

Exploratory Data Analysis Based on Instagram reach using Python

Exploratory data analysis (EDA) is a Data Science concept where we analyze a dataset to discover patterns, trends, and relationships within the data. It helps us better understand the information contained in the dataset and guides us in making informed decisions and formulating strategies to solve real business problems. If you want to understand Exploratory Data Analysis practically, this article is for you. In this article, I will take you through an implementation of Exploratory Data Analysis using Python.

To show how to perform Exploratory Data Analysis using Python, I will use a dataset based on an Instagram reach.

```
In [10]: import pandas as pd
import plotly.express as px
import plotly.graph_objects as go
import plotly.io as pio
pio.templates.default = "plotly_white"

data = pd.read_csv(r"C:\Users\shail\Downloads\archive (4)\Instagram data.csv", encoding='latin1')
```

Now let's have a look at the first five rows of the data:

```
In [11]: print(data.head())
```

	Impressions	From Home	From Hashtags	From Explore	From Other	Saves	\
0	3928	2586	1828	619	56	98	
1	5394	2727	1838	1174	78	194	
2	4021	2895	1188	0	533	41	
3	4528	2789	621	932	73	172	
4	2518	1764	255	279	37	96	

	Comments	Shares	Likes	Profile Visits	Follows	\
0	9	5	162	35	2	
1	7	14	224	48	18	
2	11	1	151	62	12	
3	10	7	213	23	8	
4	5	4	123	8	0	


```
0 Here are some of the most important data visu...
1 Here are some of the best data science projec...
2 Learn how to train a machine learning model an...
3 Here's how you can write a Python program to d...
4 Plotting annotations while visualizing your da...
```



```
0 #finance #money #business #investing #investme...
1 #healthcare #health #covid #data #data science...
2 #data #datascience #dataanalysis #dataanalytic...
3 #python #pythonprogramming #pythonprojects #py...
4 #data #visualization #datascience #data #dataana...
```

Now let's have a look at all the columns the dataset contains:

```
In [12]: print(data.columns)
```

```
Index(['Impressions', 'From Home', 'From Hashtags', 'From Explore', 'From Other', 'Saves', 'Comments', 'Shares', 'Likes', 'Profile Visits', 'Follows', 'Caption', 'Hashtags'],
      dtype='object')
```

Now let's have a look at the column info:

```
In [14]: print(data.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 119 entries, 0 to 118
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  --
0   Impressions            119 non-null    int64
1   From Home              119 non-null    int64
2   From Hashtags          119 non-null    int64
3   From Explore           119 non-null    int64
4   From Other             119 non-null    int64
5   Saves                  119 non-null    int64
6   Comments               119 non-null    int64
7   Shares                 119 non-null    int64
8   Likes                  119 non-null    int64
9   Profile Visits         119 non-null    int64
10  Follows                119 non-null    int64
11  Caption                119 non-null    object
12  Hashtags               119 non-null    object
dtypes: int64(11), object(2)
memory usage: 12.2+ KB
None
```

Next, we look at the descriptive statistics of the data:

```
In [15]: print(data.describe())
```

	Impressions	From Home	From Hashtags	From Explore	From Other	\
count	119.000000	119.000000	119.000000	119.000000	119.000000	
mean	5703.991597	2475.789916	1887.512605	1078.108840	171.092437	
std	4843.780105	1489.588348	1894.361443	2613.876132	299.431933	
min	1941.000000	1133.000000	116.000000	0.000000	9.000000	
25%	3467.000000	1945.000000	725.000000	157.500000	38.000000	
50%	4289.000000	2207.000000	1275.000000	326.000000	74.000000	
75%	6138.000000	2802.500000	2363.500000	689.500000	196.000000	
max	38919.000000	13473.000000	11817.000000	17414.000000	2547.000000	

	Saves	Comments	Shares	Likes	Profile Visits	\
count	119.000000	119.000000	119.000000	119.000000	119.000000	
mean	153.188224	6.46388909	9.361145	373.781513	50.621849	
std	156.317731	3.544576	10.089205	82.378947	87.088402	
min	22.000000	0.000000	0.000000	72.000000	4.000000	
25%	65.000000	3.000000	121.500000	15.000000	15.000000	
50%	109.000000	6.000000	6.000000	151.000000	23.000000	
75%	169.000000	8.000000	13.500000	204.000000	42.000000	
max	1895.000000	19.000000	75.000000	549.000000	611.000000	

	Follows
count	119.000000
mean	20.756303
std	40.921580
min	0.000000
25%	4.000000
50%	8.000000
75%	18.000000
max	260.000000

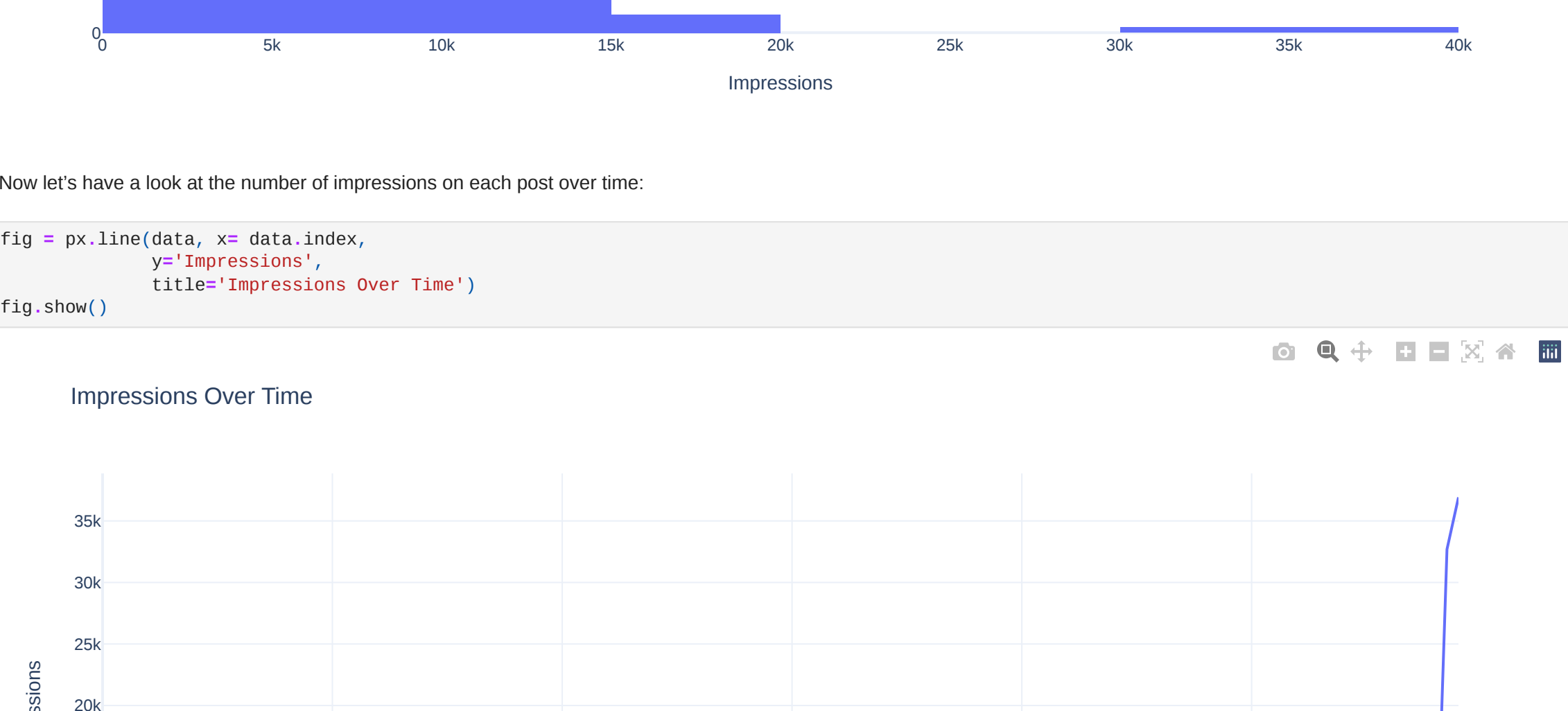
Now, before moving forward, always have a look if your data contains any missing values or not:

```
In [16]: print(data.isnull().sum())
```

```
Impressions    0
From Home      0
From Hashtags  0
From Explore   0
From Other     0
Saves          0
Comments       0
Shares         0
Likes          0
Profile Visits 0
Follows        0
Caption         0
Hashtags       0
dtype: int64
```

Luckily, this dataset doesn't have any missing values.

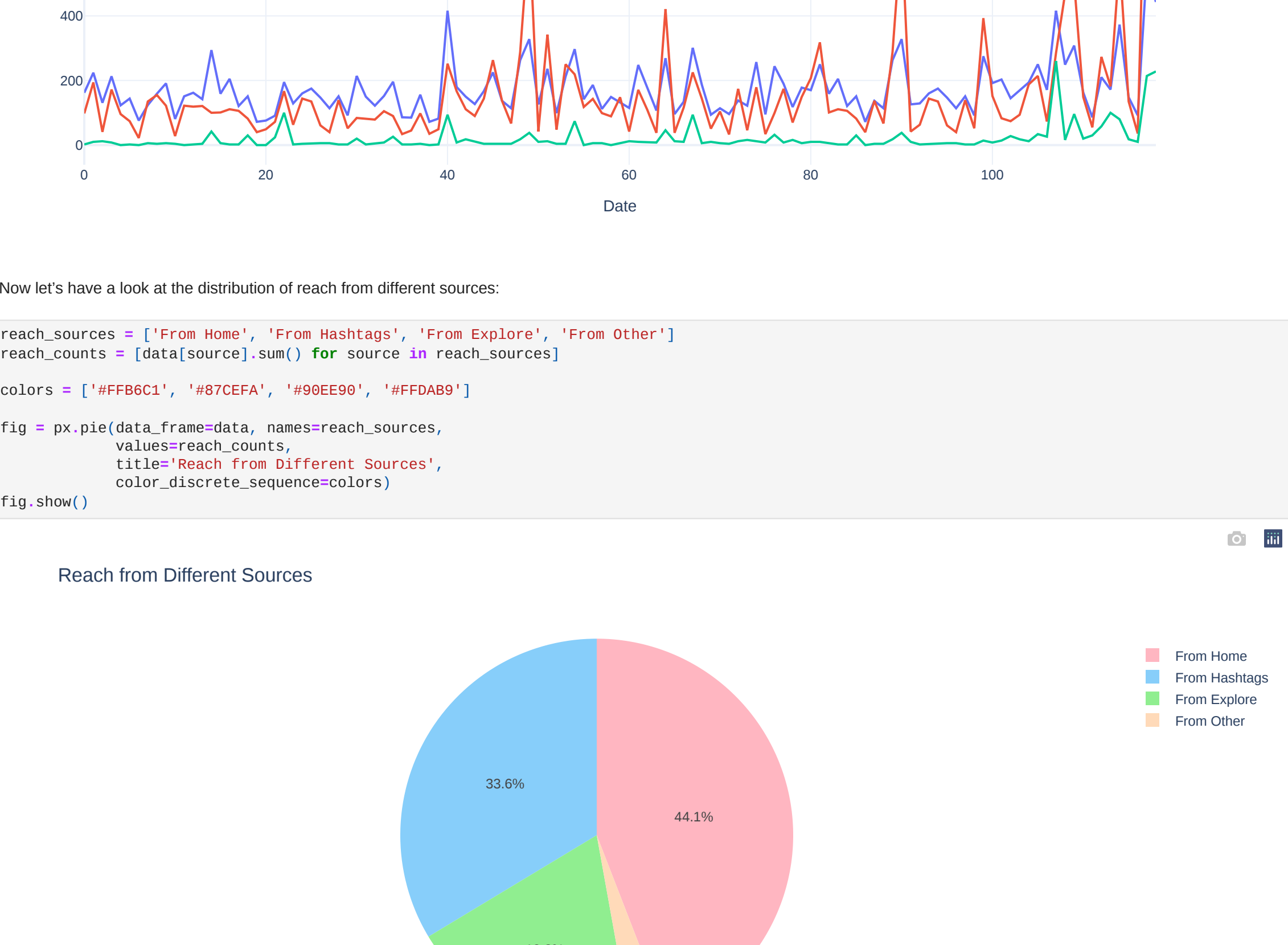
When you start exploring your data, always start by exploring the main feature of your data. For example, as we are working on a dataset based on Instagram Reach, we should start by exploring the feature that contains data about reach. In our data, the Impressions column contains the data about the reach of an Instagram post. So let's have a look at the distribution of the Impressions:



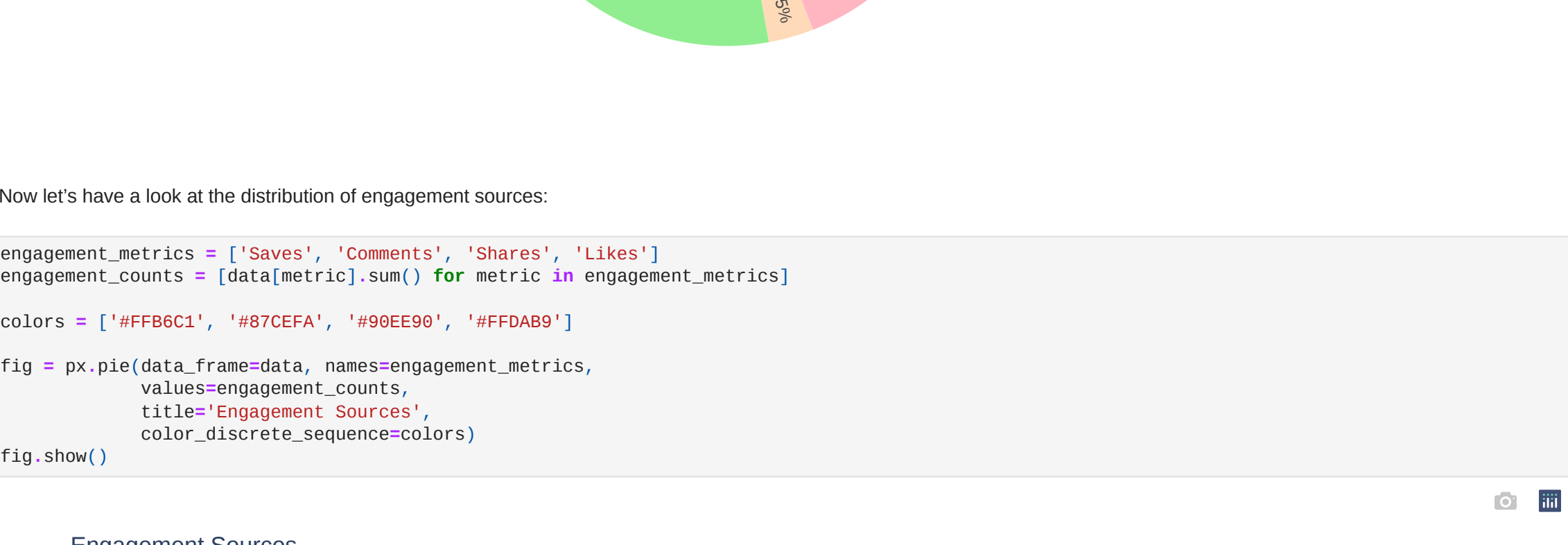
Now let's have a look at the number of impressions on each post over time:



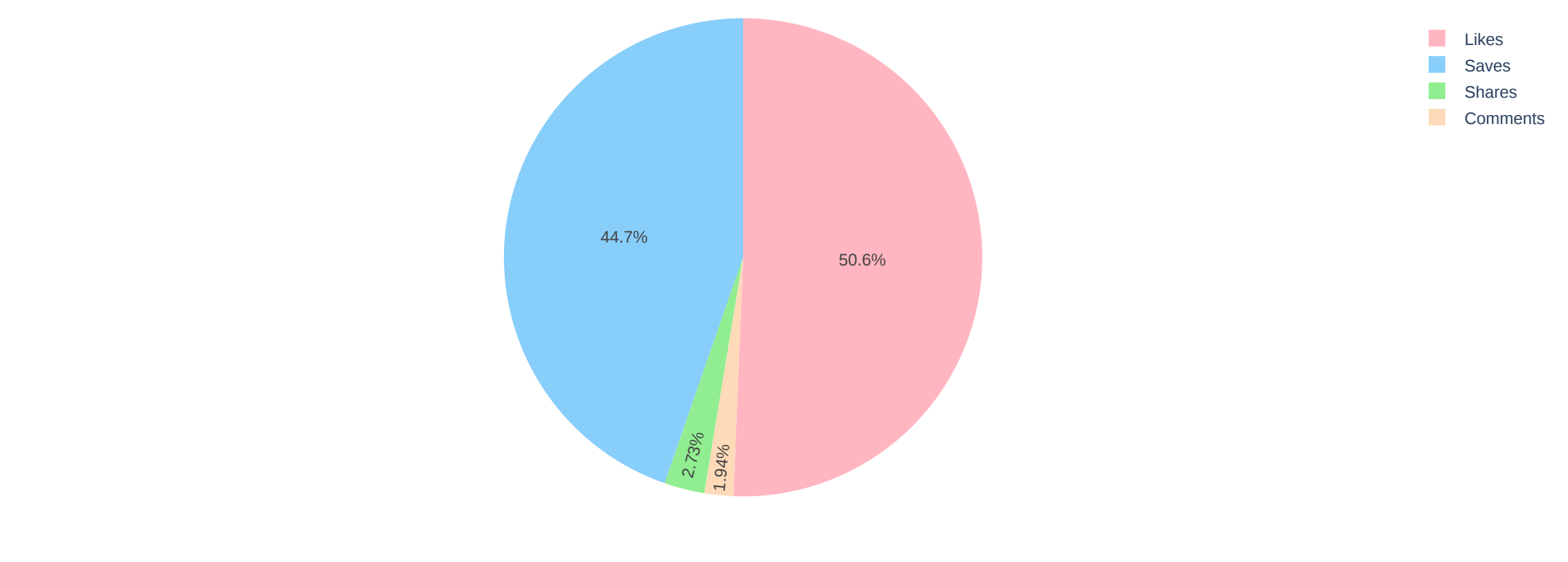
Now let's have a look at all the metrics like Likes, Saves, and Follows from each post over time:



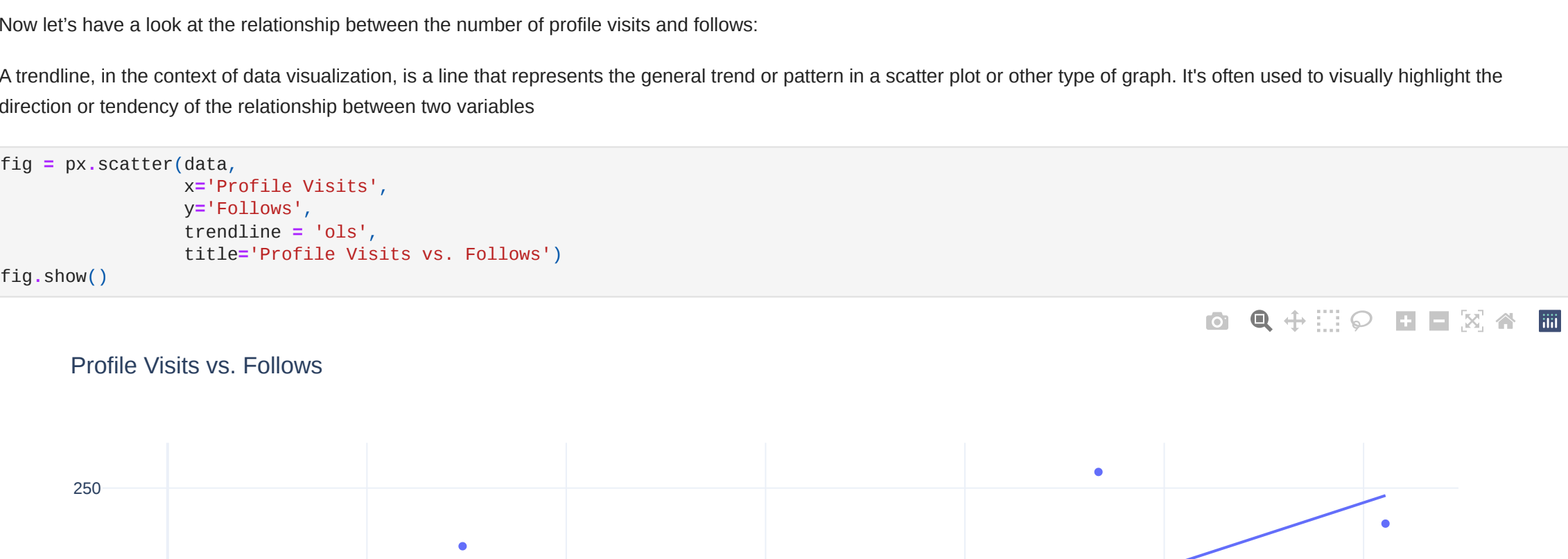
Now let's have a look at the distribution of reach from different sources:



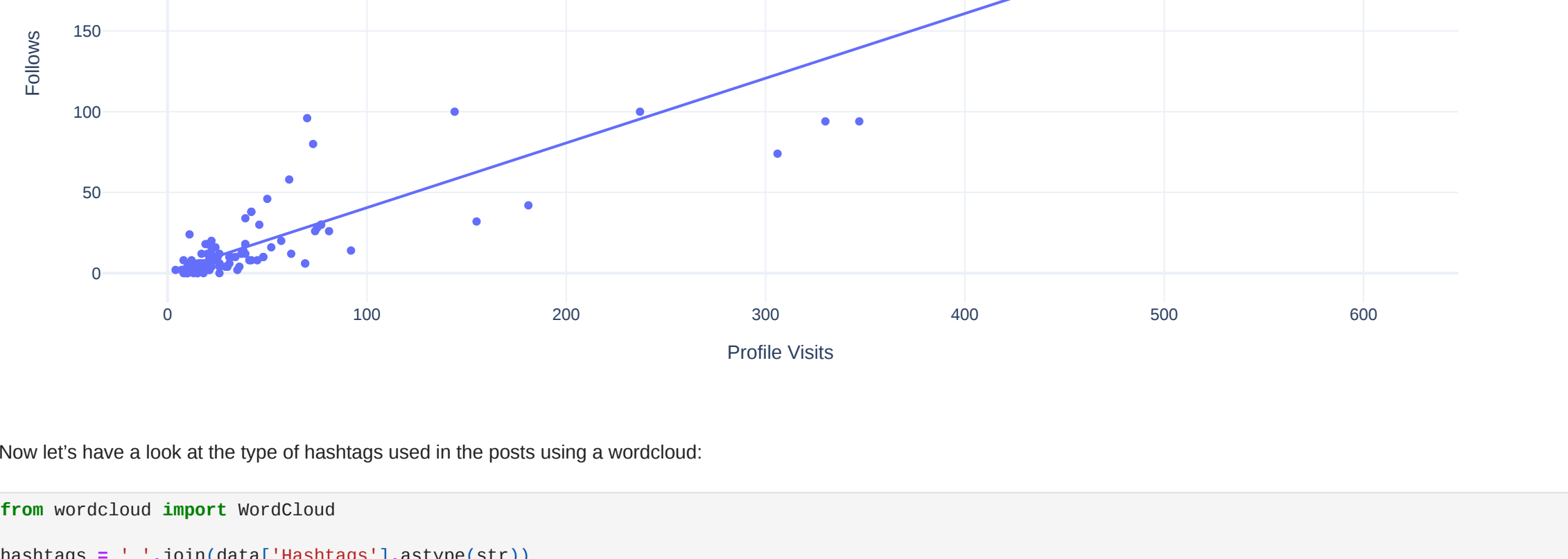
Now let's have a look at the distribution of engagement sources:



Now let's have a look at the relationship between the number of profile visits and follows:



Now let's have a look at the type of hashtags used in the posts using a wordcloud:

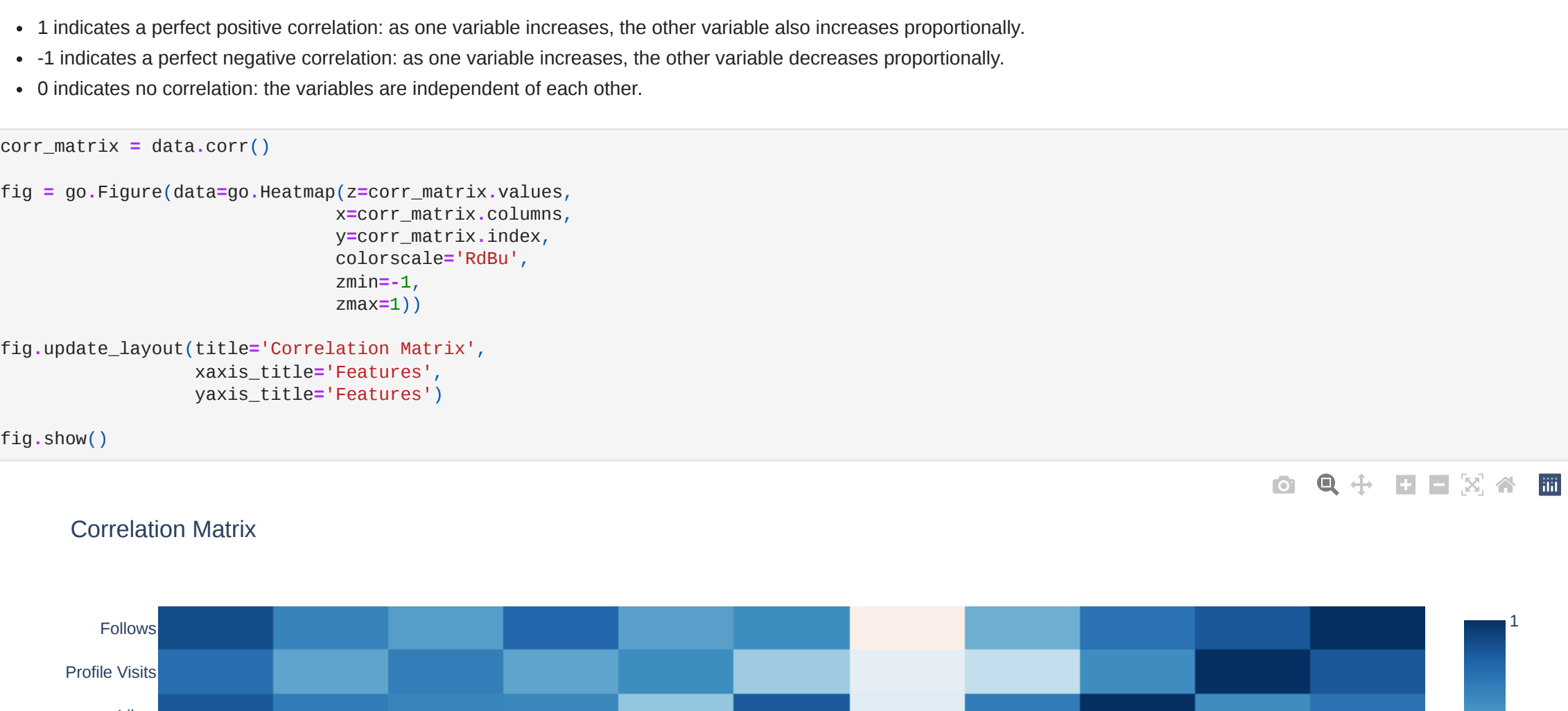


Now let's have a look at the correlation between all the features:

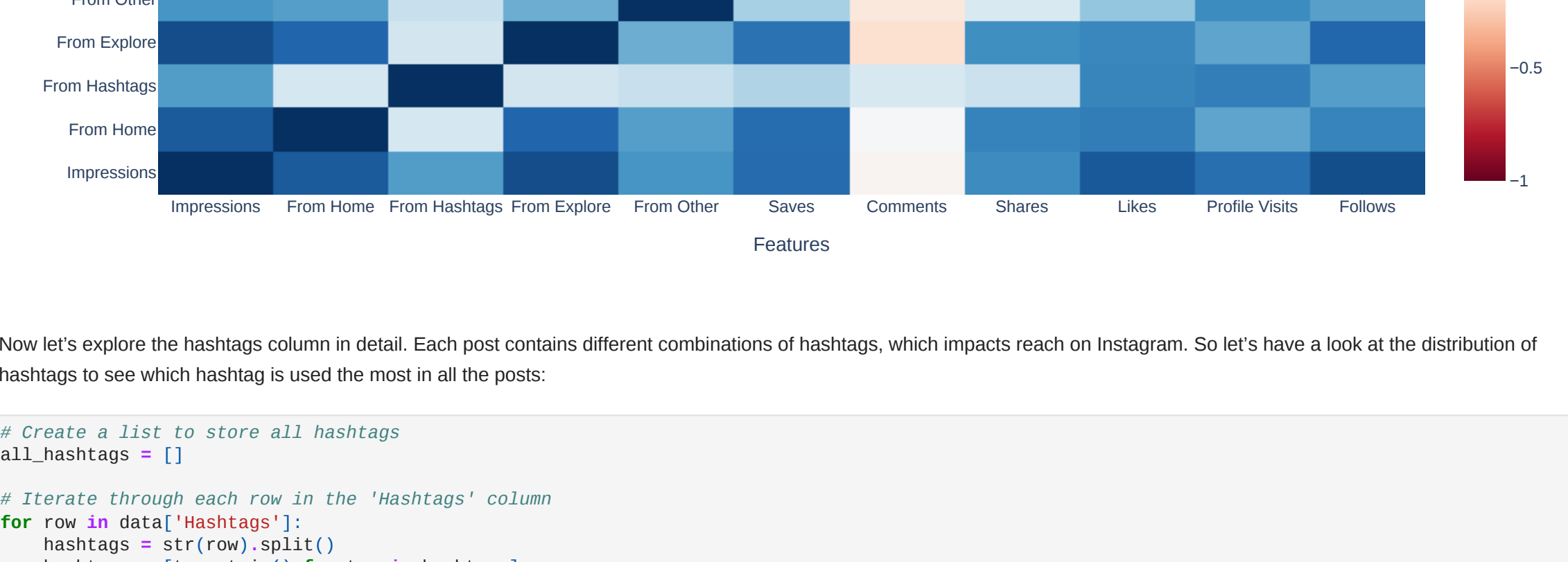
The correlation coefficient is a statistical measure that describes the strength and direction of a relationship between two variables.

The correlation coefficient can range from -1 to 1:

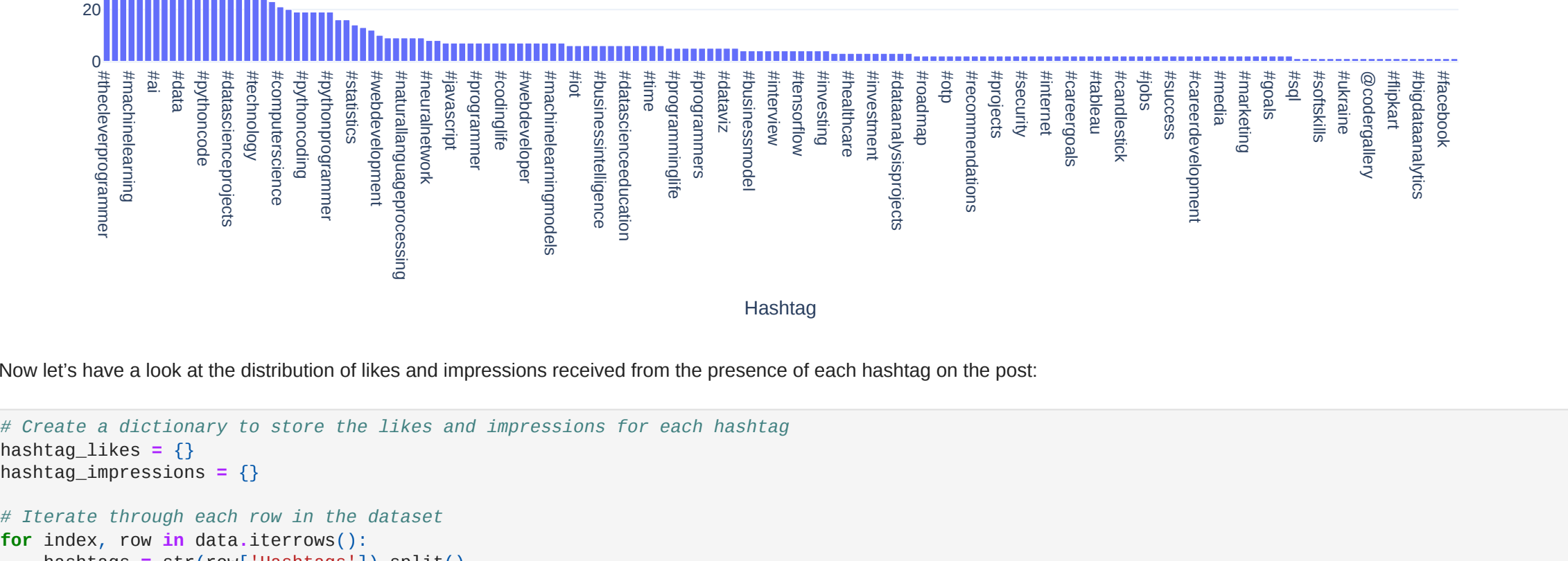
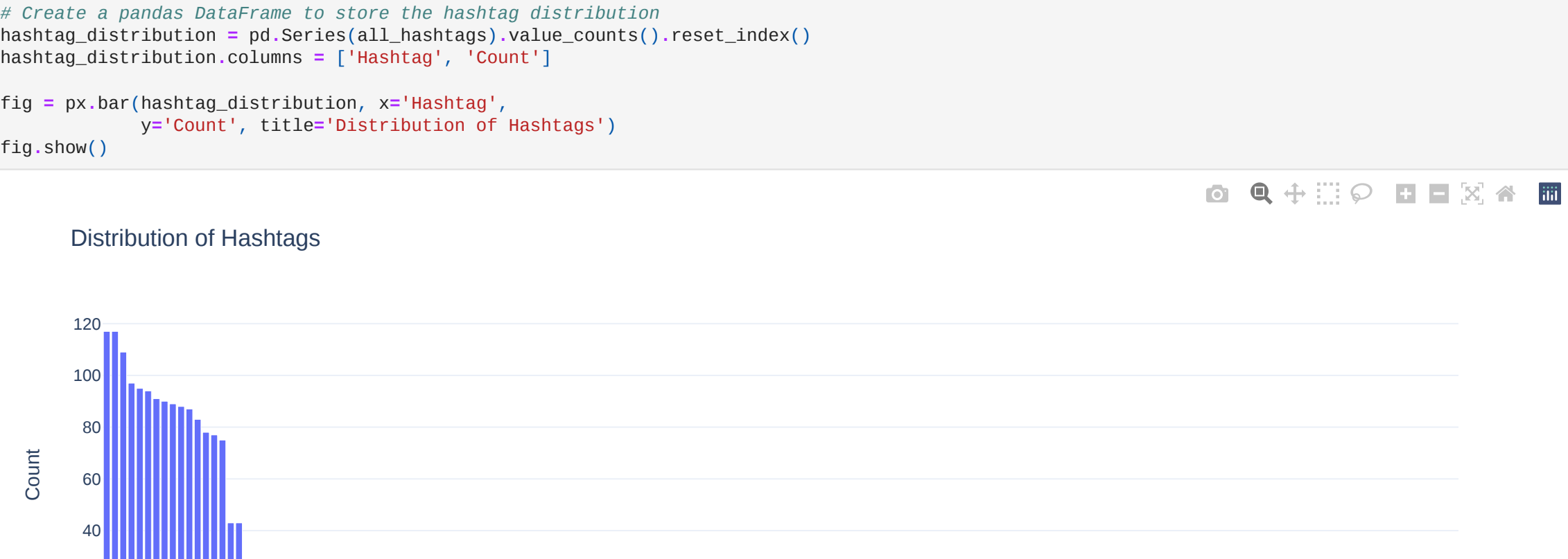
- 1 indicates a perfect positive correlation: as one variable increases, the other variable also increases proportionally.
- -1 indicates a perfect negative correlation: as one variable increases, the other variable decreases proportionally.
- 0 indicates no correlation: the variables are independent of each other.



Now let's explore the hashtags column in detail. Each post contains different combinations of hashtags, which impacts reach on Instagram. So let's have a look at the distribution of hashtags to see which hashtag is used the most in all the posts:



Now let's have a look at the distribution of likes and impressions received from the presence of each hashtag on the post:



Summary

Exploratory data analysis (EDA) is a Data Science concept where we analyze a dataset to discover patterns, trends, and relationships within the data. It helps us better understand the information contained in the dataset and guides us in making informed decisions and formulating strategies to solve real business problems. I hope you liked this article on Exploratory Data Analysis using Python. Feel free to ask valuable questions in the comments section below.

In []: