

```
In [1]: ### >> EXPLORATORY DATA ANALYSIS OF STACK OVERFLOW DATA << ###
```

```
In [2]: # Importing Modules
import pymongo
import wordcloud
import numpy as np
import pandas as pd
import seaborn as sb
import collections as C
import matplotlib.pyplot as plt
import matplotlib inline

In [3]: # Connecting to Local MongoDB Client
client = pymongo.MongoClient('localhost:27017')
```

```
db = client['StackOverflow']
```

```
In [4]: ### EDA on Tags Collection ###
```

```
In [5]: # Loading the Tags Collection
        cursor = db['Tags'].find({})[1]
        df = pd.DataFrame(list(cursor))
```

```
del df['id']

# Selecting relevant columns
df = df[['TagName', 'Count']]
df['Count'] = df['Count'].astype(int)
df.dropna(inplace=True)
df.reset_index(drop=True, inplace=True)

df.sort_values(inplace=True, axis=0, by=['Count'], ascending=False)
df.head()
```

Out[5]:

	TagName	Count
2	javascript	1955557
11	java	1641102
6	c#	1385220
10	python	1359126

```
In [6]: # Creating Word Cloud

cloud = wordcloud.L
Data = { tag : cou
cloud.generate_from
plt.figure(figsize=
plt.axis('off')
plt.imshow(cloud)
plt.savefig('Image
```

[illegible]

```
In [7]: ### EDA on Posts Collection ###
```

```
In [8]: # Traversing the Posts Collection and analysing the distribution of Tags
```

```
cursor = db['Posts'].find({})[1:]
Post_Tags = {}
Total_Posts = 0
Tagged_Posts = 0

for x in cursor:
    if 'Tags' in x:
```

```
for t in Tag_List:
    if t in Post_Tags:
        Post_Tags[t] +=
```

```

        else:
            Post_Tags[t] = 1
            Tagged_Posts += 1

```

```
print("Total Posts :",Total_Posts)
print("Tagged Posts :",tagged_Posts)
print("Unique Tags :",len(Post_Tags))

Total_Posts : 3380601
Tagged Posts : 1389860
Unique Tags : 22512

In [9]: # Plotting Top 10 Tags

Tag_List = [(Post_Tags[i],i) for i in Post_Tags]
Tag_List.sort(reverse=True)
Top_Tags = Tag_List[:10]
```

```
sb.set_theme(context='notebook')
plt.figure(figsize=(20, 10))
```

```
y = [i[0] for i in Top_Tags]
plt.bar(x,y,color='springgreen')
for i in range(1,len(y)):
    plt.text(i-0.25,y[i],str(y[i]),fontsize=15)
plt.xlabel("Tags",fontsize = 30)
plt.ylabel("No. of Questions",fontsize = 30)
plt.title("Top 10 Tags on Stack Overflow Posts")
plt.savefig('Images/Top_10_Tags.png')
```

Tag	No. of Questions
python	1358860
java	~100000
javascript	~50000
php	~50000
css	~50000
html	~50000
sql	~50000
ruby	~50000
perl	~50000
go	~50000

Tags	No. of Questions
python	0.75569
python-3.x	0.125609
pandas	0.119220
django	0.111953
numpy	0.062820
python-2.7	0.060590
list	0.042883
matplotlib	0.038246
dataframe	0.033927
dictionary	0.029821

```
In [10]: # Counting unique instances of attributes in Posts to identify subsampling criteria

Attributes = set()

cursor = db['Posts'].find({})[1:]
for x in cursor:
    Attributes = Attributes.union(set(x.keys()))

Attribute_Sets = (i:set() for i in Attributes)

cursor = db['Posts'].find({})[1:]
for x in cursor:
    for i in Attributes:
        if i in x:
            Attribute_Sets[i].add(x[i])

X = []
for i in Attributes:
    X.append((len(Attribute_Sets[i]), i))
X.sort()

for i in X:
    print(i[1], i[0])

PostTypeId 2
AnswerCount 53
CommentCount 53
FavoriteCount 396
Score 1230
CommunityOwnedDate 5534
LastEditorDisplayName 5733
OwnerDisplayName 15482
ViewCount 40082
ClosedDate 84382
LastEditorUserId 262918
Tags 467829
OwnerUserId 670395
AcceptedAnswerId 730440
ParentId 1158038
Title 1358340
LastEditDate 1410408
LastActivityDate 2501726
CreationDate 3377143
Body 3380053
_id 3380601

In [11]: # Analysing Activity on Stack Overflow over the years

cursor = db['Posts'].find({})[1:]
```

```
for x in cursor:
```

```

if Year in Activity:
    Activity[Year] += 1
else:
    Activity[Year] = 1


df = pd.DataFrame()
df['Year'] = [i for i in Activity]
df['Posts'] = [Activity[i] for i in Activity]

sb.set_theme(context = "poster", style="darkgrid")
plt.figure(figsize=(20,10))
g = sb.barplot(data=df,x="Year",y="Posts")
Y = list(df['Posts'])
for i in range(len(Y)):
    g.text(-0.2,Y[i],str(Y[i]),fontsize=13)
plt.title("No. of Posts Vs. Year")
plt.savefig('Images/Year_Activity.png')

```

Year	No. of Posts
2010	~10,000
2011	~15,000
2012	~20,000
2013	~25,000
2014	~30,000
2015	~35,000
2016	~40,000
2017	~45,000
2018	~50,000
2019	570996

Category	Posts
Blue	386437
Purple	340100
Pink	292195
Teal	259253
Green	184946
Yellow	120351



Year	Posts
2008	9077
2009	51954
2010	90279
2011	
2012	
2013	
2014	
2015	
2016	
2017	
2018	
2019	
2020	94933

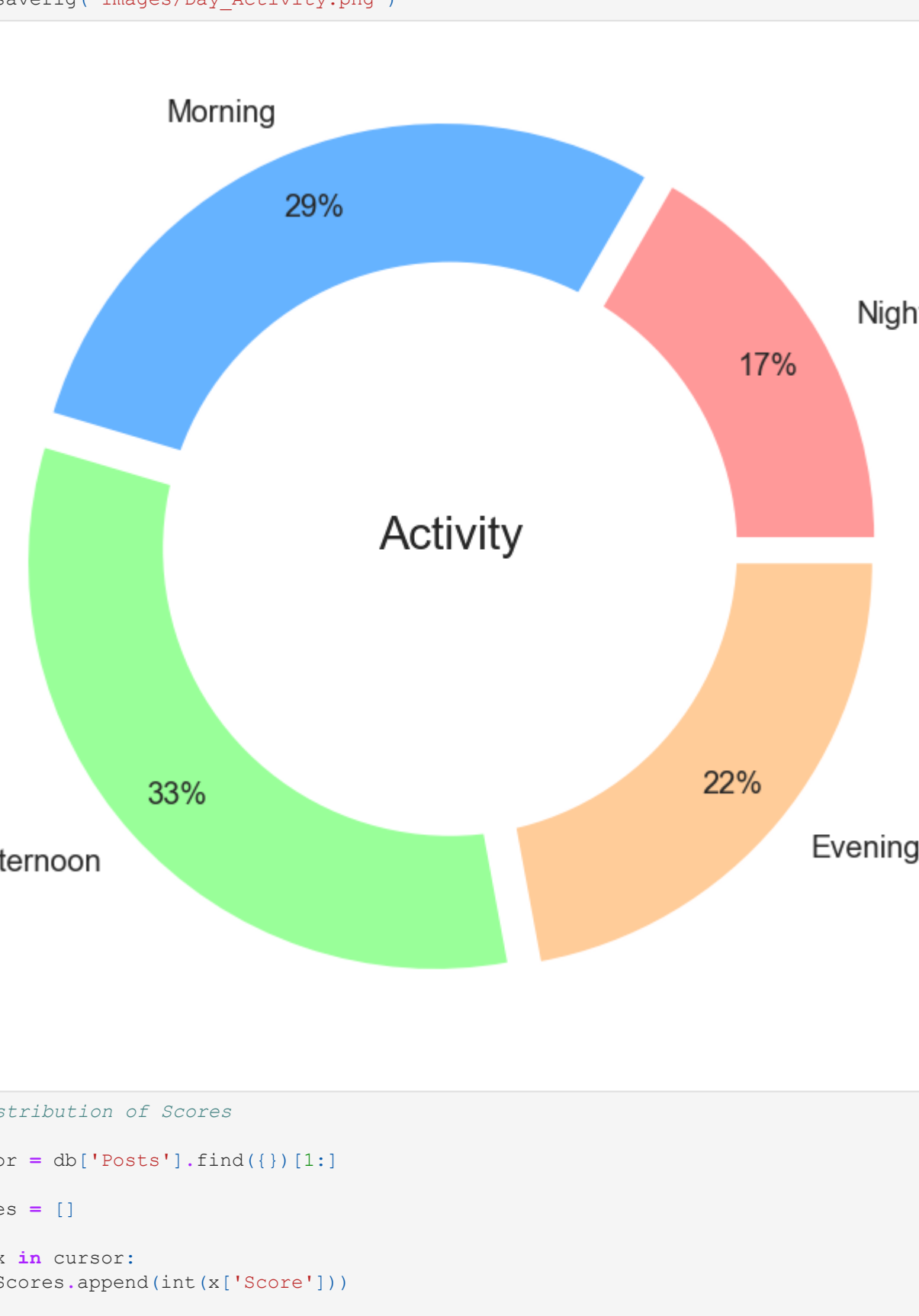
```
In [20]: # Tracking Activity on StackOverflow through the day
         cursor = db['Posts'].find({})[1]
         print('Posts: %s' % db['Posts'].find({}).count())
```

```
for x in cursor:
    A = int(x['CreationDate'][:11:13])
    if A>=0 and A<6: Periods["Night"]+=1
    elif A>=6 and A<12: Periods["Morning"]+=1
    elif A>=12 and A<18: Periods["Afternoon"]+=1
    elif A>=18 and A<24: Periods["Evening"]+=1

df = pd.DataFrame()
df['Period'] = [i for i in Periods]
df['Posts'] = [Periods[i] for i in Periods]

sizes = list(Periods.values())
labels = list(Periods.keys())
explode = (0.05,0.05,0.05,0.05)
colors = ['#f99999','#66b3ff','#99ff99','#ffcc99']

sb.set_theme(context = "poster", style="darkgrid")
plt.figure(figsize=(10,10))
plt.pie(sizes, colors = colors, labels=labels, autopct='%1.0f%%', pctdistance=0.85, explode = explode, textpct=0.5)
fig = plt.gcf()
fig.gca().add_artist([plt.Circle((0,0),0.70,fc='white')])
fig.gca().annotate("Activity", xy=(0, 0), fontsize=30,ha="center")
plt.tight_layout()
```



```
df = pd.DataFrame()
df['Scores'] = Scores
Temp = df.describe().apply(lambda s: s.apply('{0:2f}'.format))
X = [Temp.index, Temp['Scores']]
X = list(map(list, zip(*X)))

plt.figure(figsize=(5,2))
plt.axis('off')
T = plt.table(X, cellLoc='center', colLabels=['Statistic', 'Value'], colColours=['deepskyblue', 'deepskyblue'])
T.set_fontsize(20)
T.scale(1, 3)
```

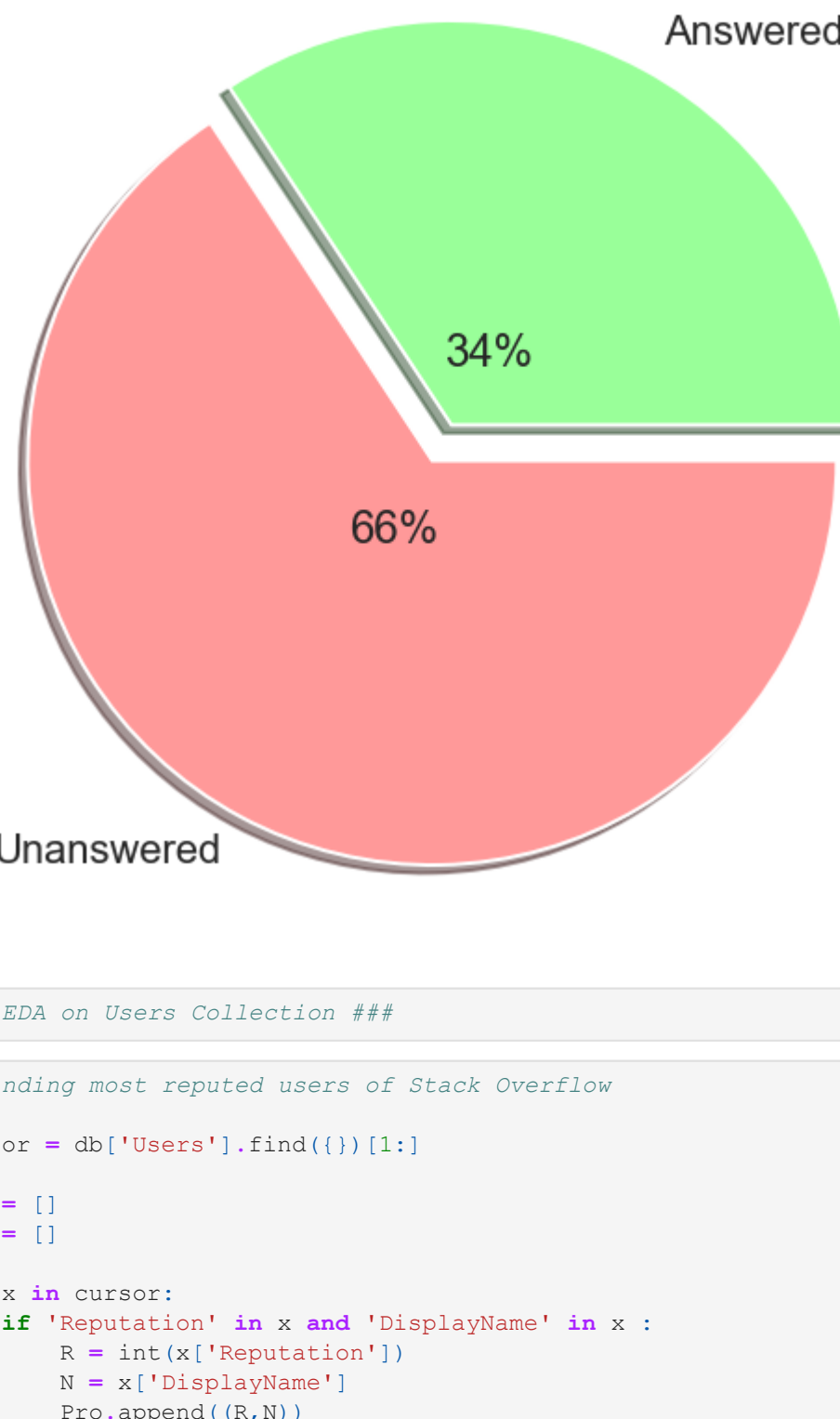
Statistic	Value
count	3380601.00
mean	3.04
std	30.27
min	-64.00
25%	0.00
50%	1.00
75%	2.00
max	14282.00

```
# Determining distribution of Answers
cursor = db['Posts'].find({})[:1]
```

Answered = 0
Unanswered = 0

```
for x in cursor:
    if 'AnswerCount' in x and int(x['AnswerCount'])>0:
        Answered+=1
    else:
        Unanswered+=1

sb.set_theme(context = "poster", style="darkgrid")
plt.figure(figsize=(20,10))
x = [Answered,Unanswered]
plt.pie(x, labels="Answered", "Unanswered", explode=(0.05,0.05), colors=['#99ff99','#ff9999'], autopct='%1.0f')
plt.savefig('Images/Answers.png')
```



```
Rep.append(R)
```



```
Pro.sort(reverse=True)
```

```
df = pd.DataFrame()
df["Name"] = [i[1] for i in Data]
```

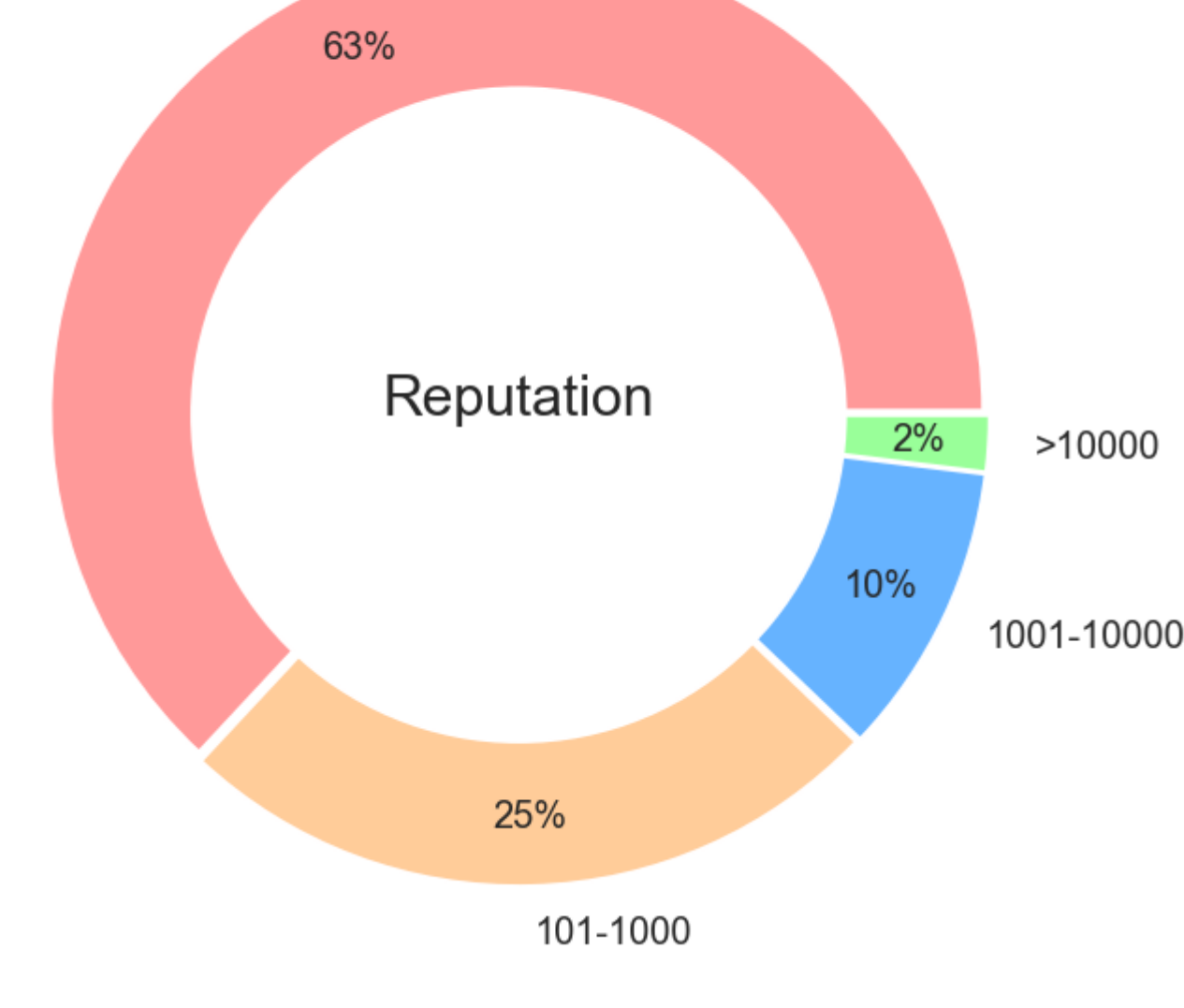
```
df["Reputation"] = [i[0] for i in Pro]

sb.set_theme(context = "poster", style="darkgrid")
plt.figure(figsize=(30,10))
g = sb.barplot(data=df,x="User",y="Reputation")
Y = list(df["Reputation"])
for i in range(len(Y)):
    g.text(i-0.2,Y[i],str(Y[i]),fontsize=15)
plt.title("Top 10 Users")
plt.savefig('Images/Top_Users.png')
```



```
explode = (0.01,0.01,0.01,0.01)
colors = ['#ff9999','#ffcc99','#66b3ff']
```

```
plt.figure(figsize=(10,10))
plt.pie(sizes, colors = colors, labels=labels, autopct='%1.0f%%', pctdistance=0.85, explode = explode, textp=
fig = plt.gcf()
fig.gca().add_artist(plt.Circle((0,0),0.70,fc='white'))
fig.gca().annotate("Reputation", xy=(0, 0), fontsize=30,ha="center")
plt.tight_layout()
plt.savefig('Images/Reputation.png')
```



```
In [ ]: ### END OF NOTEBOOK ###
```