# Click Fraud Detection Project Documentation

## Submitted By-

## Team 2 BFS Track

**Rahul Tomar – Team Lead**

**Dipanshu Upreti – Python Track**

**Vaishnavi Subramanya Desai – Python Track**

**Girish Kumar Reddy Tokala – Python Track**
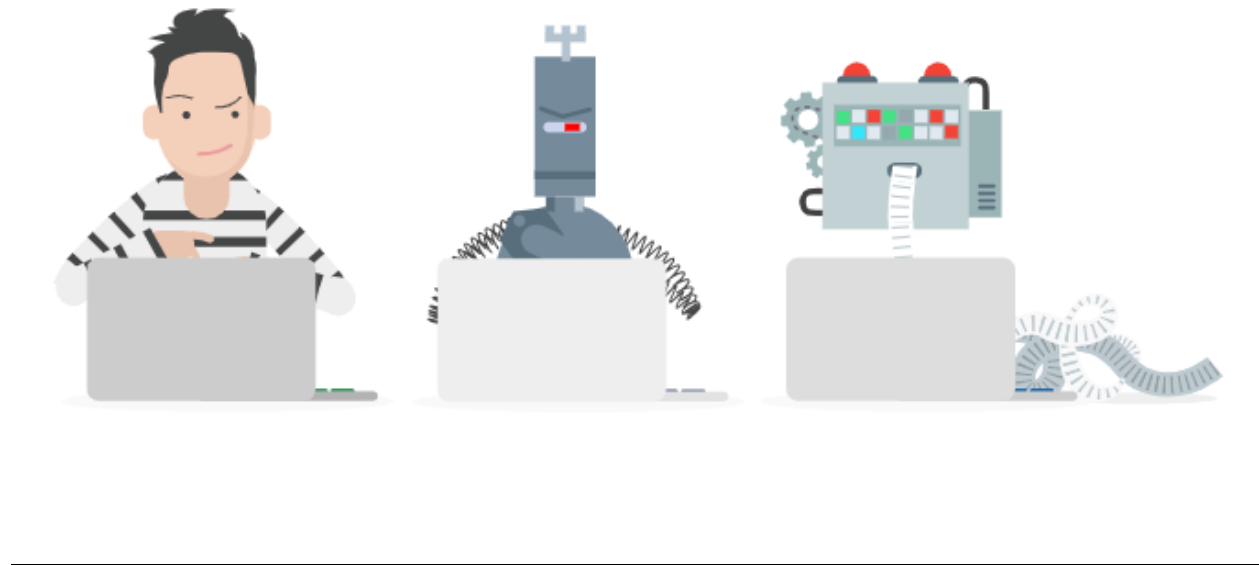
**Ragini Gaurav - UI Track**

# Table of Contents

# CLICK FRAUD DETECTION PROJECT

## Introduction

**Click fraud** is a type of fraud that occurs on the Internet in pay-per-click (PPC) online advertising. In this type of advertising, the owners of websites that post the ads are paid an amount of money determined by how many visitors to the sites click on the ads. Fraud occurs when a person, automated script, or computer program imitates a legitimate user of a web browser, clicking on such an ad without having actual interest in the target of the ad's link.
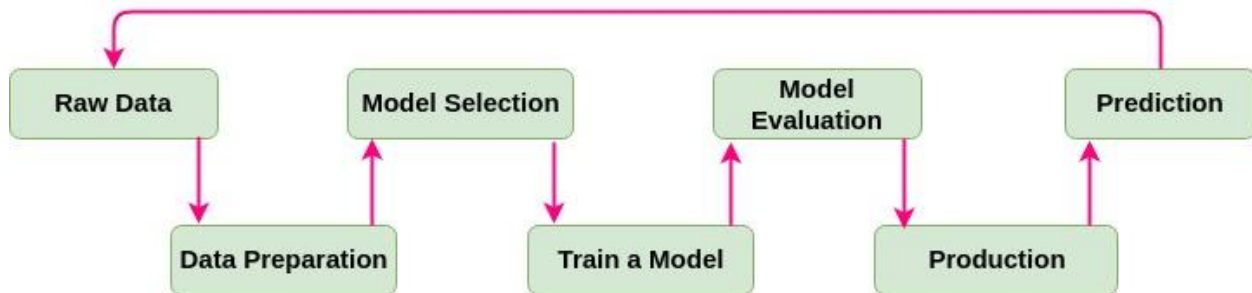
Fraud risk is everywhere, but for companies that advertise online, click fraud can happen at an overwhelming volume, resulting in misleading click data and wasted money. Ad channels can drive up costs by simply clicking on the ad at a large scale.



Although Google Adsense has checked in place it is still advisable to monitor your traffic details and analyzing them from time to time to monitor the data.

That's where our product come in the picture. It serves as a way to predict the spamming IP and Devices and find the rate at which Channels are generating clicks to find and isolate wasteful channels and increase spending in the better ones.

# Model Development



# Dataset

Dataset used by us was picked from Kaggle TalkingData AdTracking Fraud Detection Challenge(https://www.kaggle.com/c/talkingdata-adtracking-fraud-detection/data )

Each row of the training data contains a click record, with the following features.

- IP: IP address of click.
- app: app id for marketing.
- device: device **type** id of user mobile phone (e.g., iPhone 6 plus, iPhone 7, Huawei mate 7, etc.)
- os: os version id of user mobile phone
- channel: channel id of mobile ad publisher
- click_time: timestamp of click (UTC)
- attributed_time: if user download the app for after clicking an ad, this is the time of the app download
- is_attributed: the target that is to be predicted, indicating the app was downloaded

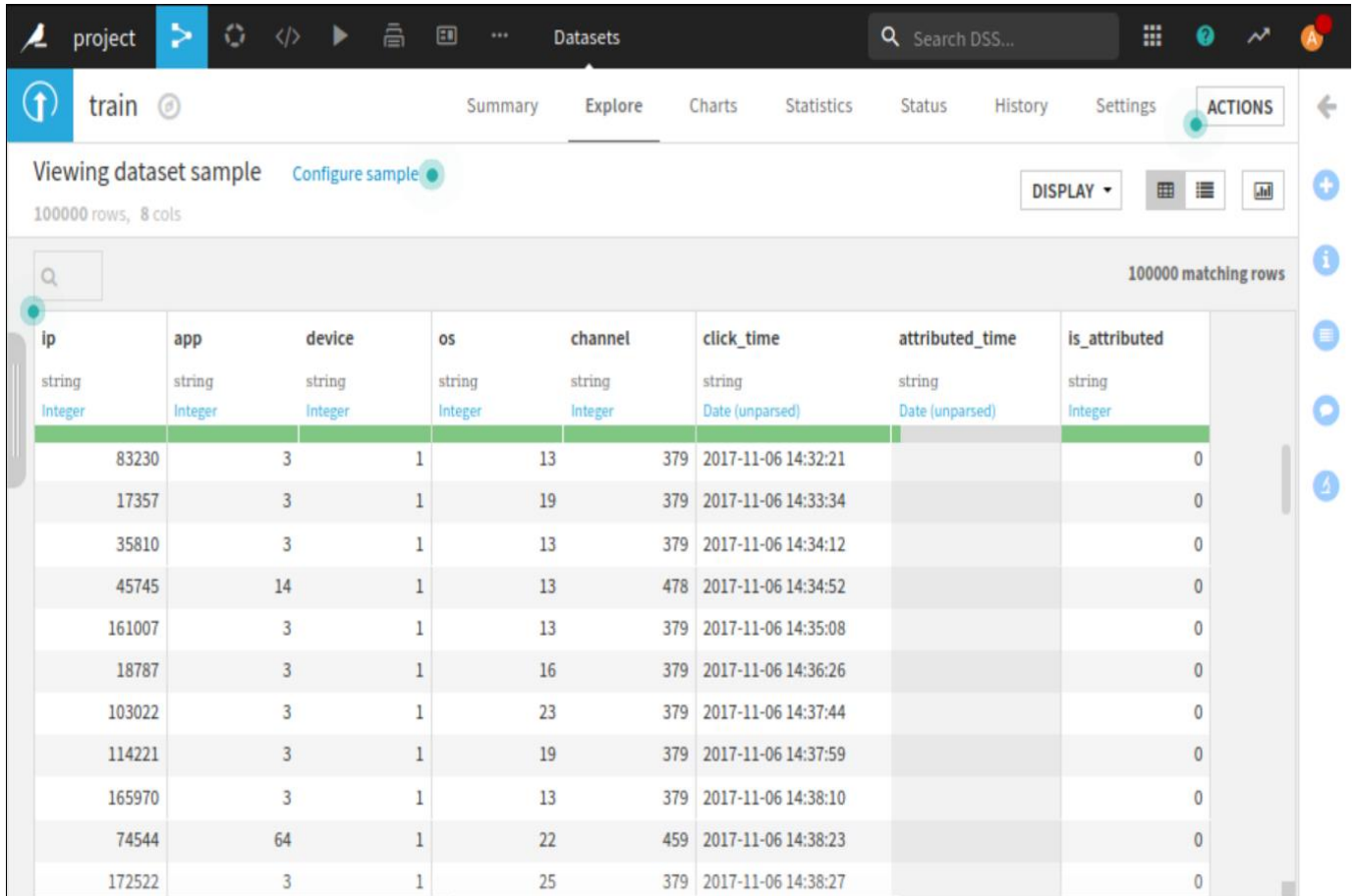Note that IP, app, device, os, and channel are encoded.

The test data is similar, with the following differences:

- click_id: a reference for making predictions
- is_attributed: not included

Our task was to build a machine learning model using **Dataiku DSS**(Data Science Software) which is a collaborative data science software platform for teams of data scientists, data analysts, and engineers to explore, prototype, build and deliver their data products more efficiently.

# Loading dataset in Dataiku

This is how our dataset looked like



The actual train dataset was very large around **7GB** in size containing 180 million entries. With our systems, it was not feasible to work in such big dataset.
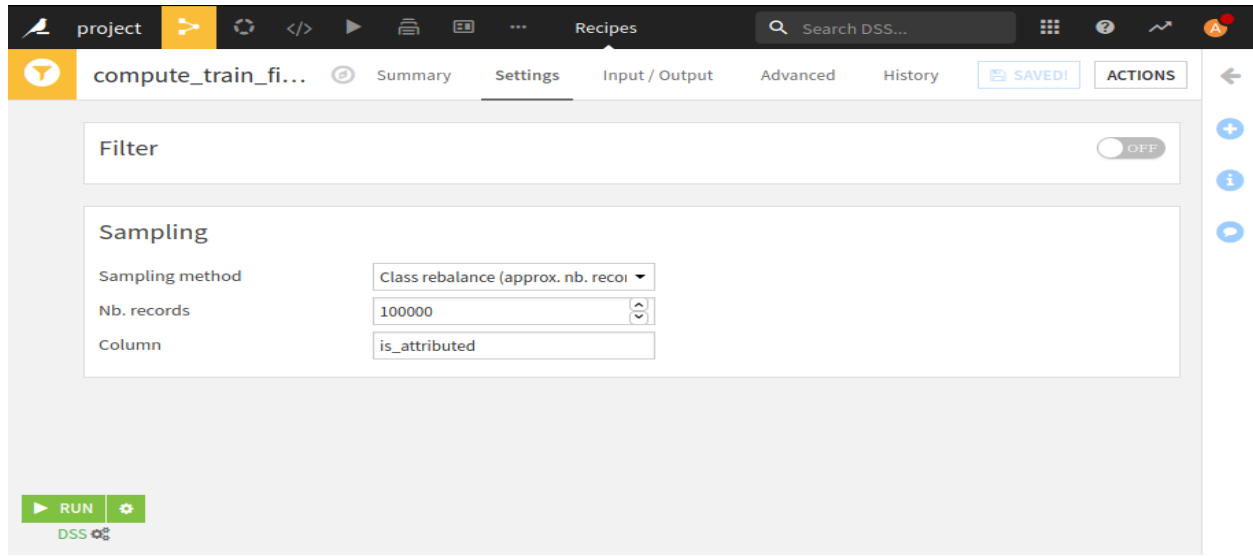
So, we worked only with 100k entries considering our system capabilities. Although the model is completely scalable and can be trained with a larger dataset in future.

# Basic Exploratory Data Analysis

It can be seen how skewed our dataset was from a first look. It was very highly unbalanced and biased towards the 0 value that corresponds to no attribution.
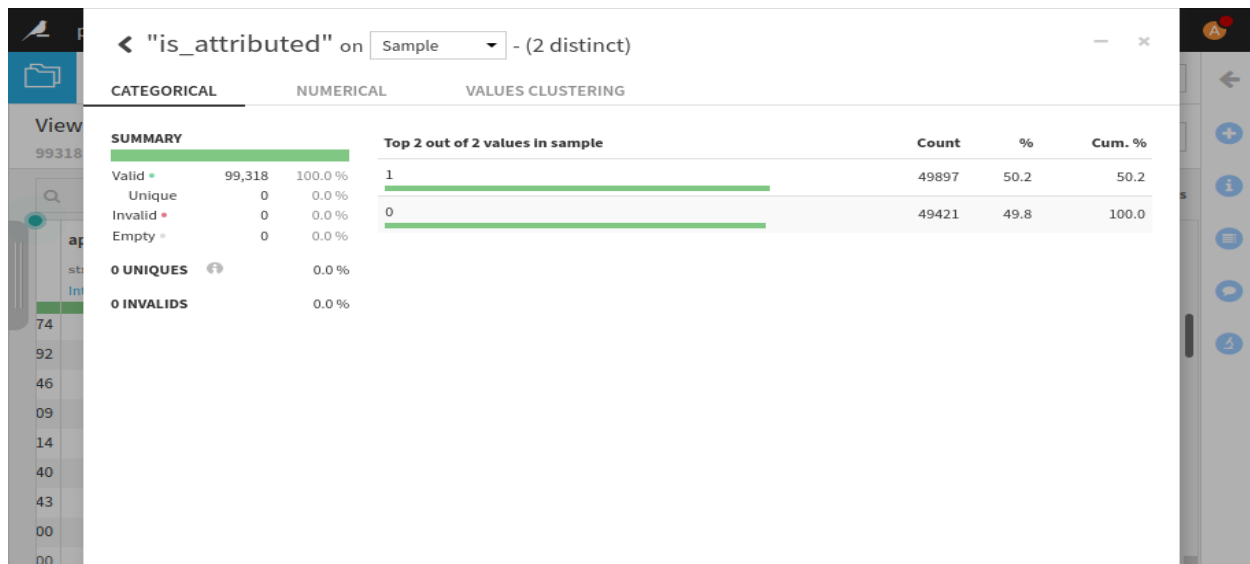


Our task was to balance the dataset for creating a more accurate model which we achieved using Dataiku sampling recipe.



We used class rebalance sampling method for balancing the is_attributed column and for filtering out 100k such entries from the dataset.

After performing sampling we got a new balanced dataset with almost equal values of 0 and 1 in is_attributed column.

Analysis of attributed_time showed that it contains mostly null or no values and hence was not helpful for prediction.



So, this feature was dropped and was not included for model creation.

# Feature Engineering

Dataset had some features through which certain properties could be extracted. It had click_time column through which date components were selected.
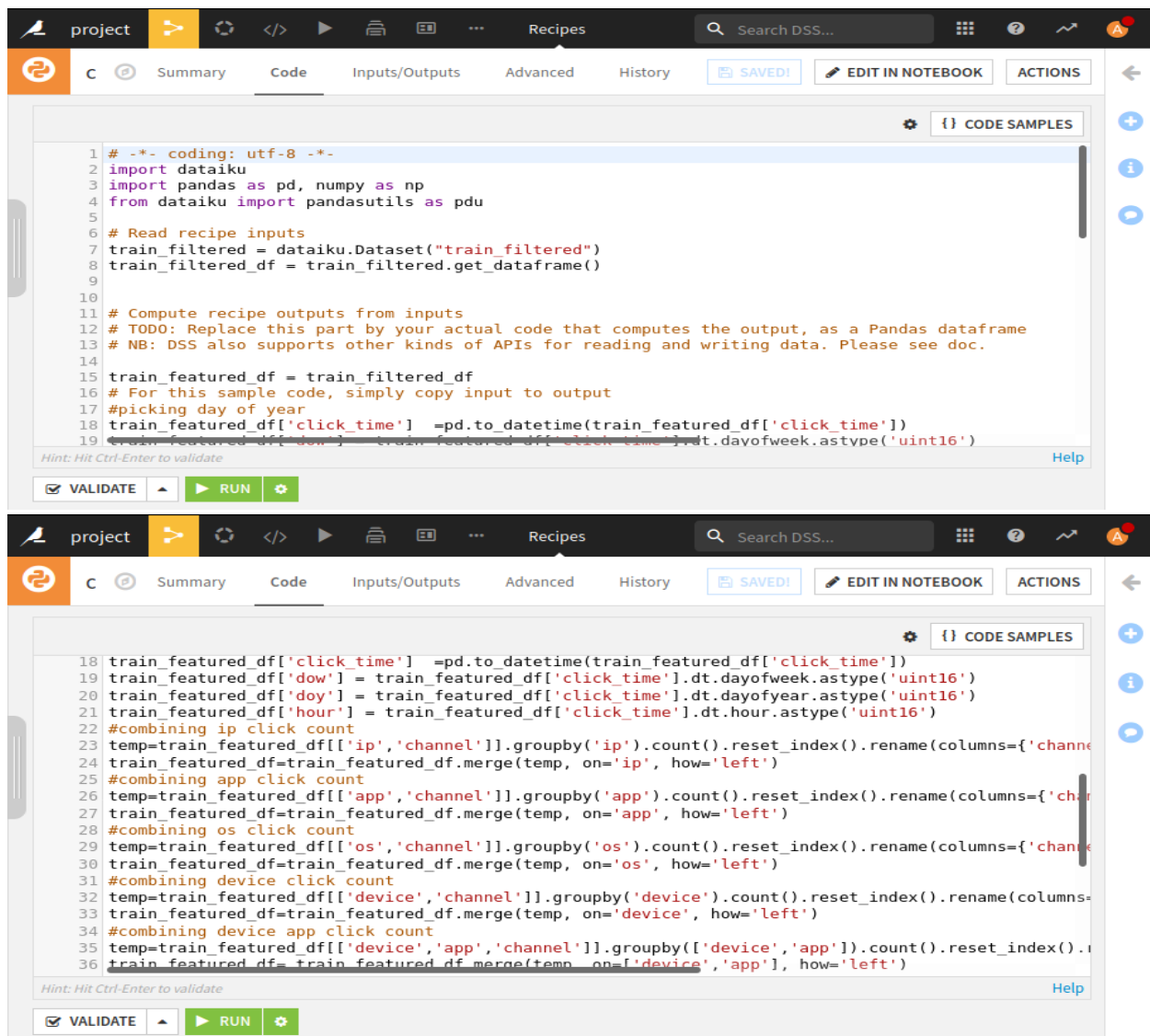
Date Components selected were:-

hour,doy(day of year),dow(day of week)

Groupby Aggregations done are by click count of:-

ip, os, device, app , ip and app, os and app, device and app

The python code recipe used was:-

```
33  train_featured_df=train_featured_df.merge(temp, on='device', how='left')
34  #combining device app click count
35  temp=train_featured_df[['device','app','channel']].groupby(['device','app']).count().reset_index().
36  train_featured_df= train_featured_df.merge(temp, on=['device','app'], how='left')
37  #combining ip app click count
38  temp=train_featured_df[['ip', 'app', 'channel']].groupby(['ip', 'app']).count().reset_index().rename
39  train_featured_df= train_featured_df.merge(temp, on=['ip','app'], how='left')
40  #combining os app click count
41  temp=train_featured_df[['os','app','channel']].groupby(['os','app']).count().reset_index().rename(co
42  train_featured_df=train_featured_df.merge(temp, on=['os','app'], how='left')
43  #dropping click_time
44  train_featured_df.drop(['click_time','attributed_time'],axis=1,inplace=True)
45
46
47  # Write recipe outputs
48  train_featured = dataiku.Dataset("train_featured")
49  train_featured.write_with_schema(train_featured_df)
50
```

After running this code the dataset looked like:-



This dataset was now ready to be served to ML models.

# Choosing Model

For building model, we took the advantage of Dataiku DSS auto ML capabilities. We chose 6 different classifiers for building our models which are:-

- Random Forest Classifier
- Logistic Regression
- XGBoost Classifier
- Decision Tree
- K Nearest Neighbor
- Artificial Neural Network

The training looked like:-

We got the best results for Random Forest Classifier, detailed metrics of all the classifiers as:-



So, we chose Random Forest Classifier for our final model.

**Random Forest Model Decision Tree:-**



**Confusion Matrix and Detailed Metrics:-**

# Final Flow

Our final flow after test predictions looked like this:-



## Test Dataset predictions:-

# How to use the model as an API?

For using Dataiku API capabilities we must have DSS License. So, we had to figure out a way through which we could use this model for predictions outside Dataiku DSS.

**Steps involved were**:-

1. Click the create API option.



2. Download it as an API package; we could not use API deployer because of license limitations.

3. After downloading it, we found a pickle file. This file can be used to make predictions for the test dataset outside DSS.

# Flask API Documentation

## Introduction

**Flask** is a micro web framework which has a minimal footprint compared to others like Django and allows us to select which components and modules we want to integrate with it. Flask is highly reliable and performant.

The Flask Restful was chosen on the suggestion of Mentor when our team was unable to get an API from Dataiku

## Project Layout

Project layout looks like the following folder structure.



API module will host our application code, from models-

- App.py serves as the main file and hosts the server whereas models.py serves as the backend for the training.

- App.py has 3 classes each serving separate calls to provide a different response to calling application based on request

- Train_filtered.csv is the dataset used for training the module

- Test_new.csv is created from the received post request

- Usage and Scores file contains the data about the system resources

```python
def feature_engineering(df):
    df['dow'] = df['click_time'].dt.dayofweek.astype('uint16')
    df['doy'] = df['click_time'].dt.dayofyear.astype('uint16')
    df['hour'] = df['click_time'].dt.hour.astype('uint16')
    features_clicks = ['ip', 'app', 'os', 'device']

    for col in features_clicks:
        col_count_dict = dict(df[[col]].groupby(col).size().sort_index())
        df['{}_clicks'.format(col)] = df[col].map(col_count_dict).astype('uint16')

    features_comb_list = [('app', 'device'), ('ip', 'app'), ('app', 'os')]
    for (col_a, col_b) in features_comb_list:
        df1 = df.groupby([col_a, col_b]).size().astype('uint16')
        df1 = pd.DataFrame(df1, columns=['{}_{}_comb_clicks'.format(col_a, col_b)]).reset_index()
        df = df.merge(df1, how='left', on=[col_a, col_b])
    return df
```

Fig - Feature Engineering based on findings from Dataiku to get best results for prediction



Fig – Api.py window while the server is running

16

```python
class MakePrediction(Resource):
    @staticmethod
    def post():
        posted_data = request.files['upload_file']
        posted_data.save(posted_data.filename)

        df1 = pd.read_csv('test_new.csv', parse_dates=['click_time'])

        test_feature = feature_engineering(df1)
        time = test_feature[['click_time']]
        test_feature.drop(['click_time', 'click_id'], axis=1, inplace=True)

        prediction = model.predict(test_feature)
        test_feature['predicted'] = prediction
        test_feature['click_time'] = time
        test_predicted = test_feature[['ip', 'app', 'device', 'os', 'channel', 'click_time', 'predicted']]

        test_predicted.to_csv('test_final.csv', index=False)
        #print(test_predicted.head())

        file1 = open("usage.txt", "r+")
        f=file1.readline()

        return test_predicted.to_json()
```

Fig – Api.py – Make Prediction class which takes the incoming file from post request convert it to CSV and perform predictions and then return file as a json file

```python
class Usage(Resource):
    @staticmethod
    def get():

        file1 = open("usage.txt", "r+")
        f=file1.readline()

        return f

class Score(Resource):
    @staticmethod
    def get():

        file1 = open("score.txt", "r+")
        f=file1.read()
        print(f)

        return f
```

Fig – Api.py – Usage and Score class which return test files created during training n response to GET request

17

```
def train_model():
    df = pd.read_csv('train_filtered.csv', parse_dates=['click_time'])
    df_feature = feature_engineering(df)
    df_feature.drop(['click_time', 'attributed_time'], axis=1, inplace=True)
    X = df_feature.drop(['is_attributed'], axis=1)
    y = df_feature['is_attributed']

    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=1337)
    classifier = RandomForestClassifier(n_estimators=50, n_jobs=-1)
    start = time.time()
    classifier.fit(X_train, y_train)
    stop = time.time()
    usage = open('usage.txt', 'w')
    usage.write('Training Time:- {} sec'.format(stop - start))
    usage.close()
    predictions = classifier.predict(X_test)

    accuracy = accuracy_score(y_test, predictions)

    print('Model Training Finished./n/tAccuracy obtained: {}'.format(accuracy))

    file = open('model.pkl', 'wb')
    pickle.dump(classifier, file)
    file.close()

    score = open('score.txt', 'w')
    score.write('Confusion_matrix:-\n{}\n\nClassification_report:-\n{}\n\nROC_AUC:-\n{}\n\nAccuracy:-\n{}\n\n'.format(
        confusion_matrix(y_test, predictions), classification_report(y_test, predictions),
        roc_auc_score(y_test, predictions), accuracy_score(y_test, predictions)))
```

Fig – Model.py train_model function which trains based on the random forest algorithm and creates a pickle file which would be used in the prediction

# Dashboard

Django based web app to upload a CSV file and give visual analytics on the given file. The team was asked to build a Django based web app to give meaningful results to the user along with usage of system resource



Fig – Directory structure of Django dashboard

```
INSTALLED_APPS = [
    'django.contrib.admin',
    'django.contrib.auth',
    'django.contrib.contenttypes',
    'django.contrib.sessions',
    'django.contrib.messages',
    'django.contrib.staticfiles',
    'firstUI'
]
```

Fig – Adding the name of the UI in installed apps in settings.py

```
from django import forms
class StudentForm(forms.Form):
    file        = forms.FileField() # for creating file input
```

Fig – Forms.py

```
def indexPage(request):
    csv = "test_new.csv"
    files = {'upload_file': open(csv, 'rb')}
    response = requests.post('http://127.0.0.1:5000/predict', files=files)
    response2 = requests.get('http://127.0.0.1:5000/usage')
    response3 = requests.get('http://127.0.0.1:5000/score')
```

Fig – Views.py - Making POST and GET requests to the API serve and saving the response

```
testfile = pd.read_csv('test_new.csv')
time_table = testfile[['click_time', 'click_id']].groupby('click_time').count().sort_values('click_time',ascending=True)
time_table_index_list = time_table.index.to_list()
time_table_values_list = time_table['click_id'].values.tolist()
seconds_minutes = [i[-4:] for i in time_table_index_list]
```

Fig – Reading the test file from insert.html and performing operations on it

```python
# stacked bar graph
testfile2 = pd.read_csv('test_res.csv')

testfile2['ip'] = testfile2['ip'].astype('category')
df3 = testfile2.groupby(['ip', 'predicted']).size().unstack(fill_value=0)
df3['counter'] = df3[1] + df3[0]

df3 = df3.sort_values('counter', ascending=False).nlargest(10, 'counter')
df3_0_list = df3[0].values.tolist()
df3_1_list = df3[1].values.tolist()
df3_label=df3.index.to_list()

# stacked bar graph 2
testfile_channel = pd.read_csv('test_res.csv')

testfile_channel['channel'] = testfile_channel['channel'].astype('category')
df5 = testfile_channel.groupby(['channel', 'predicted']).size().unstack(fill_value=0)
df5['counter'] = df5[1] + df5[0]

df5 = df5.sort_values('counter', ascending=False).nlargest(10, 'counter')
df5_0_list = df5[0].values.tolist()
df5_1_list = df5[1].values.tolist()
df5_label = df5.index.to_list()
```

Fig – Creating graphs from predicted files

```python
#Bar Graph

df4 = testfile2.groupby(['ip', 'device', 'predicted']).size().unstack(fill_value=0)
df4['counter'] = df4.sum(axis=1)
df4_suspicious = df4[df4[1] <= 1].sort_values('counter', ascending=False).nlargest(10, 'counter')
df4_0_list = df4_suspicious['counter'].values.tolist()
list_ip_device=df4_suspicious.index.tolist()
new_list = []
for i in list_ip_device:
    new_list.append(''.join(str(i)))


#''''''''''''''''''''''''''''''''''''''

usage=response2.text[1:-2]
score =response3.iter_lines()
print(score)

dataset = pd.read_json(response.json())
dataset.to_csv('test_res.csv', index=False)

cpu=psutil.cpu_percent()
cpu_util=[cpu,100-cpu]
memory=psutil.virtual_memory().percent
memory_util=[memory,1-memory]
```

Fig – Saving the values of resource utilization in variables

```
testfile2 = pd.read_csv('test_res.csv')
time_table2 = testfile2[['click_time', 'predicted']].groupby('click_time').sum().sort_values('click_time',
                                                                                              ascending=True)

time_table_index_list2 = time_table2.index.to_list()
time_table_values_list2 = time_table2['predicted'].values.tolist()
totala = testfile2[testfile2['predicted'] == 1]
totaln = testfile2[testfile2['predicted'] == 0]
pied=[len(totala.index), len(totaln.index)]
seconds_minutes2 = [i[-4:] for i in time_table_index_list2]

vara=seconds_minutes
varb=time_table_values_list
varc=seconds_minutes2
vard=time_table_values_list2


context={'varb': varb, "vara":vara,"varc":varc,"vard":vard, 'usage':usage, 'pied':pied, 's0':df3_0_list, 's1':df3_1_list,
         'ip':df3_label, 'ip_device':new_list, 'ip_device_values':df4_0_list, 'c0':df5_0_list, 'c1':df5_1_list,
         'channel':df5_label, 'cpu_util':cpu_util, 'memory_util':memory_util, 'score':score}
return render(request, 'index.html', context)
```

Fig – Passing various variables as context ad rendering the page

```
def index(request):
    if request.method == 'POST':
        s = SForm(request.POST, request.FILES)
        if s.is_valid():
            handle_uploaded_file(request.FILES['file'])
            return indexPage(request)
    else:
        s = SForm()
        return render(request, "insert.html", {'form':s})
```

Fig – Create a form for the first page which accepts CSV file as input

```
def handle_uploaded_file(f):
    print(os.getcwd())
    with open(''+f.name, 'wb+') as destination:
        for chunk in f.chunks():
            destination.write(chunk)
```

Fig –functions.py – function to save the file as a CSV file on the server

```
from django.contrib import admin
from django.urls import path, include
from django.conf.urls import url
from firstUI import views

urlpatterns = [
    path('admin/', admin.site.urls),
    path('',views.index,name='index'),
]
```

Fig –url.py – urlpatterns save the name of files to be called

```
d>python manage.py runserver
```
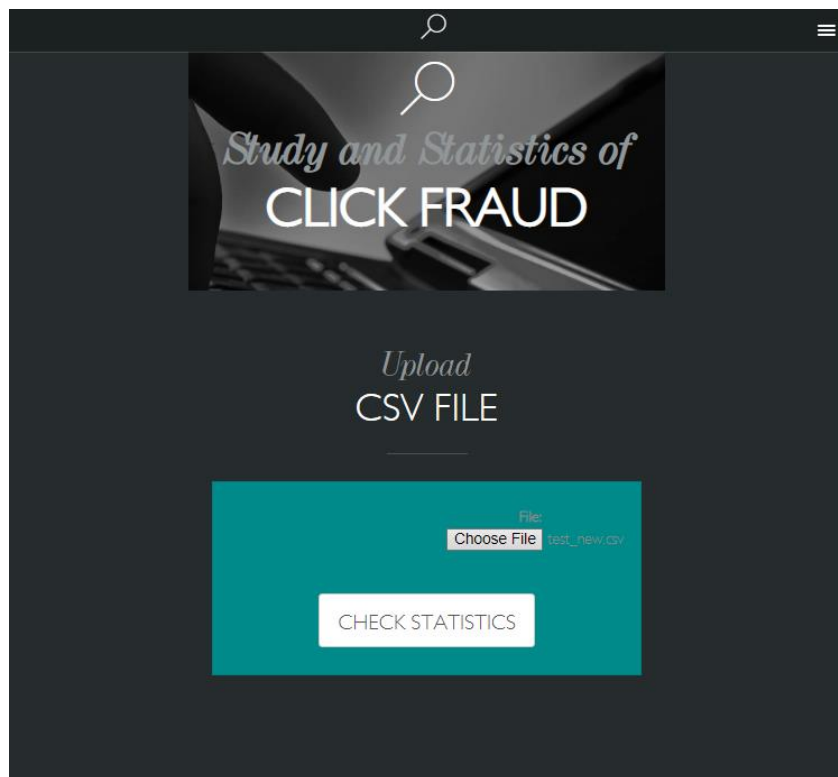
Fig – terminal command to run the server



Fig – First page on opening the URL – it takes the CSV file as input

# Data Visualization

The next page shows the training time for the model

# Dashboard

**Training Prediction** –

**Training Time:- 0.5966494083404541 sec**

Fig –This part shows the training time required for the model
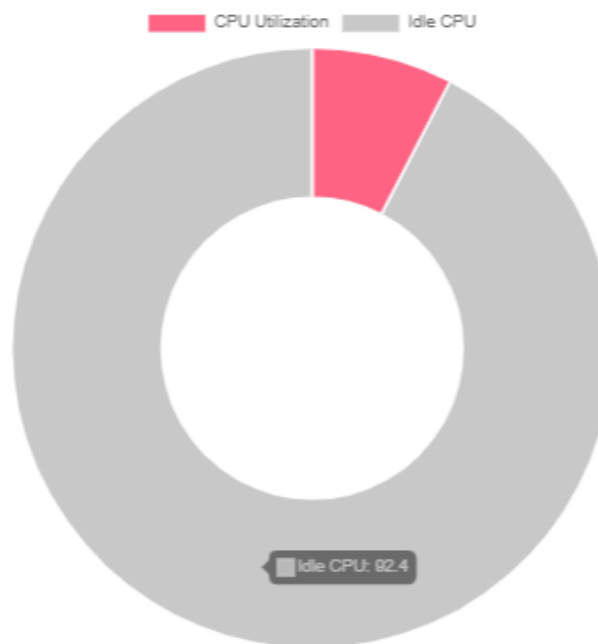
## CPU Utilization



Fig. - This doughnut-shaped graph shows the CPU Utilization while running the server in percentage
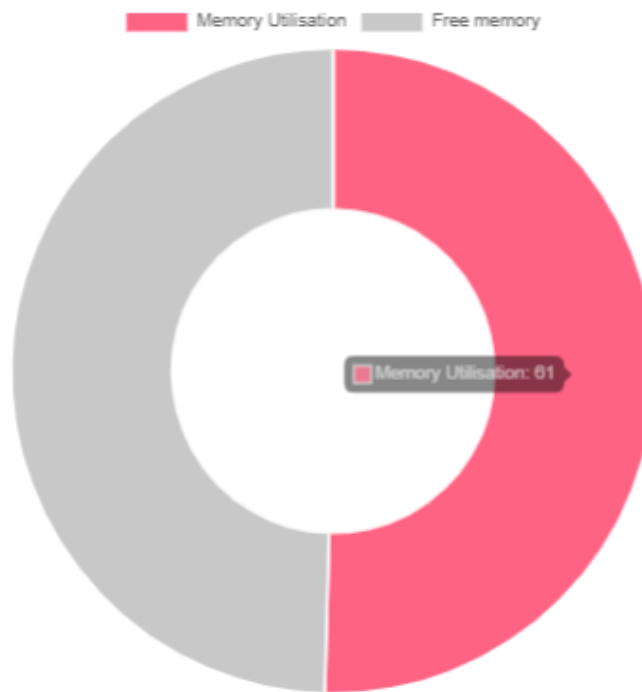
## Memory Utilization



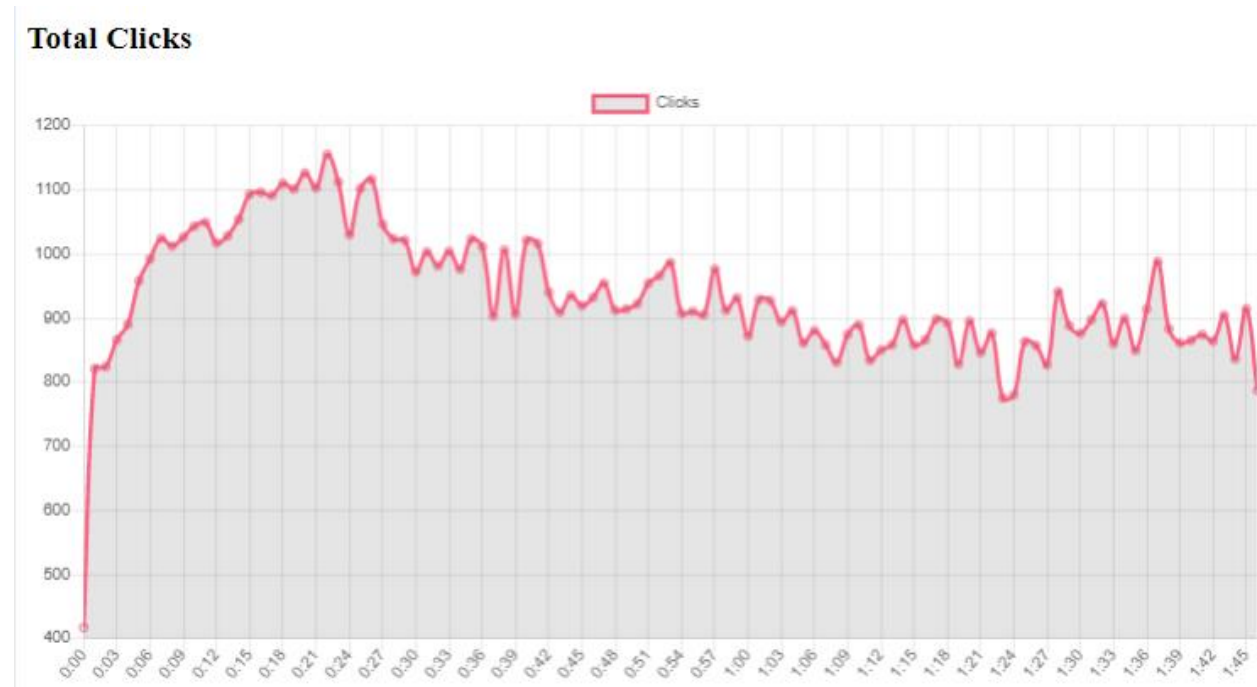Fig - This doughnut shows the Memory Utilization while running the server in percentage

## Total Clicks



Fig – Graph showing the number of clicks with time on the x-axis

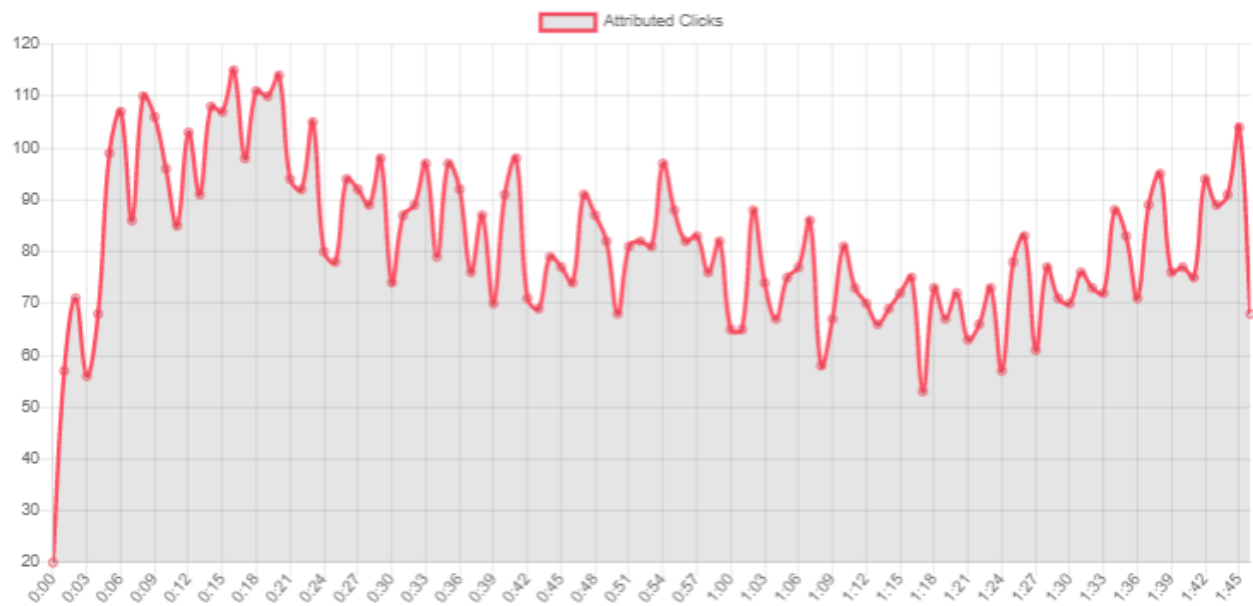## Predicted Attributed Clicks



Fig - Number of predicted Attributions(app download)
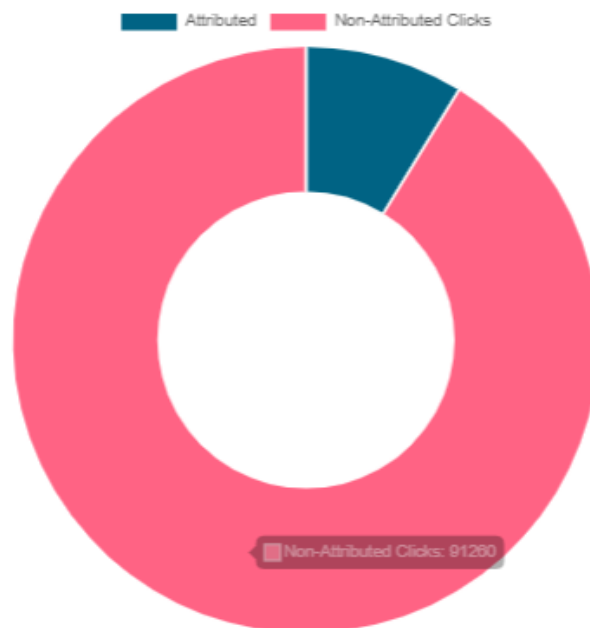
## Atrubuted vs Non-attributed Clicks

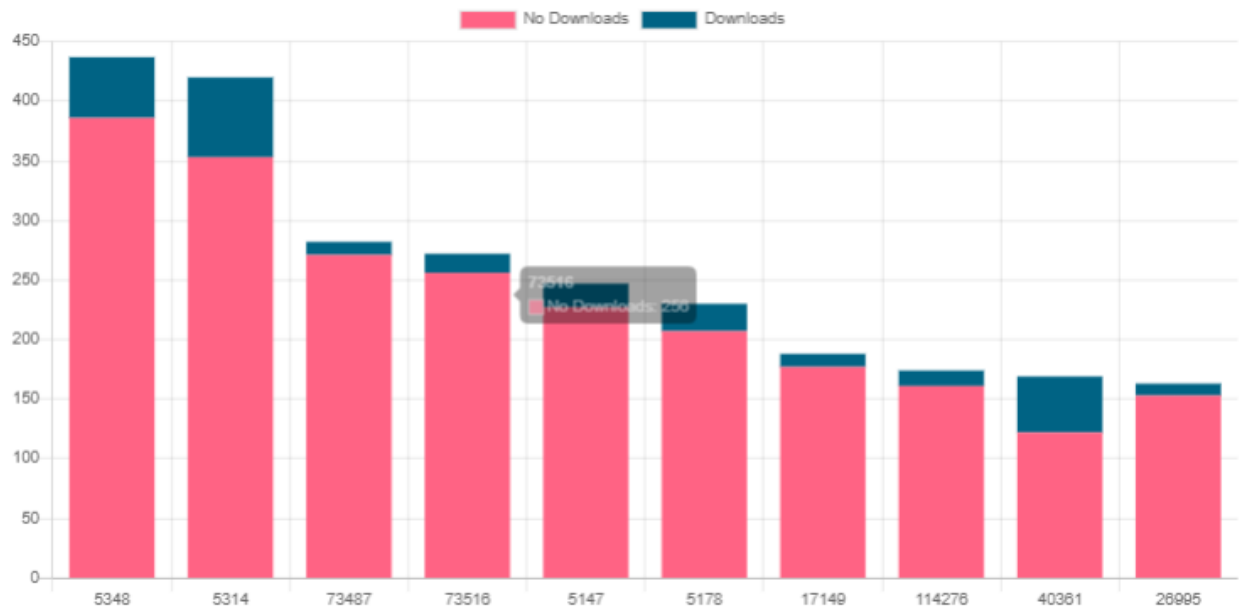Fig - Attributed to non attributed clicks

## IPs with most clicks



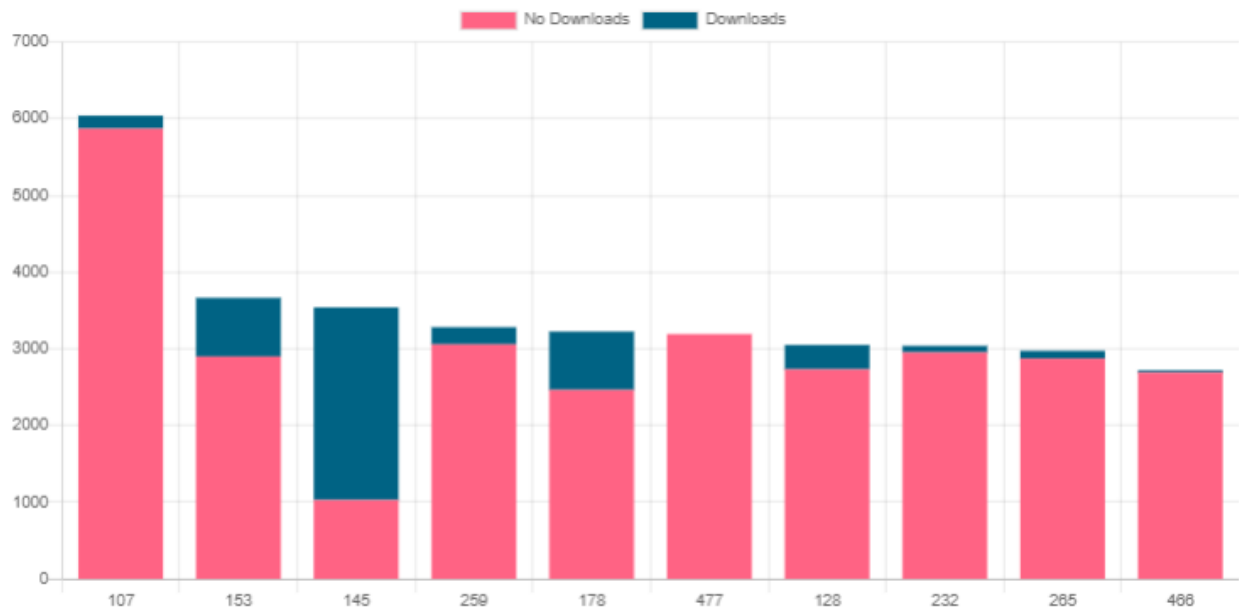Fig - IPs with the most number of clicks
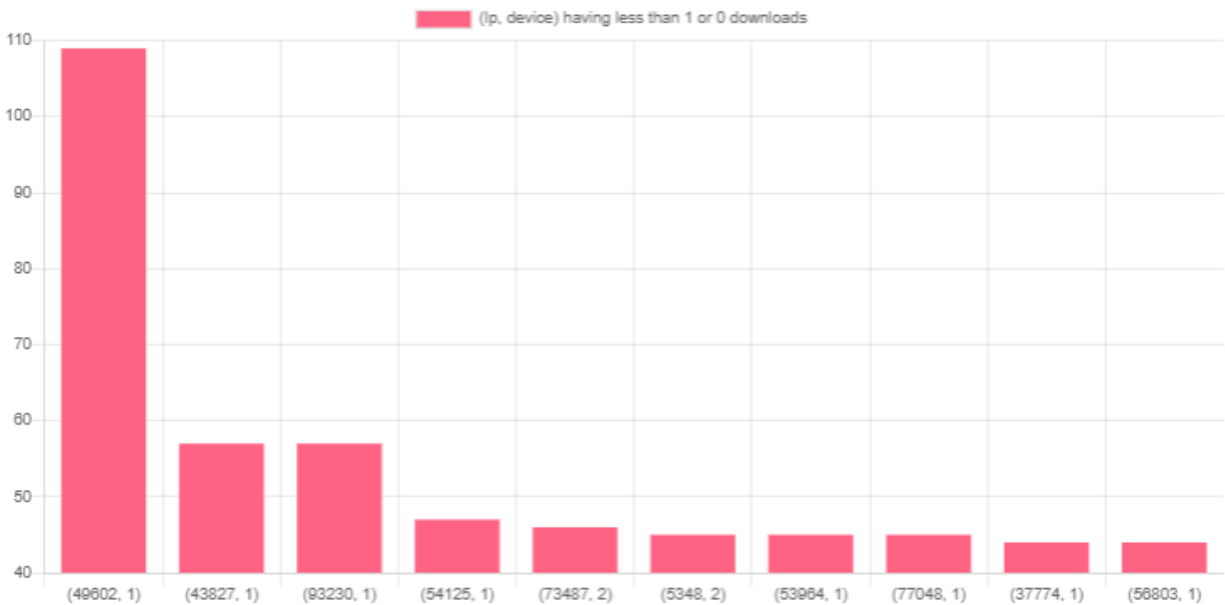
## Channel with most clicks

Fig - Suspected Frauds with 0 or 1 attributions

## Practical Implications

- Businesses can measure the success rate per channel and invest as per the requirements on each channel
- Businesses can reduce the spending on the channels which are no giving sufficient returns
- Certain IPs can be monitored for suspicious activities and added security measures can be implemented
- Resource utilization will give an idea of resources being consumed by servers