# **Group name:** Arjohi
# **Week 10 deliverables**

**Group members:**

**Name:** Arda Baris Basaran

**Email:** ardabarisbasaran@hotmail.com

**Country:** Turkey

**College/Company:** Middle East Technical University

**Specialization:** NLP


**Name:** Jose Luis Castañeda Terrones

**Email:** joseluiscastanedat@gmail.com

**Country:** Perú

**College/Company:** IFT-UNESP (São Paulo)

**Specialization:** NLP


**Name:** Hiten Chadha

**Email:** hitenchadha1995@gmail.com

**Country:** Denmark

**College/Company:** Technical University of Denmark

**Specialization:** NLP

**GitHub repository link:**
https://github.com/JoseLuisCastanedaT/dataglacier-week7-13

**Problem description & Business understanding:**

The term hate speech is understood as any type of verbal, written or behavioural communication that attacks or uses derogatory or discriminatory language against a person or group based on what they are, in other words, based on their religion, ethnicity, nationality, race, colour, ancestry, sex or another identity factor. In this problem, We will take you through a hate speech detection model with Machine Learning and Python.

Hate Speech Detection is generally a task of sentiment classification. So for training, a model that can classify hate speech from a certain piece of text can be achieved by training it on a data that is generally used to classify sentiments. So for the task of hate speech detection model, We will use the Twitter tweets to identify tweets containing Hate speech.

We will analyze a dataset CSV file from Kaggle containing 31,935 tweets. The dataset is heavily skewed with 93% of tweets or 29,720 tweets containing non-hate labeled Twitter data and 7% or 2,242 tweets containing hate-labeled Twitter data. We will try different classification algorithms after the preprocessing and data cleaning steps.

**Type of data:**
  - 1 boolean column (1 representing hate speech tweet and 0 non-hate speech tweet)
  - 1 string column (the tweet itself)
  - 1 numerical column (index column, representing the id)

| id | label | tweet |
|---|---|---|
| 1 | 0 | @user when a father is dysfunctional and is s... |
| 2 | 0 | @user @user thanks for #lyft credit i can't us... |
| 3 | 0 | bihday your majesty |
| 4 | 0 | #model i love u take with u all the time in ... |
| 5 | 0 | factsguide: society now #motivation |
| 6 | 0 | [2/2] huge fan fare and big talking before the... |
| 7 | 0 | @user camping tomorrow @user @user @user @use... |
| 8 | 0 | the next school year is the year for exams.ð... |
| 9 | 0 | we won!!! love the land!!! #allin #cavs #champ... |
| 10 | 0 | @user @user welcome here ! i'm it's so #gr... |
| 11 | 0 | â #ireland consumer price index (mom) climb... |
| 12 | 0 | we are so selfish. #orlando #standwithorlando ... |
| 13 | 0 | i get to see my daddy today!! #80days #getti... |
| 14 | 1 | @user #cnn calls #michigan middle school 'buil... |
| 15 | 1 | no comment! in #australia #opkillingbay #se... |
| 16 | 0 | ouch...junior is angryð#got7 #junior #yugyo... |
| 17 | 0 | i am thankful for having a paner. #thankful #p... |
| 18 | 1 | retweet if you agree! |
| 19 | 0 | its #friday! ð smiles all around via ig use... |
| 20 | 0 | as we all know, essential oils are not made of... |

**EDA:**

```
df_train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 31962 entries, 1 to 31962
Data columns (total 2 columns):
 #   Column  Non-Null Count  Dtype
---  ------  --------------  -----
 0   label   31962 non-null  int64
 1   tweet   31962 non-null  object
dtypes: int64(1), object(1)
memory usage: 749.1+ KB
```

**label = 0 is non-hate speech tweet, and label=1 is hate speech tweet**

```
df_train.groupby('label').count()['tweet'].reset_index().sort_values(by='tweet',ascending=False)
```
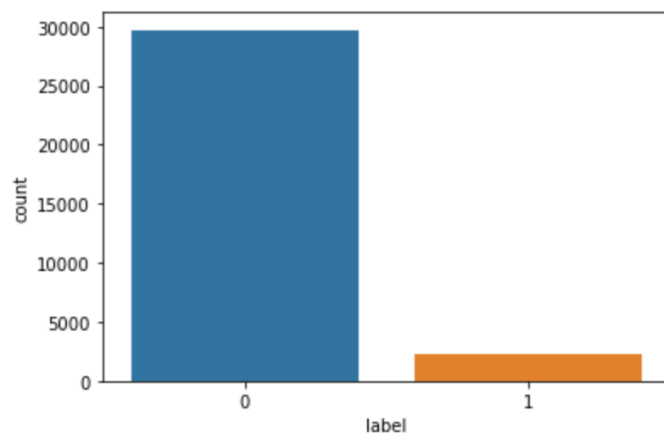
|   | label | tweet |
|---|-------|-------|
| **0** | 0 | 29720 |
| **1** | 1 | 2242 |

## Data is imbalanced, there is >90% of non-hate speech tweets

```
import seaborn as sns

sns.countplot(x='label', data=df_train)
```

```
<AxesSubplot:xlabel='label', ylabel='count'>
```

```
label_0 = len(df_train[df_train['label']==0])
label_1 = len(df_train[df_train['label']==1])
perc_0 = label_0/(label_0+label_1)*100
perc_1 = label_1/(label_0+label_1)*100

print(f'There is {label_1} hate speech tweets, which represents {perc_1:.2f}%')
print(f'There is {label_0} non hate speech tweets, which represents {perc_0:.2f}%')
```

```
There is 2242 hate speech tweets, which represents 7.01%
There is 29720 non hate speech tweets, which represents 92.99%
```

**Graphical representation of the most common words:**



**Cleaning the tweets:**

```
id
1      father dysfunctional selfish drag kid dysfunct...
2      thanks lyft credit cant use cause dont offer w...
3                                        bihday majesty
4                             model love u take u time ur
5                        factsguide society motivation
6      huge fan fare big talking leave chaos pay disp...
7                             camping tomorrow dannya
8      next school year year exam cant think school e...
9      love land allin cavs champion cleveland clevel...
10                                        welcome im gr
11     ireland consumer price index mom climbed previ...
12     selfish orlando standwithorlando pulseshooting...
13                        get see daddy today day gettingfed
14     cnn call michigan middle school build wall cha...
15     comment australia opkillingbay seashepherd hel...
16            ouchjunior angry got junior yugyoem omg
17                   thankful paner thankful positive
18                                       retweet agree
19     friday smile around via ig user cooky make people
20                    know essential oil made chemical
Name: tweet, dtype: object
```

**Most common words:**

| | Common_words | count |
|---|---|---|
| 0 | day | 2859 |
| 1 | love | 2802 |
| 2 | u | 1728 |
| 3 | happy | 1692 |
| 4 | amp | 1627 |
| 5 | time | 1244 |
| 6 | life | 1225 |
| 7 | like | 1200 |
| 8 | im | 1146 |
| 9 | today | 1085 |
| 10 | get | 1000 |
| 11 | new | 996 |
| 12 | thankful | 946 |
| 13 | positive | 931 |
| 14 | father | 920 |
| 15 | people | 875 |
| 16 | good | 869 |
| 17 | bihday | 854 |
| 18 | make | 847 |
| 19 | one | 843 |

**Most positive common words (non-hate speech):**

|    | Common_words | count |
|----|--------------|-------|
| 0  | day          | 2844  |
| 1  | love         | 2773  |
| 2  | happy        | 1680  |
| 3  | u            | 1634  |
| 4  | amp          | 1356  |
| 5  | time         | 1214  |
| 6  | life         | 1211  |
| 7  | im           | 1101  |
| 8  | today        | 1069  |
| 9  | like         | 1062  |
| 10 | get          | 949   |
| 11 | thankful     | 946   |
| 12 | positive     | 928   |
| 13 | new          | 926   |
| 14 | father       | 914   |
| 15 | bihday       | 854   |
| 16 | good         | 836   |
| 17 | make         | 809   |
| 18 | smile        | 804   |
| 19 | one          | 792   |

**Most negative common words (hate speech):**

| | Common_words | count |
|---|---|---|
| 0 | amp | 271 |
| 1 | trump | 211 |
| 2 | white | 155 |
| 3 | black | 149 |
| 4 | libtard | 149 |
| 5 | like | 138 |
| 6 | woman | 121 |
| 7 | racist | 112 |
| 8 | politics | 96 |
| 9 | u | 94 |
| 10 | liberal | 92 |
| 11 | allahsoil | 92 |
| 12 | people | 89 |
| 13 | might | 77 |
| 14 | sjw | 74 |
| 15 | hate | 73 |
| 16 | obama | 70 |
| 17 | new | 70 |
| 18 | dont | 66 |
| 19 | racism | 65 |