

Python for data analysis

MAQUET François-Xavier

MOSER Paul

NESSI Victor

I. Data pre-processing

- Check for NA

```
1 df[df.isna().any(axis=1)]
```

```
url    timedelta  n_tokens_title  n_tokens_content  n_unique_tokens  n_non_stop_words
```

- Normalisation of column name

```
1 column_headers = list(df.columns.values)
2 print("The Column Header :", column_headers)
```

```
The Column Header : ['url', 'timedelta', 'n_tokens_title', 'n_tokens_content', 'n_unique_tokens', 'n_non_stop_words', 'n_non_stop_unique_tokens', 'num_hrefs', 'num_self_hrefs', 'num_imgs', 'num_videos', 'average_token_length', 'num_keywords', 'data_channel_is_lifestyle', 'data_channel_is_entertainment', 'data_channel_is_bus', 'data_channel_is_socmed', 'data_channel_is_tech', 'data_channel_is_world', 'kw_min_min', 'kw_max_min', 'kw_avg_min', 'kw_min_max', 'kw_max_max', 'kw_avg_max', 'kw_min_avg', 'kw_max_avg', 'kw_avg_avg', 'self_reference_min_shares', 'self_reference_max_shares', 'self_reference_avg_shares', 'weekday_is_monday', 'weekday_is_tuesday', 'weekday_is_wednesday', 'weekday_is_thursday', 'weekday_is_friday', 'weekday_is_saturday', 'weekday_is_sunday', 'is_weekend', 'LDA_00', 'LDA_01', 'LDA_02', 'LDA_03', 'LDA_04', 'global_subjectivity', 'global_sentiment_polarity', 'global_rate_positive_words', 'global_rate_negative_words', 'rate_positive_words', 'rate_negative_words', 'avg_positive_polarity', 'min_positive_polarity', 'max_positive_polarity', 'avg_negative_polarity', 'min_negative_polarity', 'max_negative_polarity', 'title_subjectivity', 'title_sentiment_polarity', 'abs_title_subjectivity', 'abs_title_sentiment_polarity', 'shares']
```

```
df.columns = df.columns.str.replace(' ', '')
```

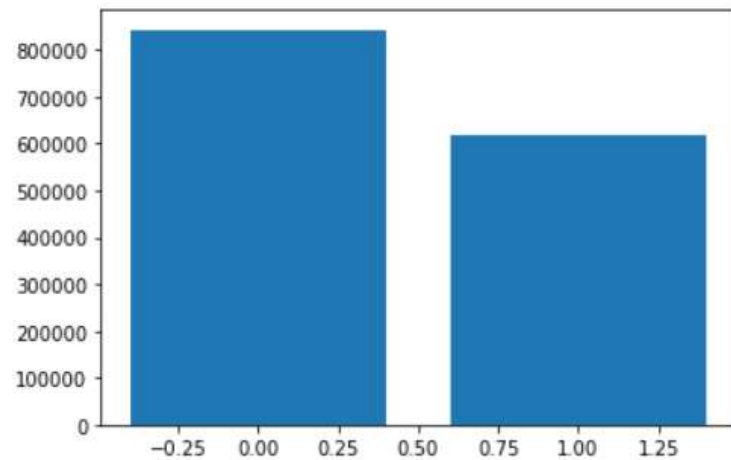
I. Data pre-processing

- Imputation

is_weekend

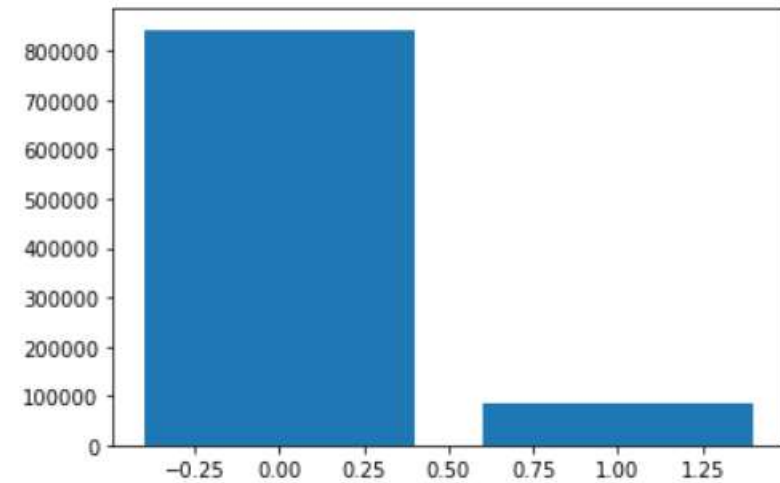
```
1 plt.bar(df['weekday_is_saturday'],df['shares'])
```

<BarContainer object of 39644 artists>



```
1 plt.bar(df['weekday_is_sunday'],df['shares'])
```

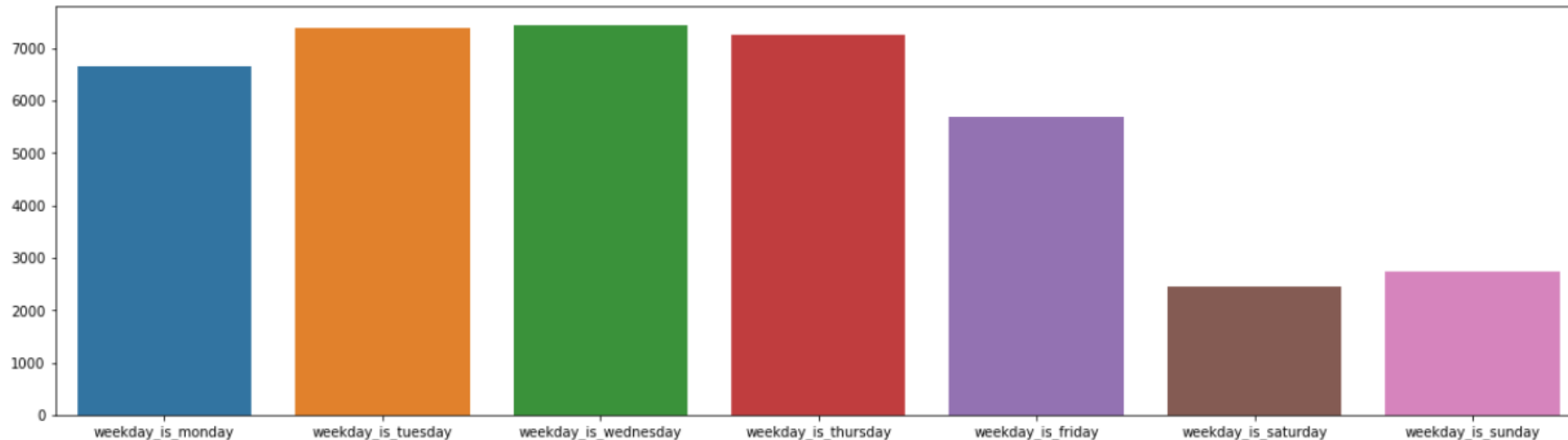
<BarContainer object of 39644 artists>



I. Data pre-processing

- Optimization of days' columns

```
1 dh = df[['weekday_is_monday','weekday_is_tuesday', 'weekday_is_wednesday', 'weekday_is_thursday', 'weekday_is_friday', 'week  
2 x = dh.index  
3 y = dh[['weekday_is_monday','weekday_is_tuesday', 'weekday_is_wednesday', 'weekday_is_thursday', 'weekday_is_friday', 'weekd  
4 plt.rcParams["figure.figsize"] = (20,5.5)  
5 sns.barplot(x,y)  
6 plt.show()
```



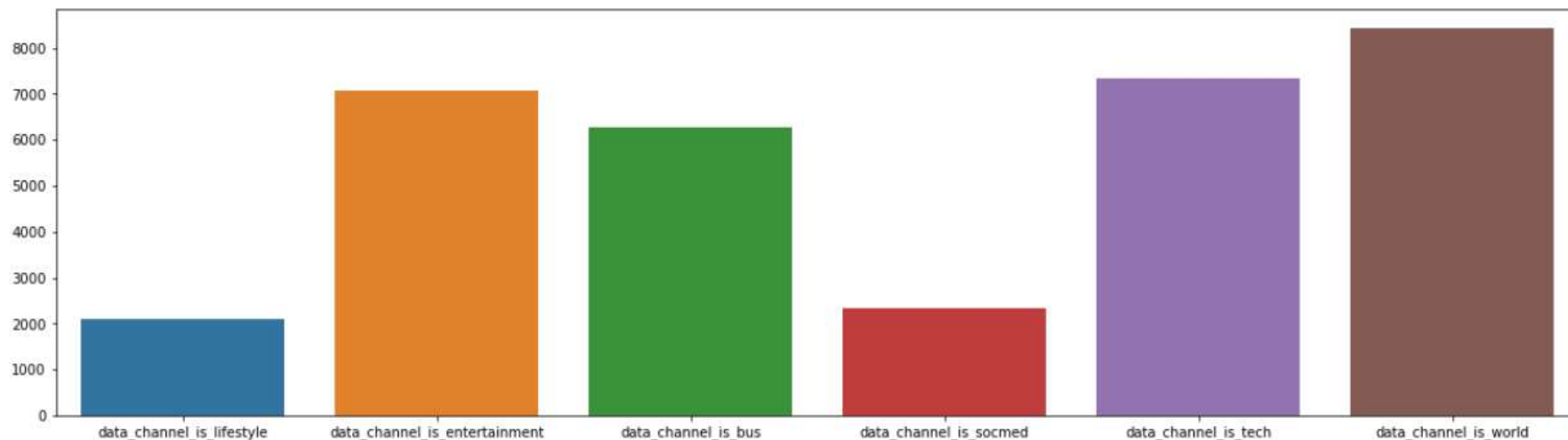
Let's optimize these days' column by creating a new column with the same amount of data

```
1 df['day'] = df['weekday_is_monday']+2*df['weekday_is_tuesday']+3*df['weekday_is_wednesday']+4*df['weekday_is_thursday']+5*df  
2 df.drop(columns=["weekday_is_monday","weekday_is_tuesday","weekday_is_wednesday","weekday_is_thursday","weekday_is_friday",
```

I. Data pre-processing

- Optimization of domain columns

```
1 dh = df[['data_channel_is_lifestyle','data_channel_is_entertainment', 'data_channel_is_bus', 'data_channel_is_socmed', 'data_
2 x = dh.index
3 y = dh[['data_channel_is_lifestyle','data_channel_is_entertainment', 'data_channel_is_bus', 'data_channel_is_socmed', 'data_
4 plt.rcParams["figure.figsize"] = (20,5.5)
5 sns.barplot(x,y)
6 plt.show()
```



Let's optimize these days' column by creating a new column with the same amount of data

```
1 df['domain'] = df['data_channel_is_lifestyle']+2*df['data_channel_is_entertainment']+3*df['data_channel_is_bus']+4*df['data_
2 df
```

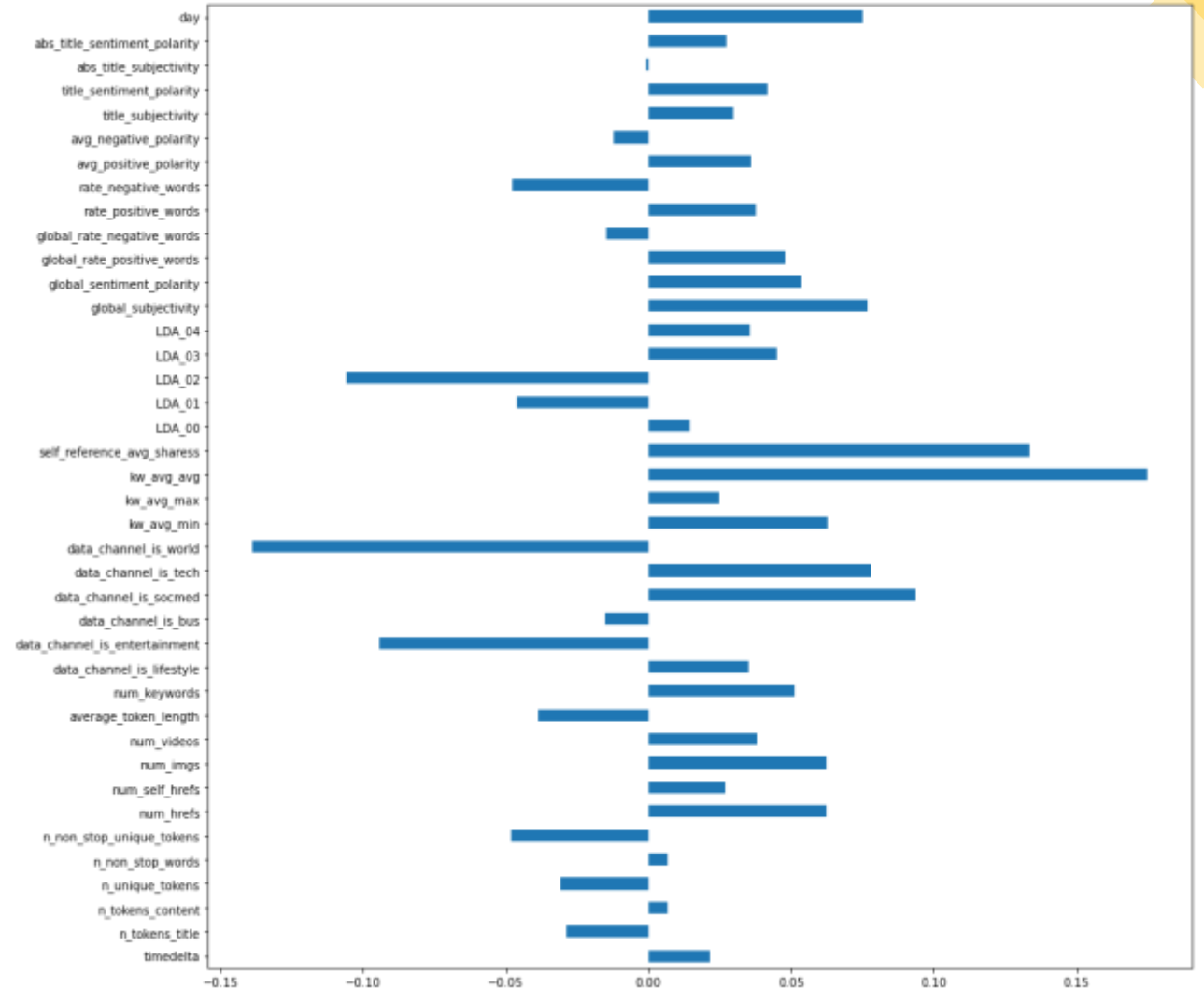

Plot the correlations

[illegible]

II. Visualization

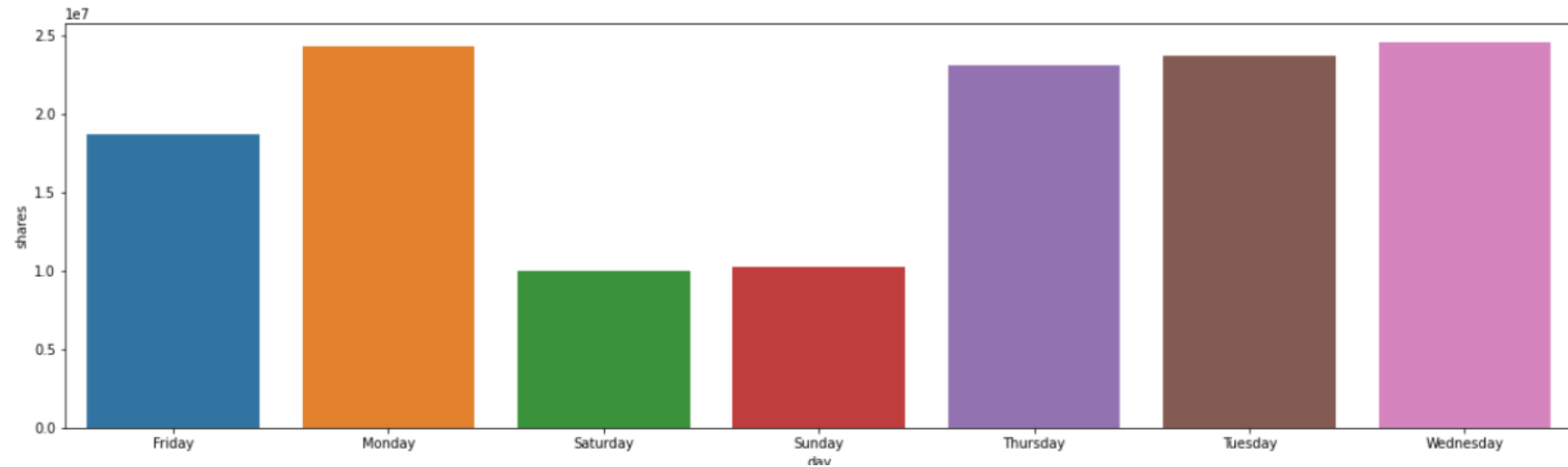
```
dcorr = df.corr(method='kendall')
```

```
dcorr.drop(columns='shares', inplace = True)  
line = dcorr.loc['shares']  
line.plot.barh(figsize=(15,15))
```

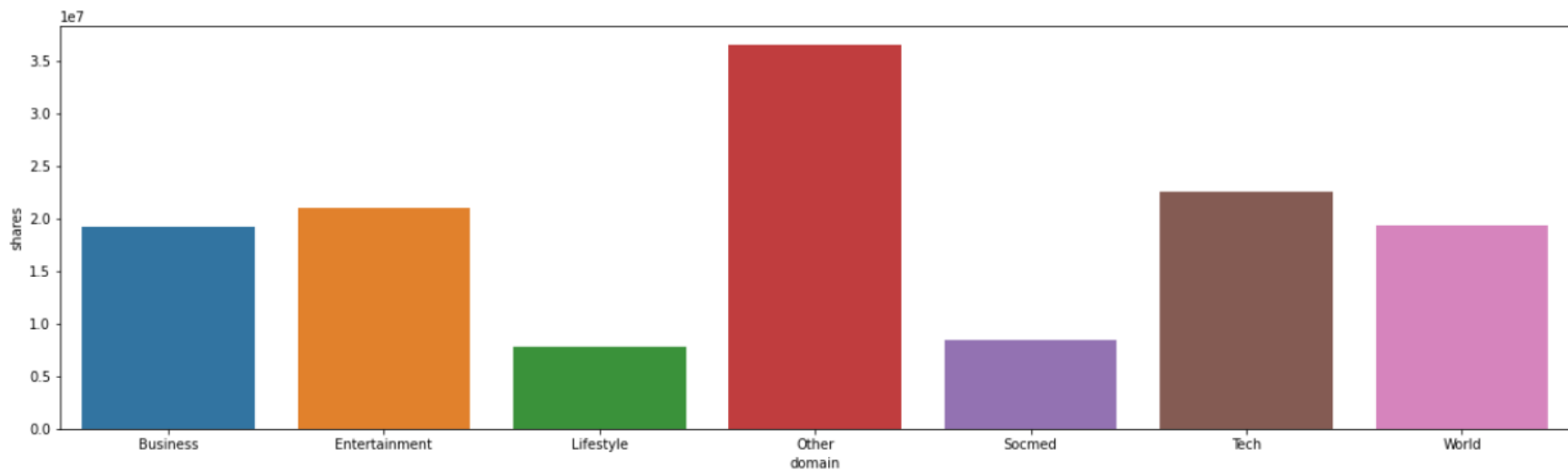


II. Visualization

- Amount of shares by day

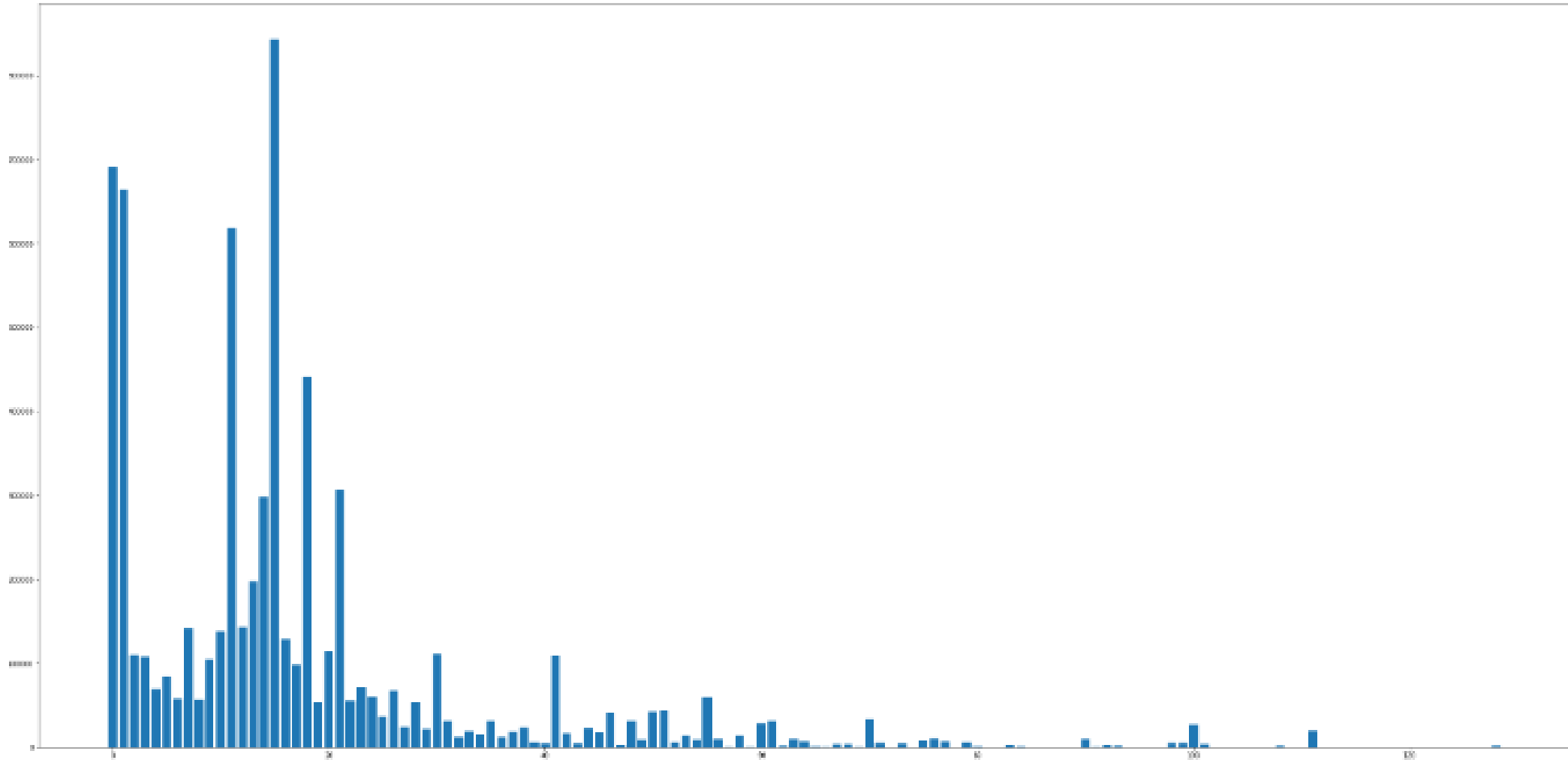


- Amount of shares by domain



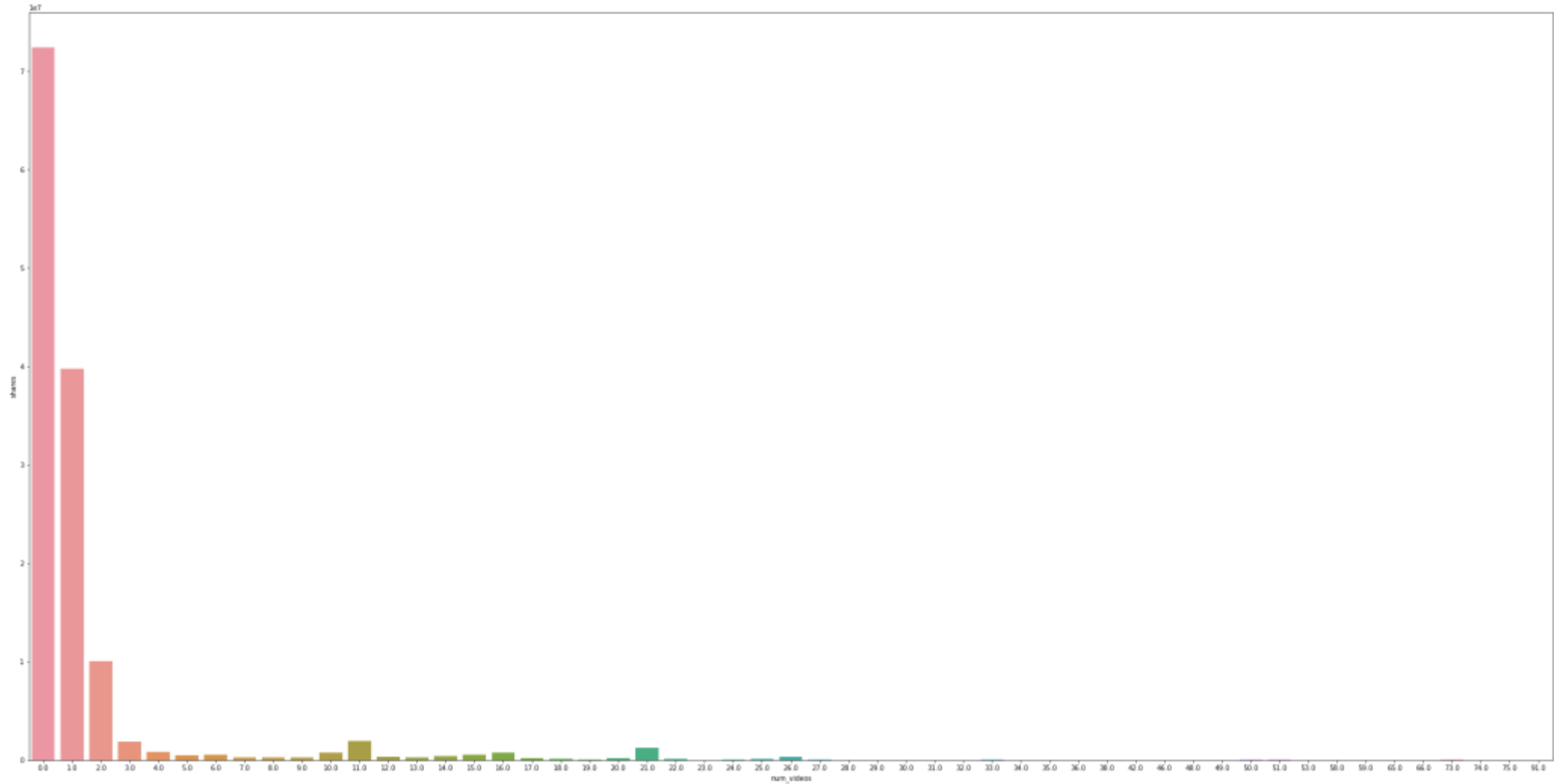
II. Visualization

- Number of shares for the number of images



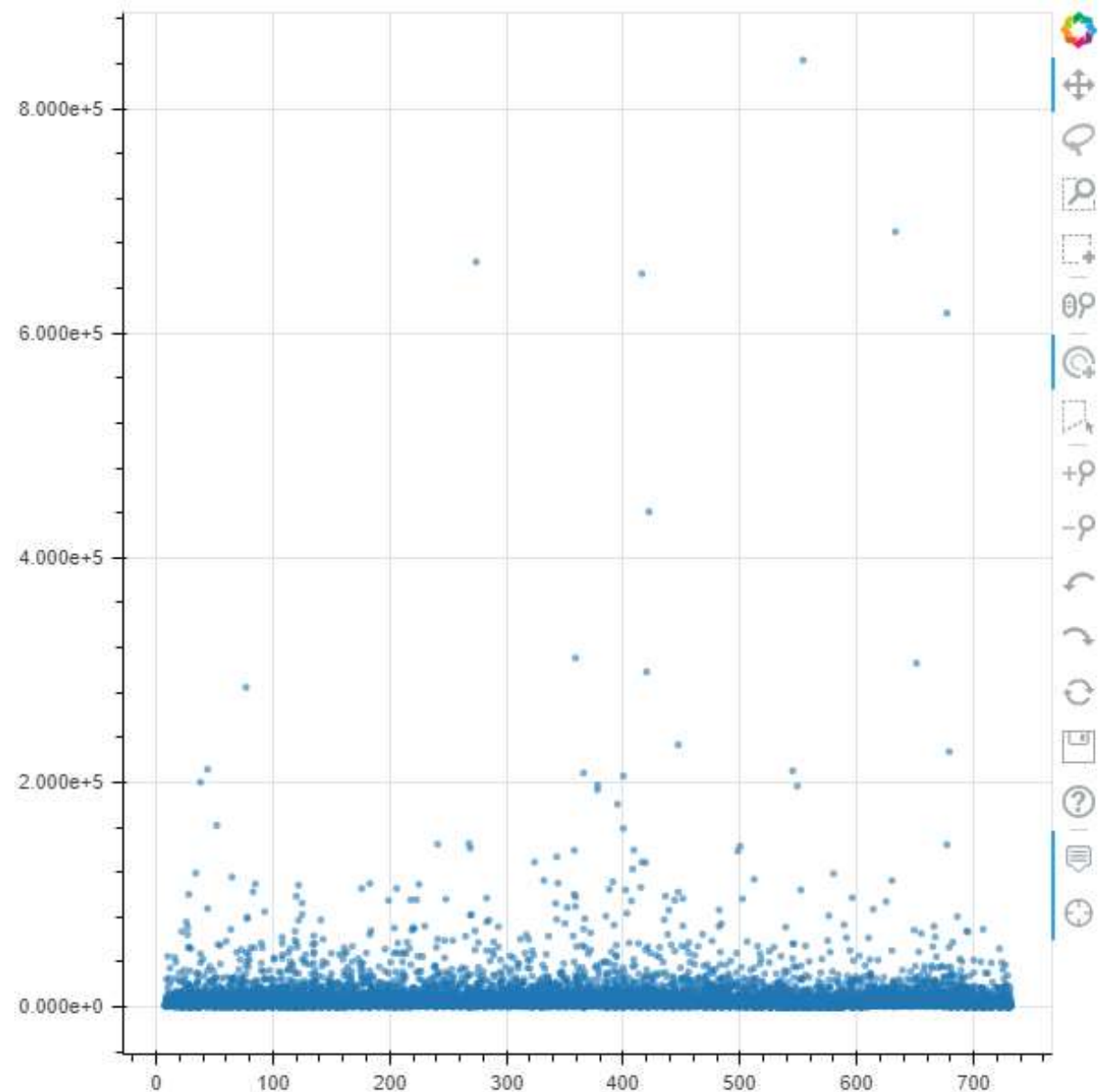
II. Visualization

- Number of shares for the number of videos



II. Visualization

- shares for the time delta using bokeh

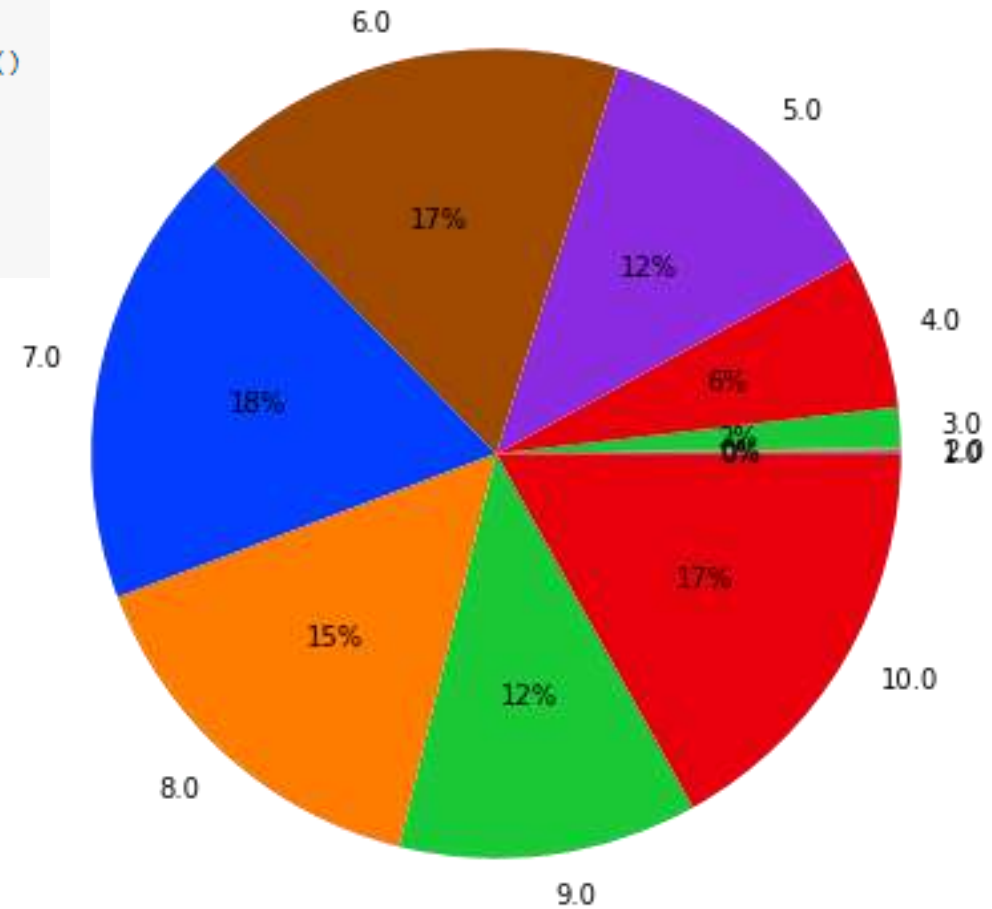


II. Visualization

- Distribution of articles according to their number of keywords

```
import matplotlib.pyplot as plt

tabgene = df[["num_keywords", "url"]].groupby("num_keywords").count()
data = tabgene["url"]
labels = tabgene.index
colors = sns.color_palette('bright')[0:6]
plt.pie(data, labels = labels, colors = colors, autopct='%.0f%%')
plt.rcParams["figure.figsize"] = (7,7)
plt.show()
```



III. Modeling

Regression

- Splitting and fitting the dataset

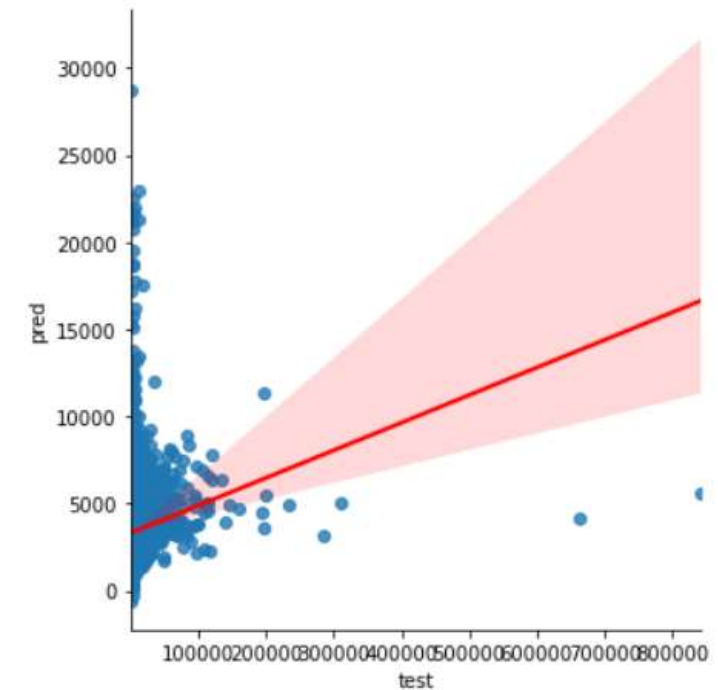
```
1 X = df.drop(['shares', 'url'], axis=1)
2 y = df['shares']
3 sc=StandardScaler()
4 X=sc.fit_transform(X)
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33)
```


III. Modeling

- Linear Regression

```
1 lm = linear_model.LinearRegression()
2 lm.fit(X_train, y_train)
3 y_pred = lm.predict(X_test)
4 print(mean_squared_error(y_test, y_pred))
5 eval_logic = pd.DataFrame()
6 eval_logic['test'] = y_test
7 eval_logic['pred'] = y_pred
8
9 sns.lmplot(x = 'test', y = 'pred', data = eval_logic, line_kws={'color':'red'})
10 plt.show()
```

168711874.01750997

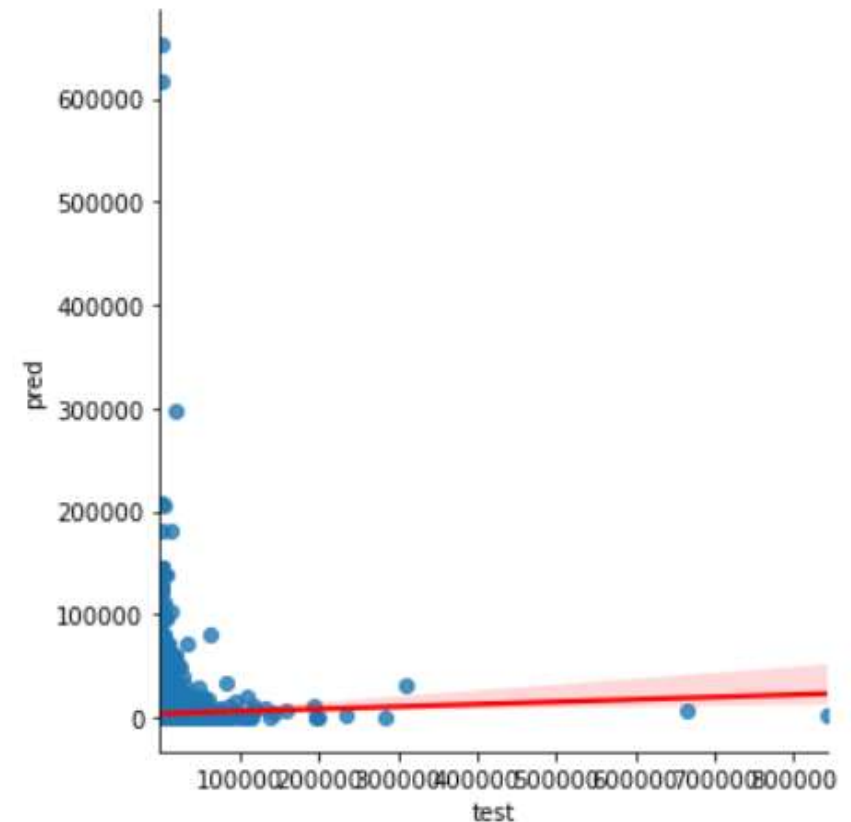


III. Modeling

- Decision Tree

```
1 dt = tree.DecisionTreeRegressor()
2 dt.fit(X_train, y_train)
3 y_pred = dt.predict(X_test)
4 print(mean_squared_error(y_test, y_pred))
5 eval_logic = pd.DataFrame()
6 eval_logic['test'] = y_test
7 eval_logic['pred'] = y_pred
8
9 sns.lmplot(x = 'test', y = 'pred', data = eval_logic, line_kws={'color':'red'})
10 plt.show()
```

316264834.3105557

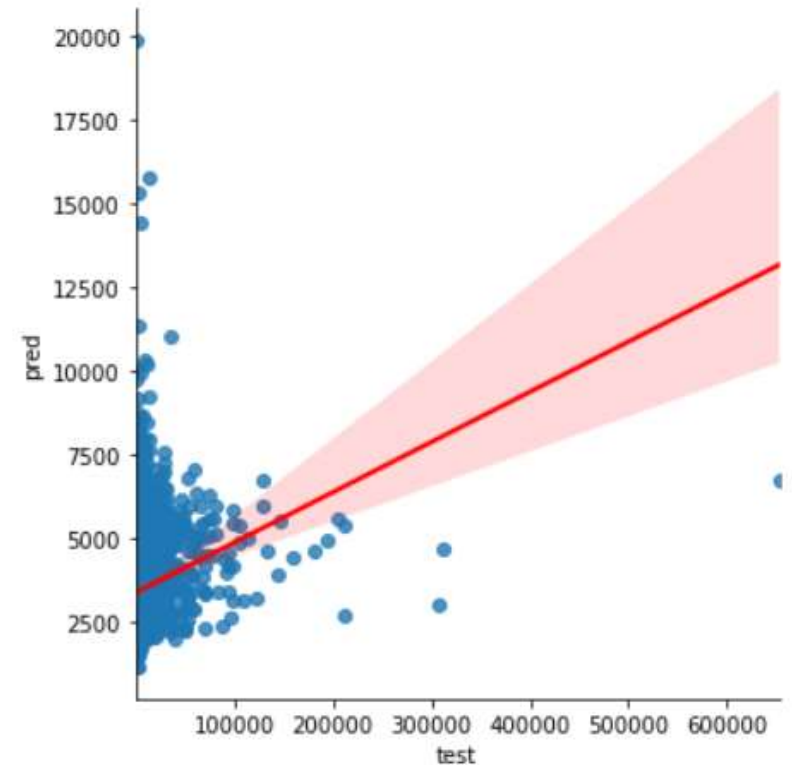


III. Modeling

- **Logic Regression**

```
1 from sklearn.linear_model import TweedieRegressor
2 lr = TweedieRegressor()
3 lr.fit(X_train, y_train)
4 y_pred = lr.predict(X_test)
5 print(mean_squared_error(y_test, y_pred))
6 eval_logic = pd.DataFrame()
7 eval_logic['test'] = y_test
8 eval_logic['pred'] = y_pred
9
10 sns.lmplot(x = 'test', y = 'pred', data = eval_logic, line_kws={'color':'red'})
11 plt.show()
```

110492485.29553321

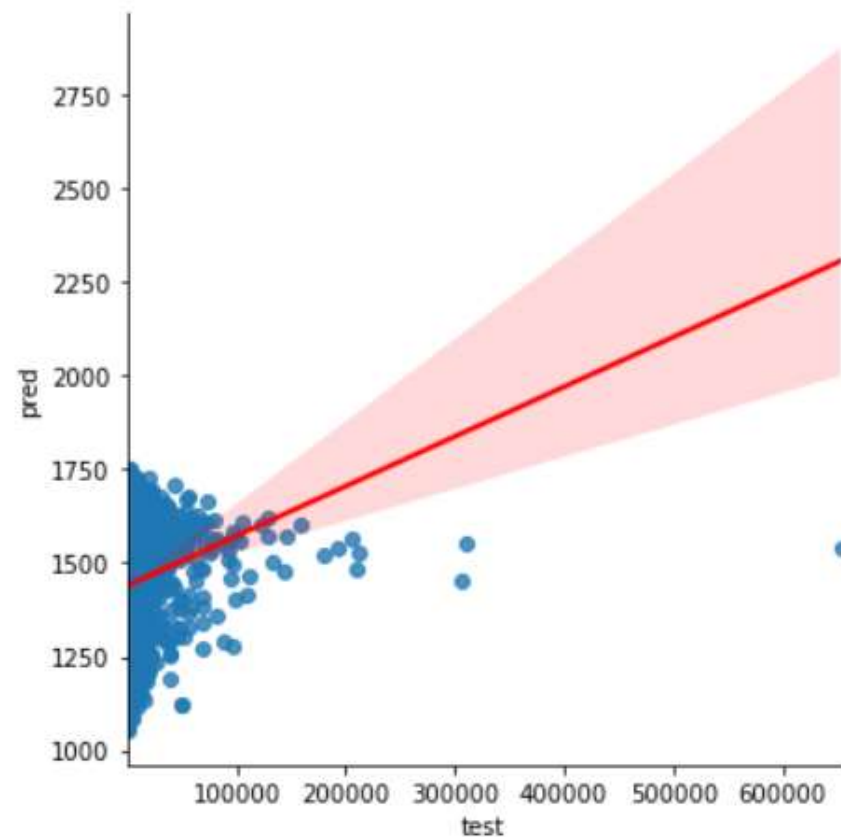


III. Modeling

- SVM

```
1 from sklearn import svm
2 regr = svm.SVR()
3 regr.fit(X_train, y_train)
4 y_pred = regr.predict(X_test)
5 print(mean_squared_error(y_test, y_pred))
6 eval_logic = pd.DataFrame()
7 eval_logic['test'] = y_test
8 eval_logic['pred'] = y_pred
9 eval_logic.sort_values(by = 'test', inplace = True)
10
11 sns.lmplot(x = 'test', y = 'pred', data = eval_logic, line_kws={'color':'red'})
12 plt.show()
```

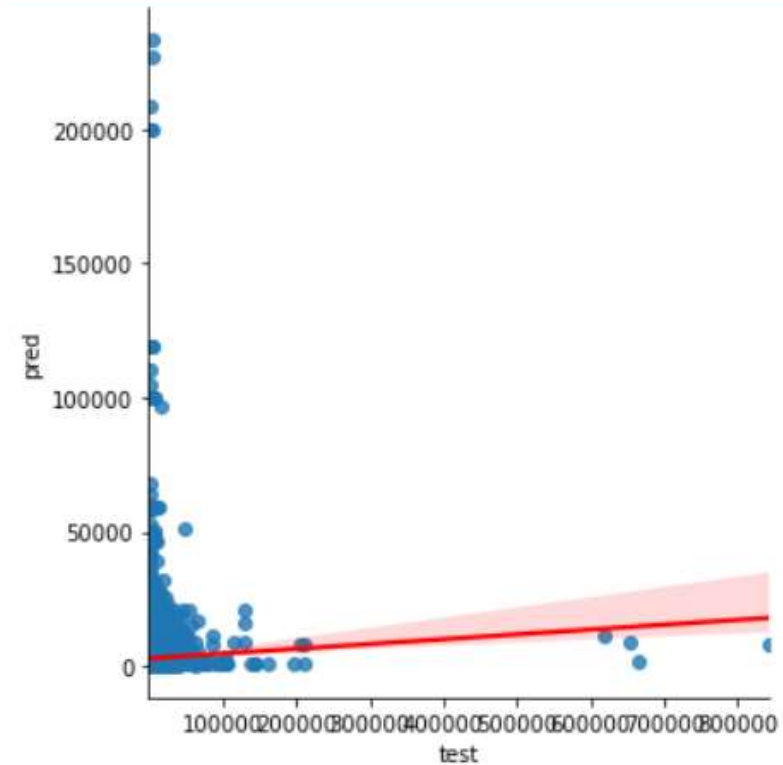
116231013.14354613



III. Modeling

- Neural Network

```
1 from sklearn.neural_network import MLPClassifier
2 clf = MLPClassifier()
3 clf.fit(X_train, y_train)
4 y_pred = clf.predict(X_test)
5 np.mean(y_test== y_pred)
6 eval_logic = pd.DataFrame()
7 eval_logic['test'] = y_test
8 eval_logic['pred'] = y_pred
9 eval_logic.sort_values(by = 'test', inplace = True)
10
11 sns.lmplot(x = 'test', y = 'pred', data = eval_logic, line_kws={'color':'red'})
12 plt.show()
```



III. Modeling

Classification

- Splitting and fitting the dataset

```
1 x = df.drop(['shares', 'url'], axis=1)
2 y = df['shares'].apply(lambda x: 0 if x < df['shares'].mean() else 1)
```

```
1 sc=StandardScaler()
2 x=sc.fit_transform(x)
3 x,y = SMOTE().fit_resample(x, y)
```

```
1 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.33, random_state=0)
```

III. Modeling

- **Logistic Regression**

```
1 LR = LogisticRegression(multi_class='auto')
2 LR.fit(x_train,y_train)
3 pred = LR.predict(x_test)
4 print("Accu", accuracy_score(y_test,pred))
5 print("MSE", metrics.mean_squared_error(y_test, pred))
```

Accu 0.6396102337637403

MSE 0.3603897662362598

III. Modeling

- **Decision Tree**

```
1 dtree = tree.DecisionTreeRegressor()  
2 dtree.fit(x_train, y_train)  
3 pred = dtree.predict(x_test)  
4 print("Accu", accuracy_score(y_test, pred))  
5 print("MSE", metrics.mean_squared_error(y_test, pred))
```

Accu 0.7727163634618154

MSE 0.22728363653818462

III. Modeling

- **Random Forest**

```
1 randf = RandomForestClassifier(n_estimators=100, max_features="sqrt", random_state=40)
2 randf.fit(x_train, y_train)
3 pred = randf.predict(x_test)
4 print("Accu", accuracy_score(y_test, pred))
5 print("MSE", metrics.mean_squared_error(y_test, pred))
```

Accu 0.8773100369605914

MSE 0.12268996303940863

III. Modeling

- **SVM**

```
1 clf = svm.SVC()  
2 clf.fit(x_train, y_train)  
3 y_pred = clf.predict(x_test)  
4 mean_squared_error(y_test, y_pred)
```

0.2874285988575817

III.Modeling

Model name	Logic Regression	Decision Tree	Random Forest	
Accuracy	0.6396102337637403	0.7727163634618154	0.8818221091537465	
MSE	0.3603897662362598	0.22728363653818462	0.11817789084625353	