# Toxic Comments Classification

—

ECE 539 Project-Fall 2018
Deepan Das(MS-ECE, UW Madison)

# Executive Summary: *Problem at hand*
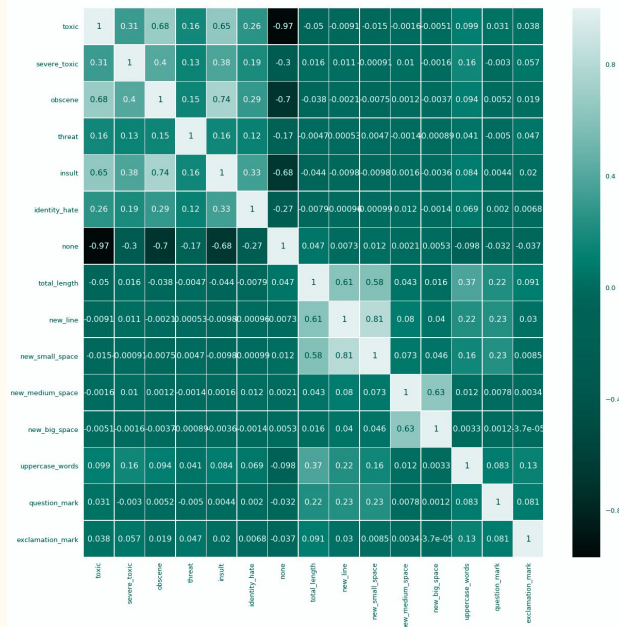
**Objective:**

Explore and develop models for **toxic comment detection** based solely on **textual information.**
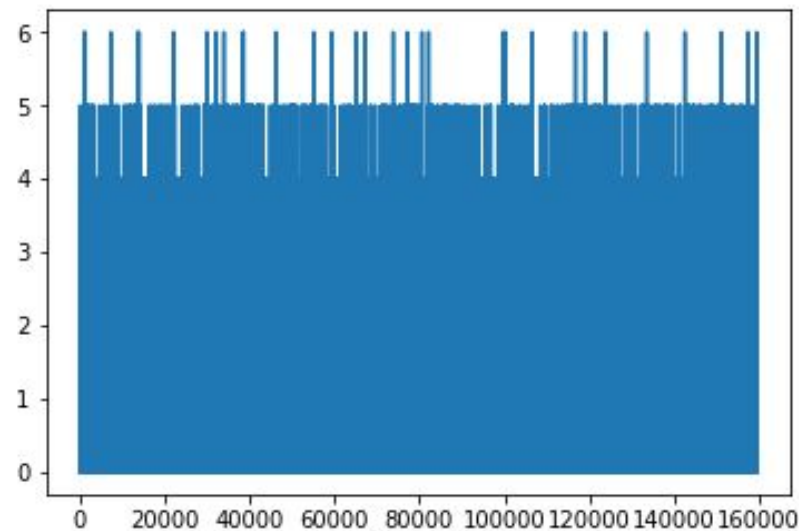
**Self-Assessment:**

- Implemented 3 distinct models
- Used open-source code to build these models, but developed on them later with visible improvements
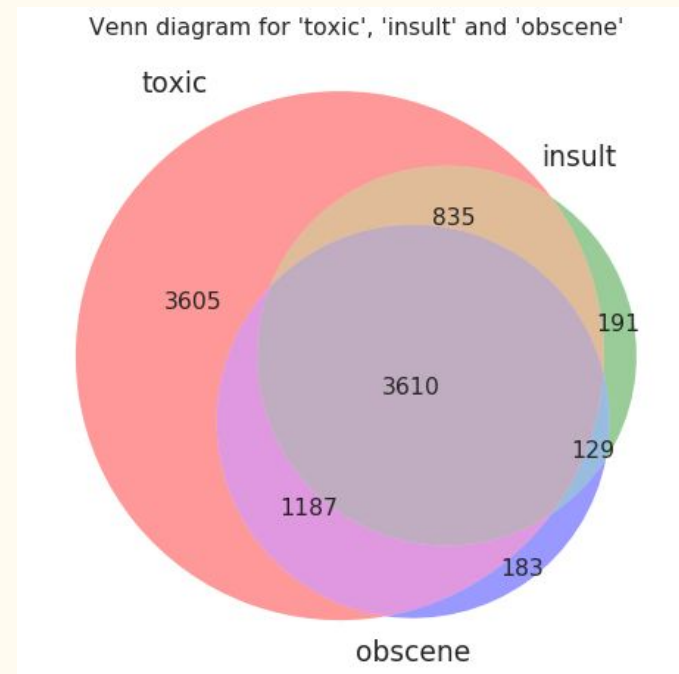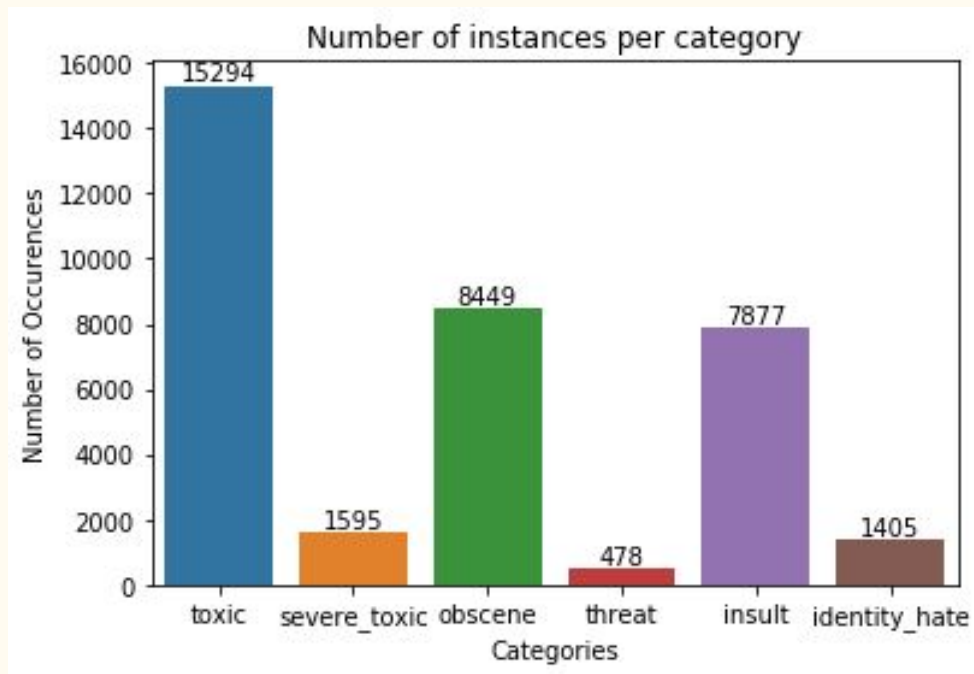
# Dataset: EDA and Insights


Correlation between attributes


Number of classes each comment belongs to

# Dataset: EDA and Insights



Number of instances per category



Venn diagram for 'toxic', 'insult' and 'obscene'
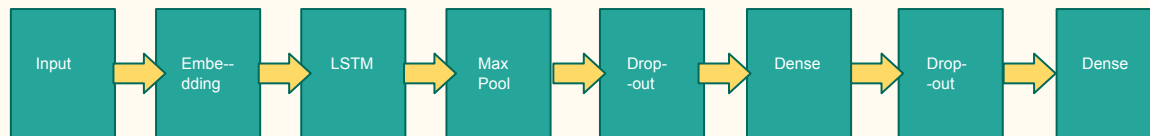
# Dataset: EDA and Insights

- Choice of **Embedding** essential
- No CNNs !
- Tendency for comments with more sentences to be toxic.

# Approach I: Logistic Regression

- Feature used: **TF-IDF**
- Used Python's in-built Scikit Library Functions
- **Kaggle Score: 0.970**

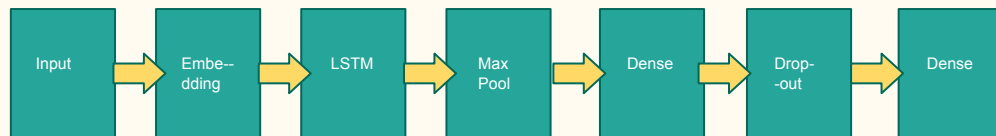# Approach II: Bi-LSTM, Random Embedding

- Training Accuracy: **97.01%**
- Validation(Testing) Accuracy: **96.89%**
- **Random Embedding Layer(128-d)**



[Inspiration: Wei-Yeoung Seow(Top 100 Kaggler). **Change in Embedding layer.]**
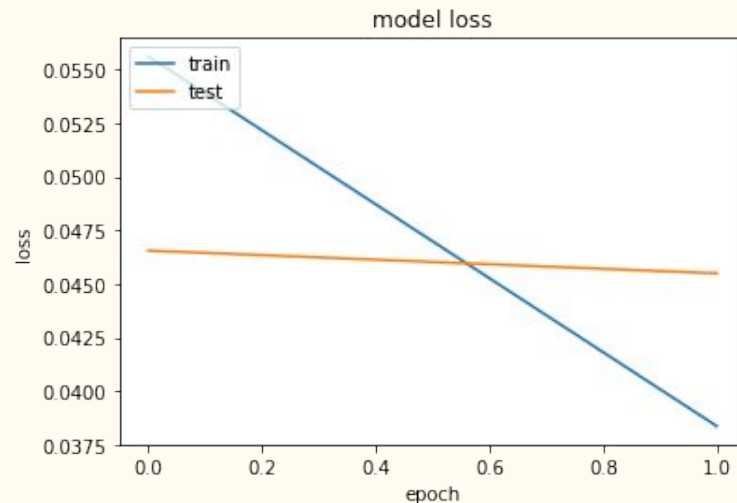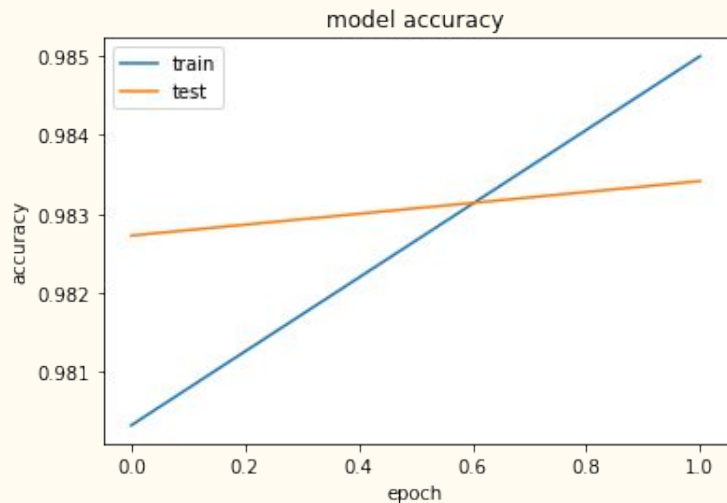
# Approach III: Bi-LSTM, GloVe Embedding

- Training Accuracy: **98.50%**
- Validation(Testing) Accuracy: **98.35%**
- **GloVe 600b20D Embedding**, Gigaword+Wikipedia Embedding



Inspiration: **Jeremy Howard, FastAI(USF)[Change in Embedding and Dropout to improve performance]**
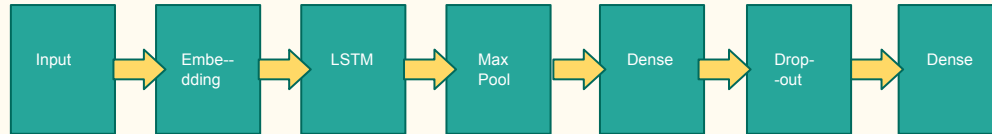
# Approach III: Bi-LSTM, GloVe Embedding

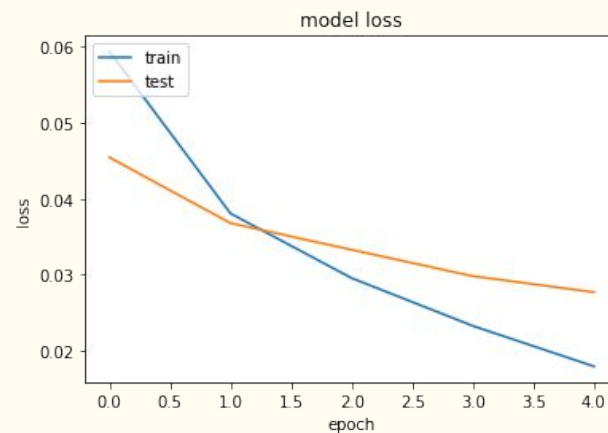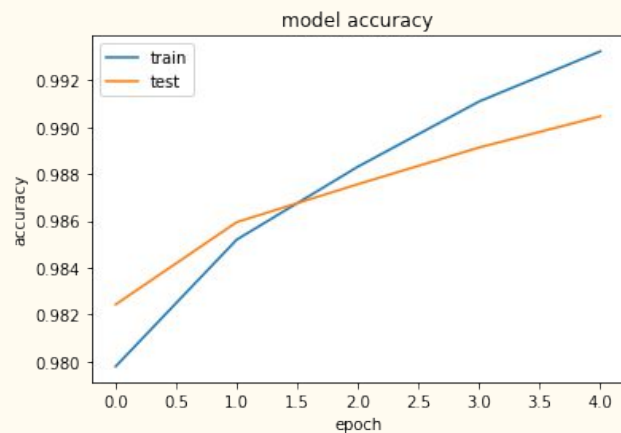Problem at hand: Overfitting by the **2nd Epoch**

# Approach IV: Bi-LSTM, GloVe, Data Augmentation

- Training Accuracy: **99.11%**
- Testing Accuracy: **98.91%**
- TextBlob's Machine translation tool.

| Input | | Embe-- dding | | LSTM | | Max Pool | | Dense | | Drop- -out | | Dense |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

**\*Independent Development**

# Approach IV: Bi-LSTM, GloVe, Data Augmentation
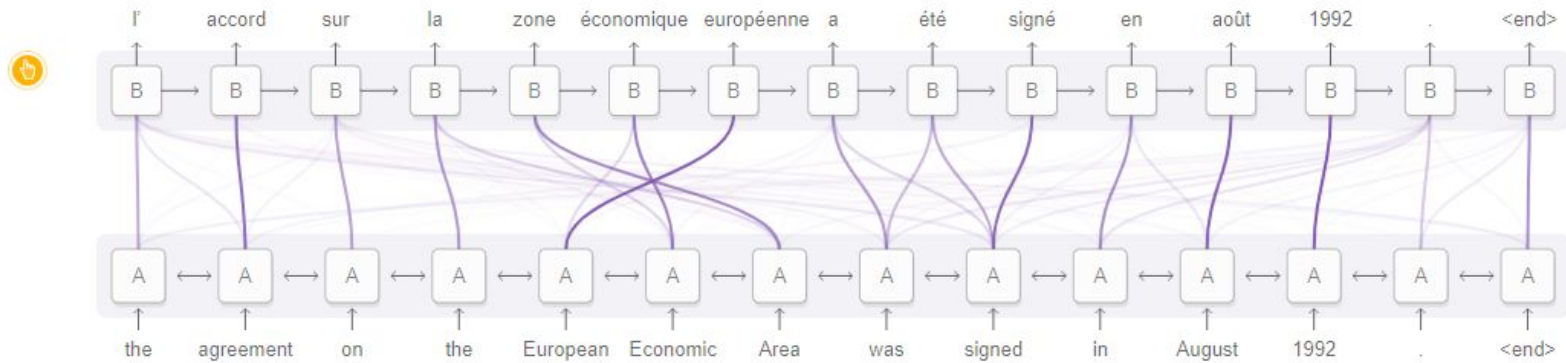
# Approach V: Hierarchical Attention Network



Diagram derived from Fig. 3 of Bahdanau, et al. 2014

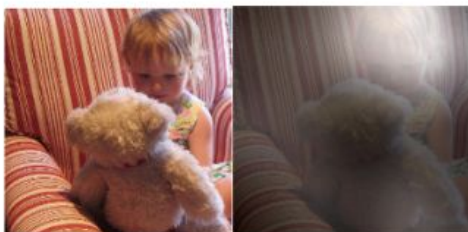# Approach V: Hierarchical Attention Network



A woman is throwing a frisbee in a park.

A dog is standing on a hardwood floor.

A stop sign is on a road with a mountain in the background.

A little girl sitting on a bed with a teddy bear.

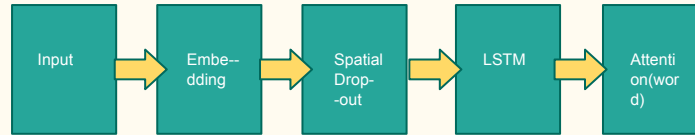A group of people sitting on a boat in the water.

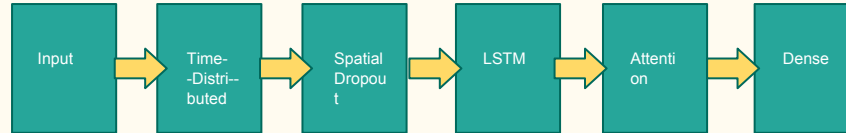A giraffe standing in a forest with trees in the background.

Figure 6: Attending to objects in an image during caption generation. The white regions indicate where the attention mechanism focused on during the generation of the underlined word. From Xu, Kelvin, et al. "Show, attend and tell: Neural image caption generation with visual attention." International Conference on Machine Learning. 2015.
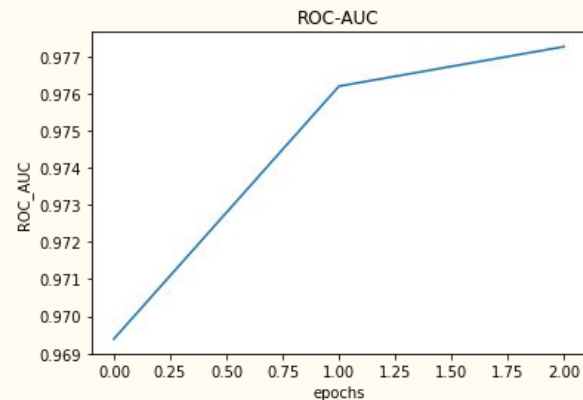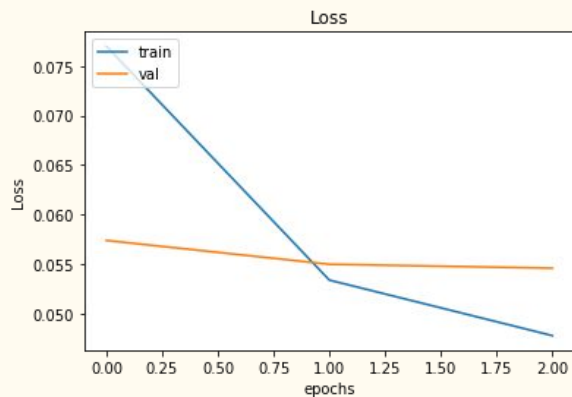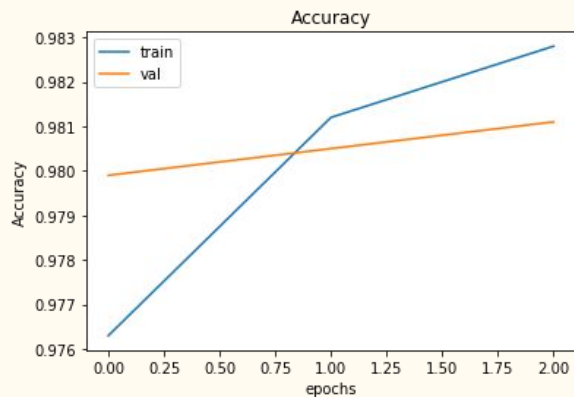
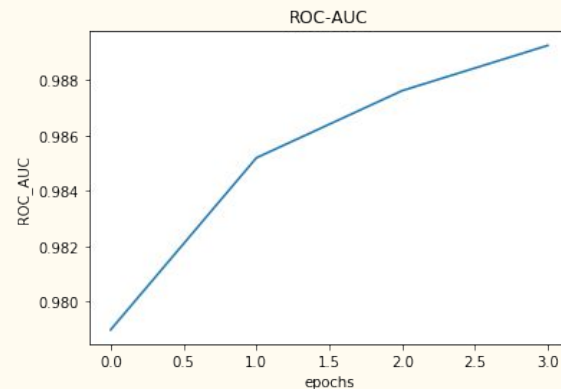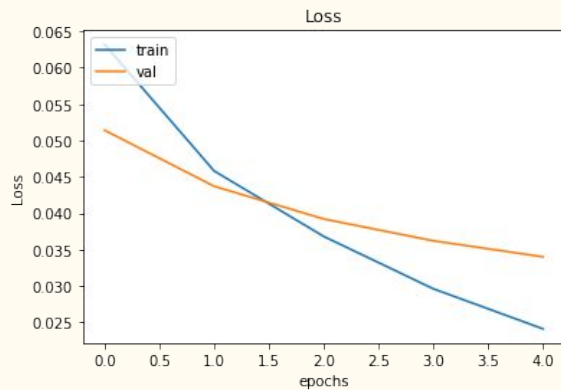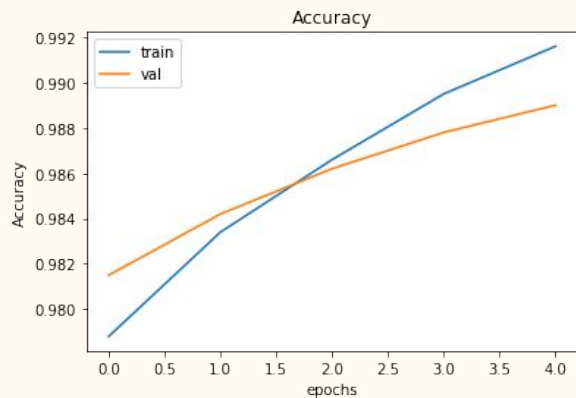# Approach V: Hierarchical Attention Network

# Approach V: Hierarchical Attention Network

Hello Again, Overfitting!

# Approach V: Data Augmentation

# Results

| Method | Tra_Loss | Val_Loss | Tra_Accuracy | Val_Accuracy |
|---|---|---|---|---|
| **Bi-LSTM+Random Embedding** | 0.0448 | 0.0489 | 97.01% | 96.89% |
| **Bi-LSTM+GloVe Embedding** | 0.0384 | 0.0455 | 98.50% | 98.34% |
| **Bi-LSTM+GloVe+Data Augmentation** | 0.0232 | 0.0298 | 99.11% | 98.91% |
| **Hierarchical Attention Network(HAN) + GloVe** | 0.0478 | 0.0546 | 98.28% | 98.11% |
| **HAN+Data Augmentation** | 0.0241 | 0.0340 | 99.16% | 98.90% |

# Acknowledgments