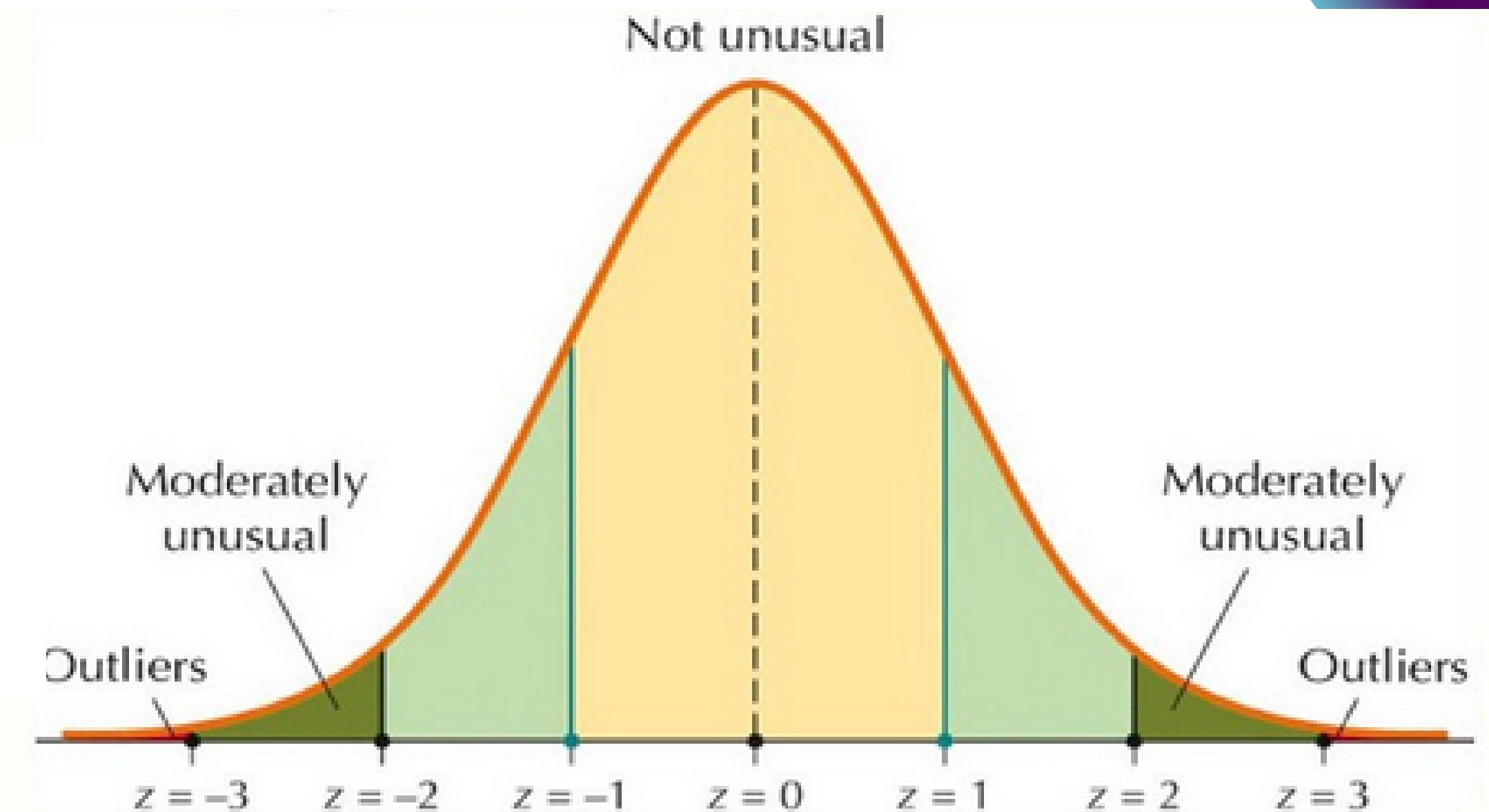
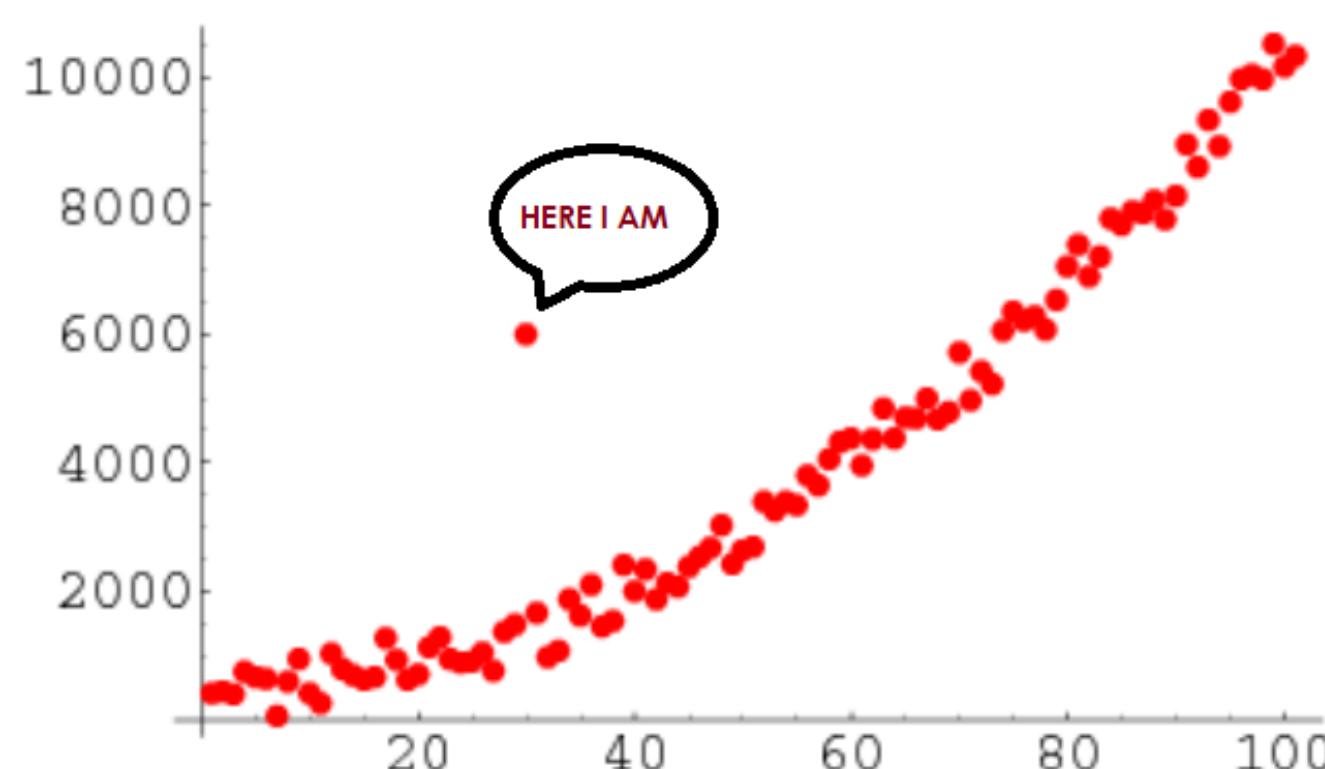
Correlations

Most correlated variables

Top 10 correlation with popularity



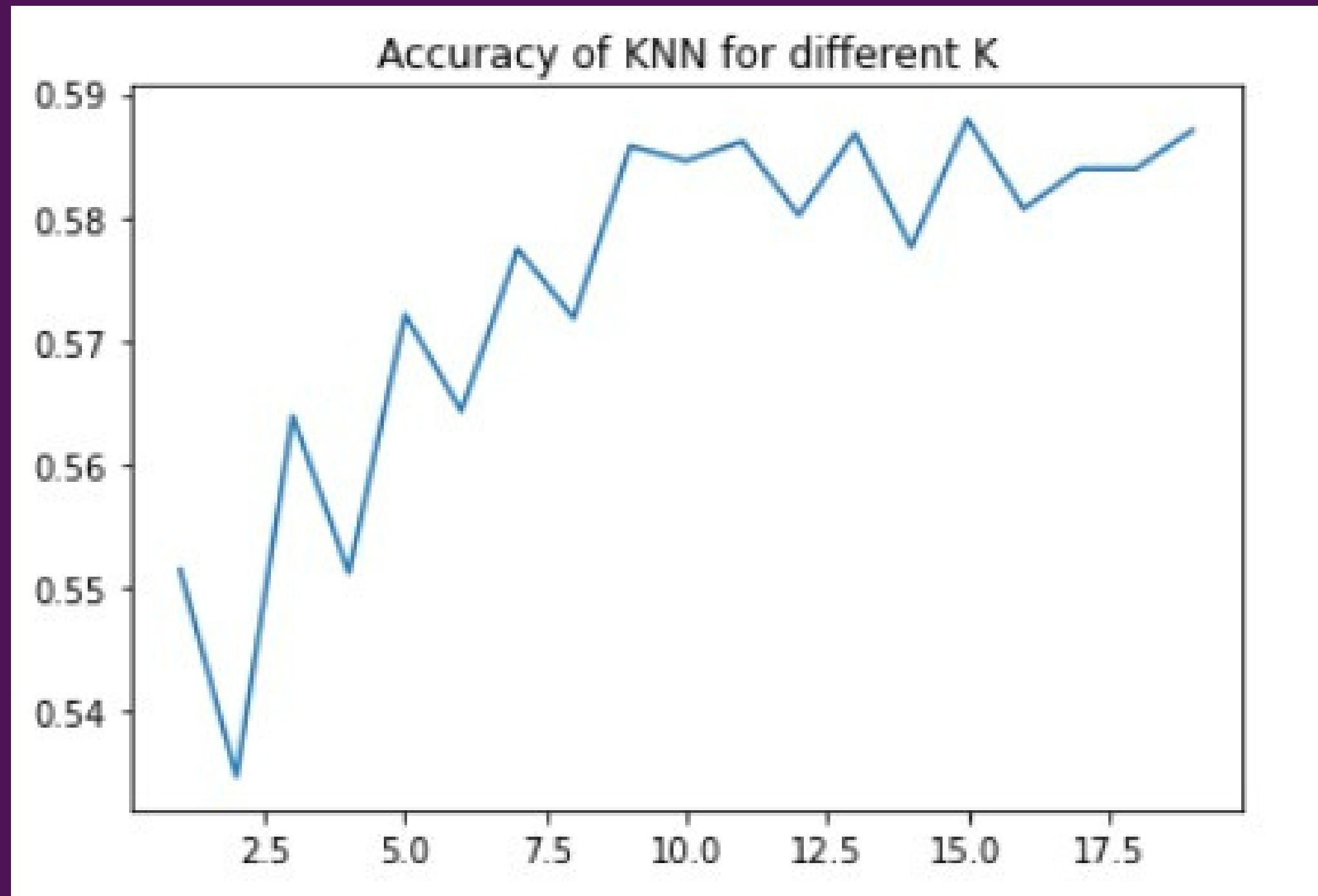
Outliers Removal and data split



Application models

(All dataset without outliers)

KNN



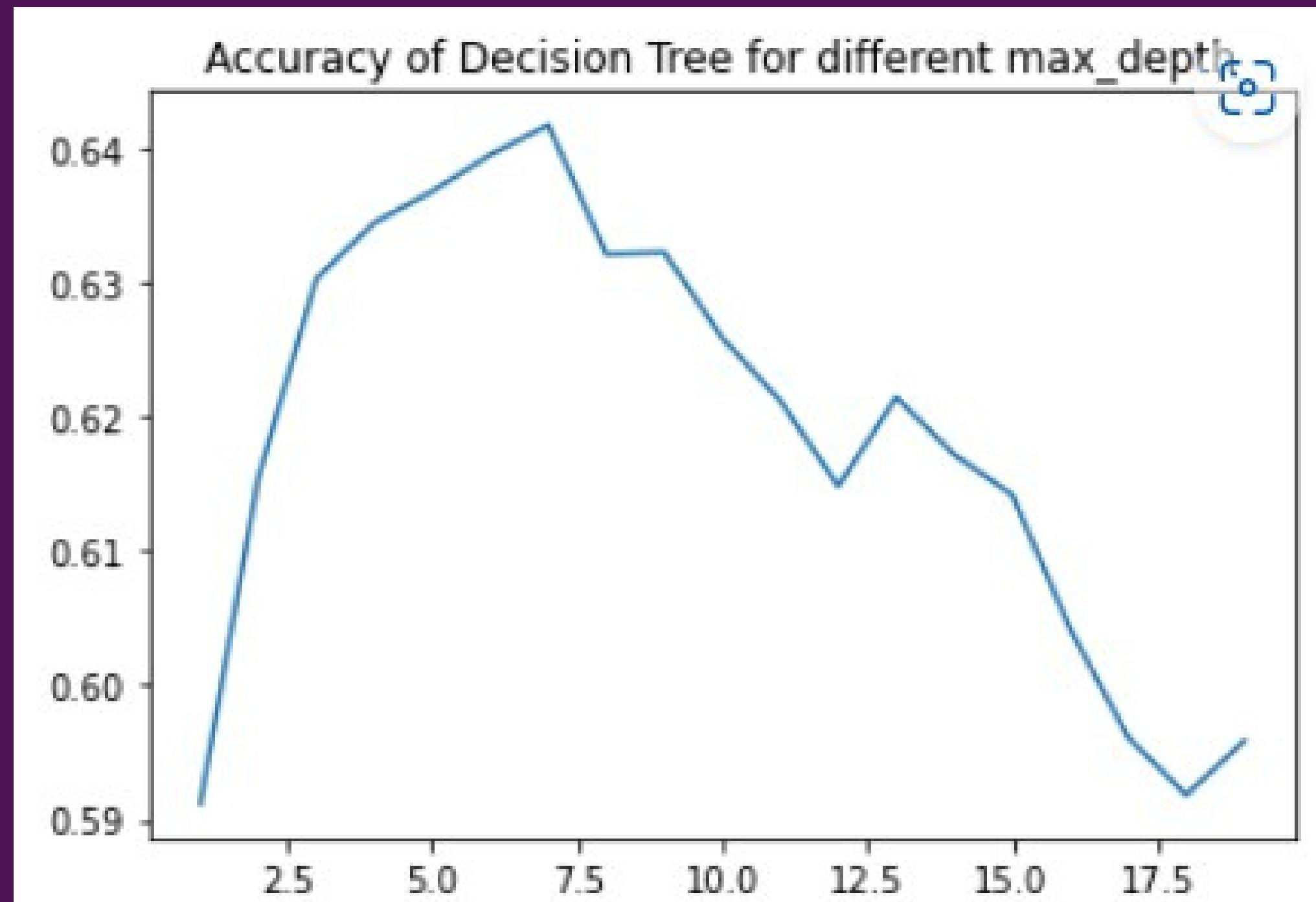
With default parameters
Accuracy = 0.561

With best K :
Accuracy = 0.575

Application models

(All dataset without outliers)

Decision Tree



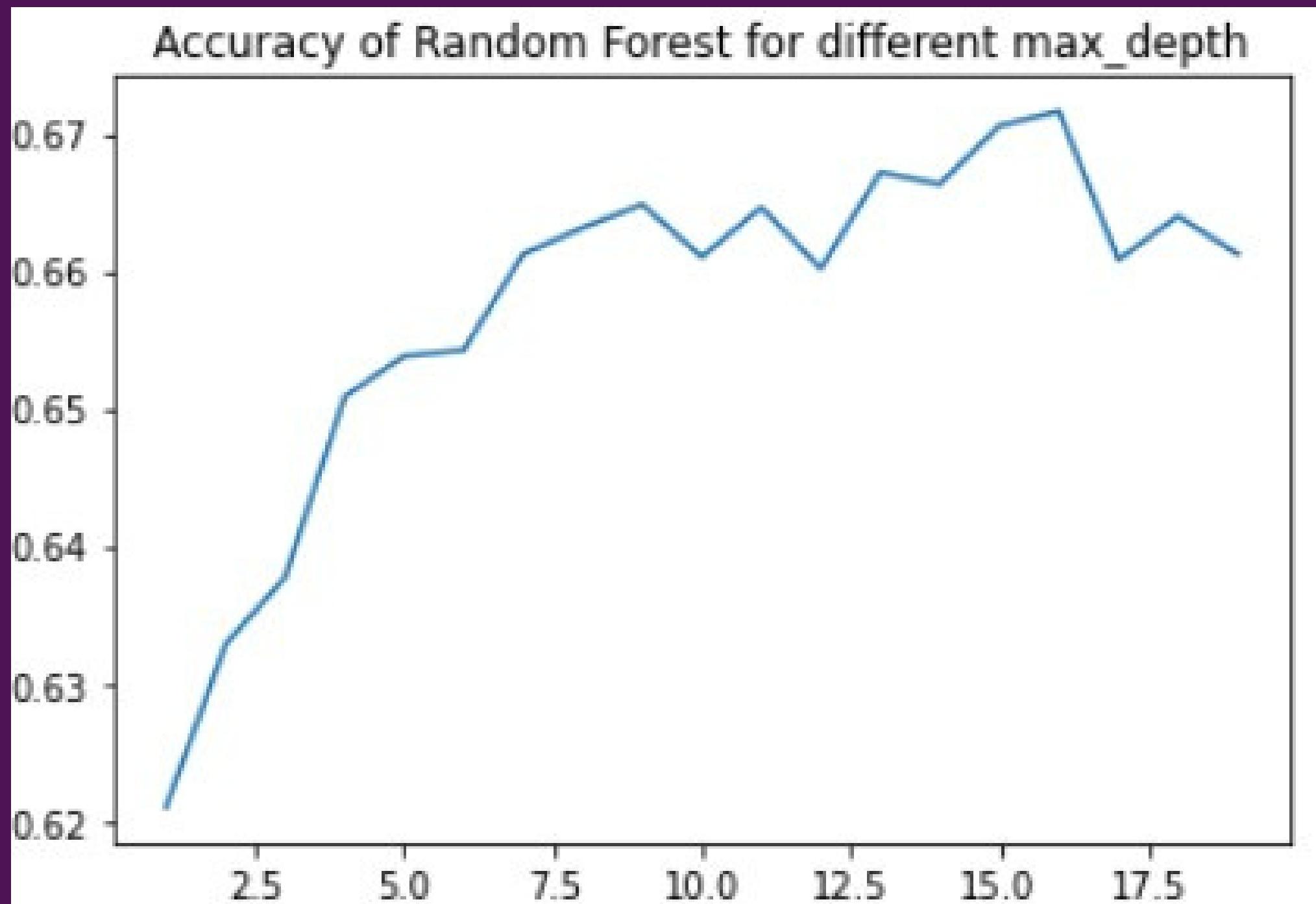
With default parameters
Accuracy = 0.583

With best max_depth :
Accuracy = 0.643

Application models

(All dataset without outliers)

Random Forest



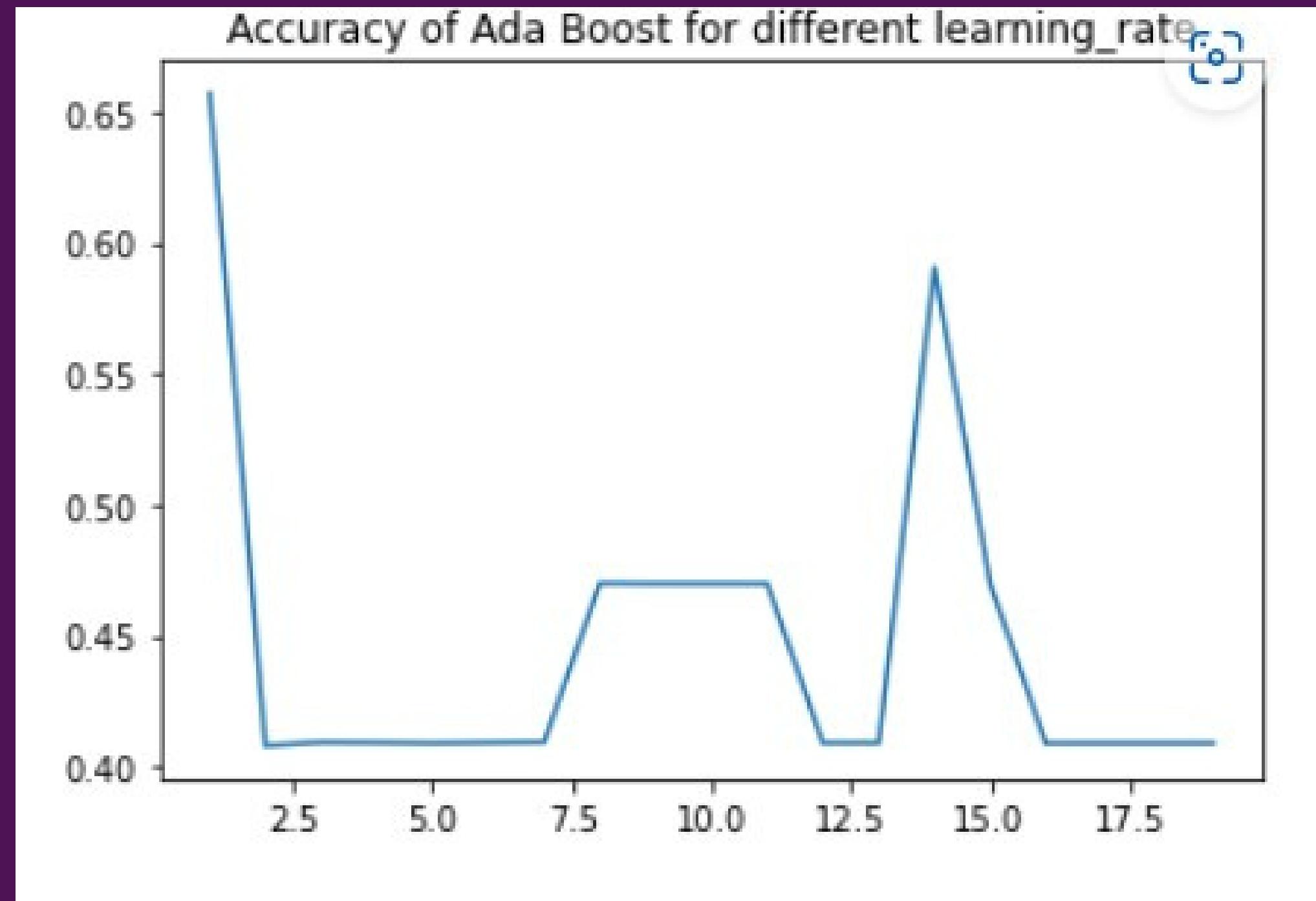
With default parameters
Accuracy = 0.632

With best max_depth :
Accuracy = 0.671

Application models

(All dataset without outliers)

Ada Boost



With default parameters

Accuracy = 0.665

With best learning rate :

Accuracy = 0.665

Application models Summary

(All dataset without outliers)

Model	Default Accuracy	Best Accuracy
Random forest	63%	67%
Ada boost	66%	66%
Desicion tree	58%	64%
Knn	56%	57%

Model	Accuracy	Precision	Recall	F1	AUC
Random Forest (RF)	0.67	0.67	0.71	0.69	0.73
Adaptive Boosting (AdaBoost)	0.66	0.68	0.67	0.67	0.72
Support Vector Machine (SVM)	0.66	0.67	0.68	0.68	0.71
K-Nearest Neighbors (KNN)	0.62	0.66	0.55	0.60	0.67
Naïve Bayes (NB)	0.62	0.68	0.49	0.57	0.65

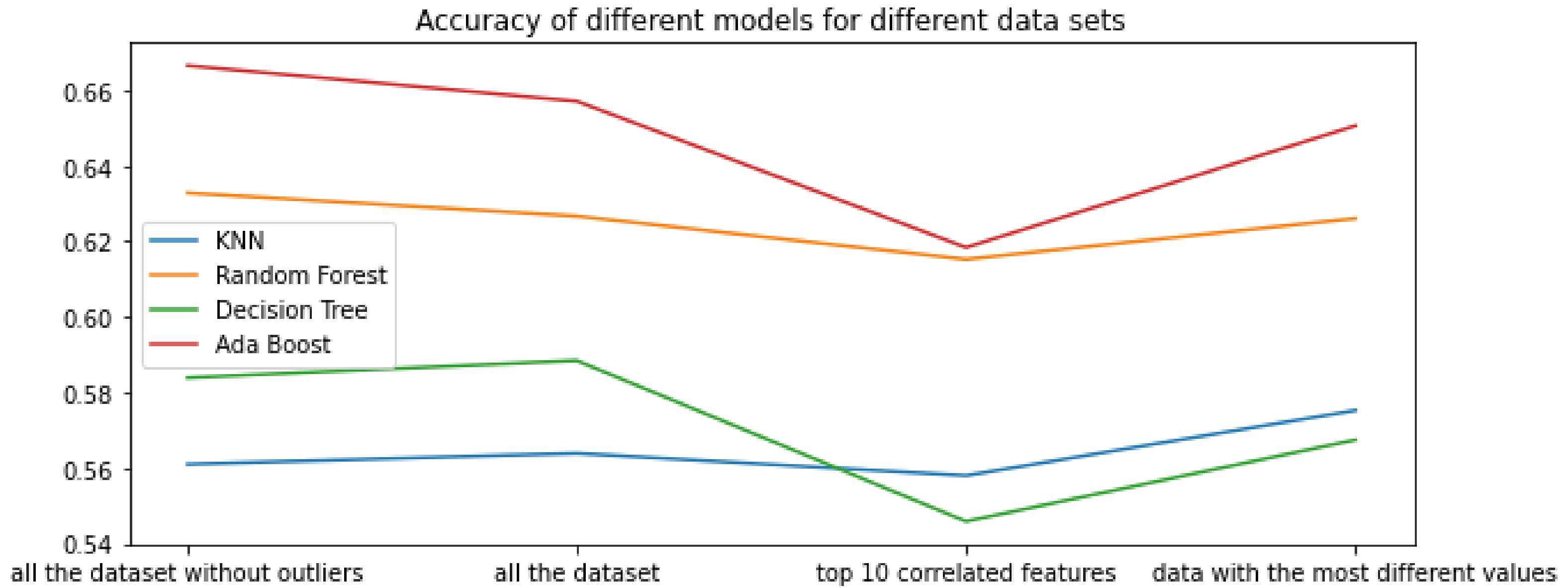
Application models Summary

(Changes of dataset, default parameters)

Model	All without outlier	All with outliers	High Correlation Variables	Most different values
Random forest	63.2%	62.6%	60.6%	62.3%
Ada boost	66.6%	65.7%	64%	64.7%
Desicion tree	58.3%	58.8%	56.1%	56.2%
Knn	56.1%	56.3%	56.5%	57.4%

Application models Summary

(Changes of dataset, default parameters)



The background features a dark purple gradient with three large, semi-transparent circles. One circle is centered in the middle-right area, another is in the bottom-left corner, and a third is in the top-right corner. All three circles overlap each other.

End