**Hochschule
Bonn-Rhein-Sieg**
University of Applied Sciences

# Aggressive Language Classifier

**Natural Language Processing**

**Date:** July 1st, 2020
**Team:** Manoj Kolpe Lingappa and Alan Gomez
**Github:** https://github.com/alanorlando95/Aggressive-language-classifier
**Google colab :** https://colab.research.google.com/drive/1AMwk1apDr-WZaKd66H00MfFB2XJcwUAU?usp=sharing
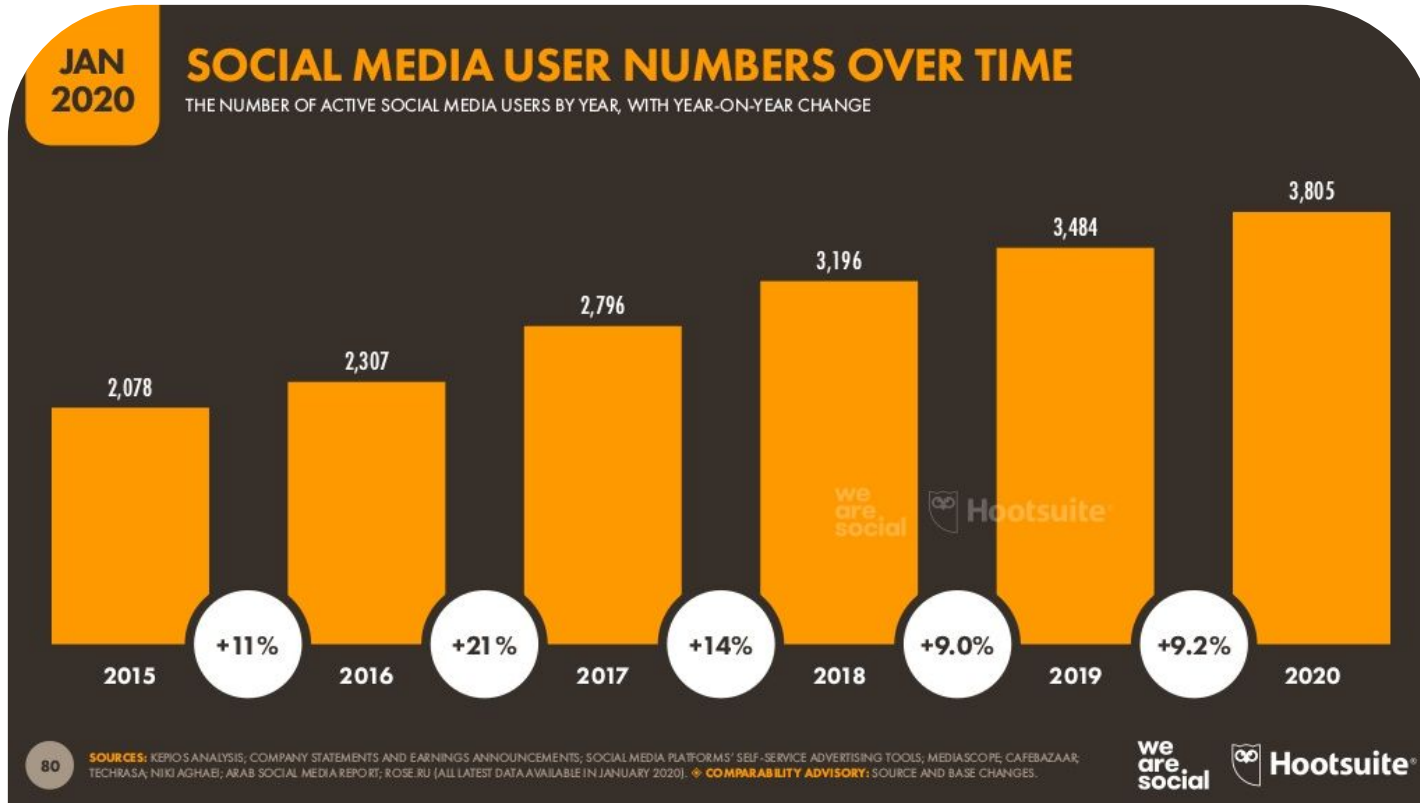
# Motivation



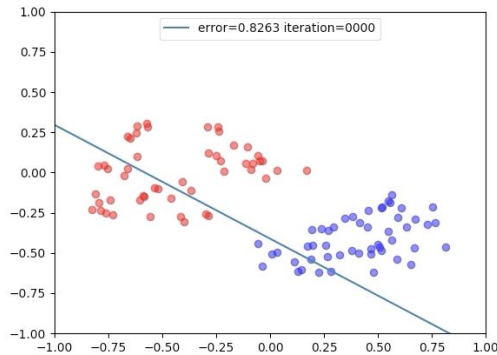Fig 1. Social media users over time (in millions) [1].

# Motivation

- 87 percent of young people in the USA have seen cyberbullying occurring online (2017 - 2019) [2].
- 36.5 percent of people in the USA feel they have been cyberbullied in their lifetime (2017 - 2019) [2].
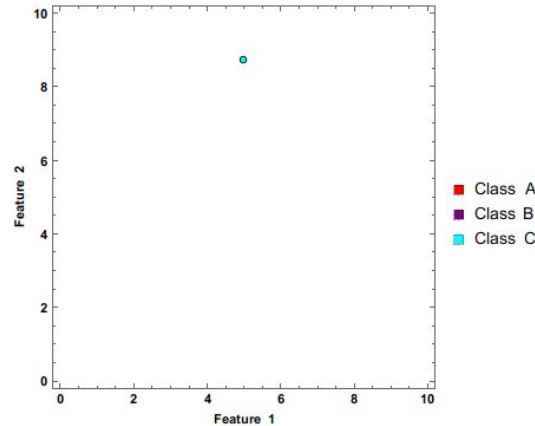
# Modern problems require modern solutions (maybe not too modern)

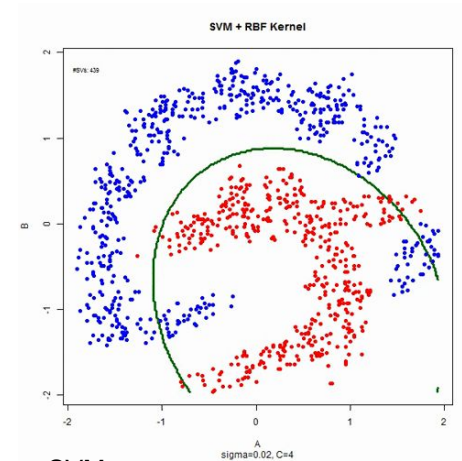- Logistic Regression
- Naive Bayes
- Support Vector Machine



Logistic regression

Credits:
https://brainbomb.org/Artificial-Intelligence/Machine-Learning/ML-Linear-Classification-Logistic-Regression/



Naive bayes

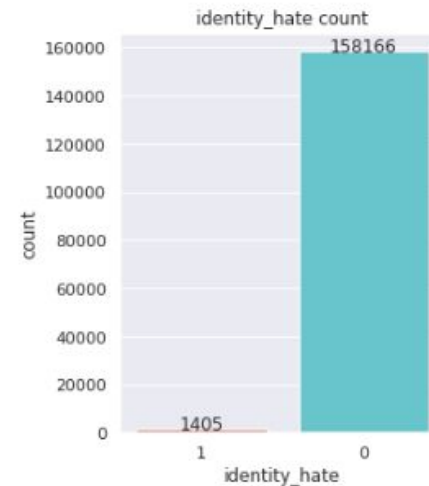Credits: https://commons.wikimedia.org/wiki/File:Naive_Bayes_Classifier.gif



SVM

Credits: https://gfycat.com/sentimentalthickbullfrog

# Dataset

- Jigsaw Multilingual Toxic Comment Classification, made up of English comments from Wikipedia's talk page edits [3].



Credits: https://mc.ai/detecting-toxic-comment/

# Toxic words

# Non toxic words

**Aggressive Language Classifier - Manoj, Alan**

Length histogram

Hochschule
Bonn-Rhein-Sieg
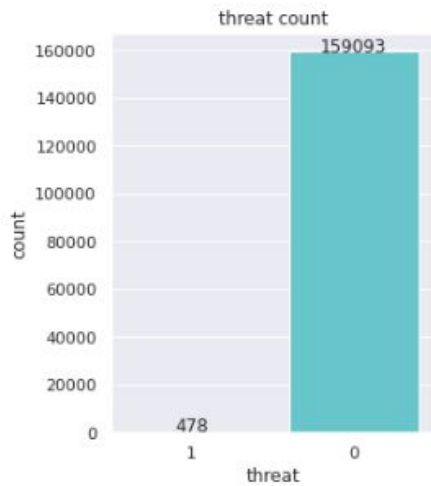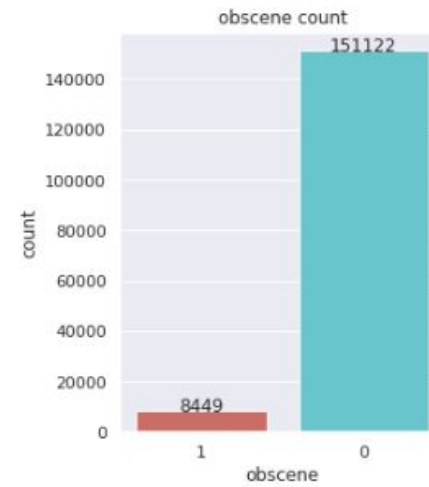University of Applied Sciences

# The implementation based on Gaydhani et. al [4]

- Train models by performing greedy search with different parameters:
    - CountVectorizer:
        - n-gram range ( (1, 1), (1, 2), and (1, 3) )
    - TfidfTransformer:
        - norm (L1 and L2)
- CountVectorizer extracts the n-gram features
- TfidfTransformer weight the n-gram features

Performance Comparison of Algorithms w.r.t different Features

Hochschule
Bonn-Rhein-Sieg
University of Applied Sciences

# The implementation based on Gaydhani et. al [4]

- Compare models
- Hyperparameter tuning:
  - MultinomialNB:
    - Additive (Laplace/Lidstone) smoothing parameter (0.01, 0.1, 1, and 10)
  - LogisticRegression:
    - Inverse of regularization strength (10 and 100)
    - Solver (newton-cg, liblinear, and saga)
  - SVM:
    - Scale of regularization term (0.00001, 0.0001, 0.001, 0.01, and 0.1)

Result of Naive Bayes for different hyperparameter values

Result of Logistic Regression for different hyperparameter values

Result of support vector machine for different hyperparameter values

Total data points: 156871
Type of comment: Toxic

| Google colab RAM: 12GB | Time taken to train w.r.t different features (min) | Time taken for Hyperparameter Tuning (min) |
|---|---|---|
| Naive Bayes | 42 | 7 |
| Logistic regression | 81 | 35 |
| Support vector machine | 52 | 10 |

Hochschule
Bonn-Rhein-Sieg
University of Applied Sciences

# SVM results (toxic)

```
Time taken: 10.101313591003418
               precision    recall  f1-score   support

           0       0.95      0.97      0.96     56684
           1       0.76      0.64      0.69      7294

    accuracy                           0.94     63978
   macro avg       0.86      0.81      0.83     63978
weighted avg       0.93      0.94      0.93     63978
```

Hochschule
Bonn-Rhein-Sieg
University of Applied Sciences

# Confusion matrix (toxic)

# Train the best model to classify the data into the other categories of the dataset

Input
sentence

Results

```
Hello, this a test sentnce. I don't want to be toxic
toxic : no
severe_toxic : no
obscene : no
threat : no
insult : no
identity_hate : no
```

# Implicit threats are not detected

Any difficulty and we will assume control but, when the looting starts, the shooting starts. Thank you!
toxic : no
severe_toxic : no
obscene : no
threat : no
insult : no
identity_hate : no

# Sensitive to unknown words

```
I kill people
toxic : yes
severe_toxic : no
obscene : no
threat : yes
insult : no
identity_hate : no

I am killing people
toxic : no
severe_toxic : no
obscene : no
threat : no
insult : no
identity_hate : no
```

Hochschule
Bonn-Rhein-Sieg
University of Applied Sciences

# Take Home Message

- The training data can influence the results on a classification task
  - Number of positive and negative examples
  - Unseen words

- One can train different models comparing different parameters using sklearn

- Different TFIDF norm and n gram range affects the validation accuracy.

- Hyperparameter tuning can give better validation accuracy than using gridsearch with respect to different features.

- Confusion matrix can be used to describe performance of a classification model.

- Traditional solver such as naive bayes, SVM and logistic regression produce reasonable performance in classification of toxic comment.

Hochschule
Bonn-Rhein-Sieg
University of Applied Sciences

# References

[1] Kemp, Simon. "Digital 2020: Global Digital Overview - DataReportal – Global Digital Insights." DataReportal. DataReportal – Global Digital Insights, January 30, 2020. https://datareportal.com/reports/digital-2020-global-digital-overview.

[2] "51 Critical Cyberbullying Statistics in 2020." BroadbandSearch.net. Accessed June 18, 2020. https://www.broadbandsearch.net/blog/cyber-bullying-statistics.

[3] "Jigsaw Multilingual Toxic Comment Classification." Kaggle. Accessed June 18, 2020. https://www.kaggle.com/c/jigsaw-multilingual-toxic-comment-classification/data.

[4] Gaydhani, Aditya, Vikrant Doma, Shrikant Kendre, and Laxmi Bhagwat. "Detecting hate speech and offensive language on twitter using machine learning: An n-gram and tfidf based approach." arXiv preprint arXiv:1809.08651 (2018).

Hochschule
Bonn-Rhein-Sieg
University of Applied Sciences

**Aggressive Language Classifier - Manoj, Alan**