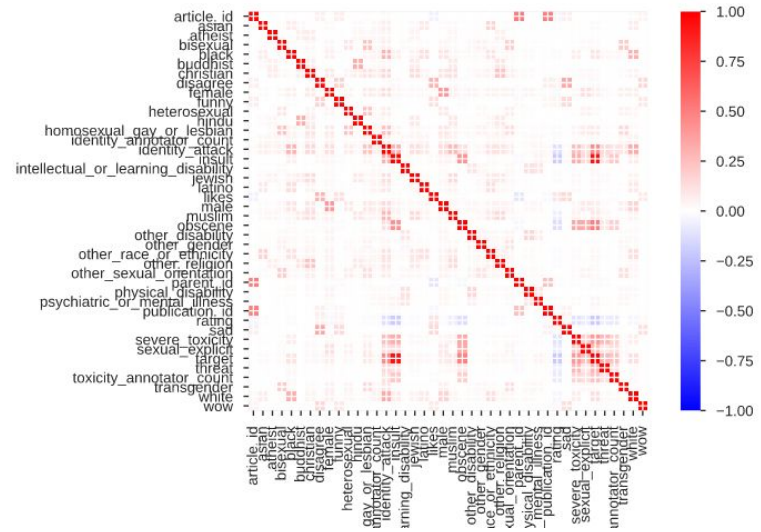# AstraZeneca AI Challenge

Team Name : CodingBeginner
Name : Rohith Srinivaas M
(SHA008208)
IIT Madras

# JigSaw Unintended Bias in Toxicity Classification

- Objective : Rate the toxicity in a given data instance
- Target prediction is a fraction denoting the toxicity level in the comment thereby classifying it as a positive or negative comment.
- Subtype attributes for toxicity:
    - Severe_toxicity
    - Obscene
    - Threat
    - Insult
    - Identity_attack
    - sexual_explicit

# Exploratory Data Analysis ( EDA )

- Pandas profiler was used for quick EDA of the dataset.
- Important observations are given below:
  - **target** is highly correlated with **insult** with ρ = 0.928206624.
  - The overlap between **white** and **black** identity comments is high.
  - A large number of comments about the **Jewish** identity is toxic towards the **Muslim** identity

# Approaches

- Using Logistic Regression
  - Using ELI5 to understand bias.
- Using textCNN
  - Using LIME to understand bias.

# Metrics

- SubGroup AUC :
  - only to the examples that mention the specific identity subgroup
  - *A low value in this metric means the model does a poor job of distinguishing between toxic and non-toxic comments that mention the identity.*
- BPSN (Background Positive, Subgroup Negative) AUC :
  - To the non-toxic examples that mention the identity and the toxic examples that do not
  - *A low value in this metric means that the model confuses non-toxic examples that mention the identity with toxic examples that do not, likely meaning that the model predicts higher toxicity scores than it should for non-toxic examples mentioning the identity.*

Reference : Nuanced Metrics for Measuring Unintended Bias with Real Data for Text Classification

# Metrics

- BNSP (Background Negative, Subgroup Positive) AUC :
  - To the toxic examples that mention the identity and the non-toxic examples that do not.
  - *A low value here means that the model confuses toxic examples that mention the identity with non-toxic examples that do not, likely meaning that the model predicts lower toxicity scores than it should for toxic examples mentioning the identity.*

Reference : Nuanced Metrics for Measuring Unintended Bias with Real Data for Text Classification

# Results and Discussion

- Using Logistic Regression
  - Used TweetTokenizer from nltk package with Tfidvectorizer to get the word-vector.
  - Performed Logistic Regression with the obtained word-vector with class = 1 for target >= 0.5 and vice versa
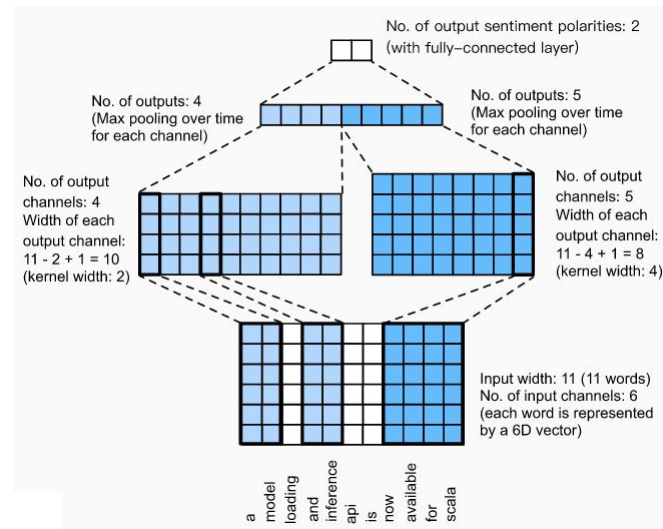  - Used ELI5 for model interpretation:

# Results and Discussion

- Using TextCNN
  - One of the important hyperparameters in textCNN is the sequence length.

|  | Number of words | Sequence length | score |
|---|---|---|---|
| 0 | 50000 | 150 | 0.9073 |
| 2 | 100000 | 150 | 0.9096 |
| 1 | 50000 | 300 | 0.9142 |
| 3 | 100000 | 300 | 0.9175 |

**y=0** (probability **1.000**, score **-8.223**) top features

| Contribution? | Feature |
|---|---|
| +7.550 | Highlighted in text (sum) |
| +0.673 | <BIAS> |

they stood was that not enough ? now , can we please put this to bed ? stop the crying and move on < sheesh >



No. of output sentiment polarities: 2
(with fully-connected layer)

No. of outputs: 4
(Max pooling over time
for each channel)

No. of outputs: 5
(Max pooling over time
for each channel)

No. of output channels: 4
Width of each output channel:
11 - 2 + 1 = 10
(kernel width: 2)

No. of output channels: 5
Width of each output channel:
11 - 4 + 1 = 8
(kernel width: 4)

Input width: 11 (11 words)
No. of input channels: 6
(each word is represented by a 6D vector)

a model loading and inference api is now available for scala

# Improvements

- Using pretrained Word Embeddings like BERT, FastText, Glove to predict the toxicity and then use weighted sum of these vectors ( with subgroup bias ) to get the meta word-vector with minimum bias which can be used for further classification.

Thank You