# Detecting Insincere Questions on Quora

*Past Kaggle competition*

**Course**: Technologies for Big Data Management and Analytics
**Fall Semester**, 2020-2021
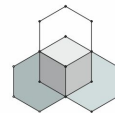**Professor**: Grigorios Tsoumakas

Theodoros Konstantinidis

George Georgiou

Panagiotis Papaemmanouil

SCHOOL OF INFORMATICS

Data & Web Science    MSc Program

# Overview

1. **Problem Statement**

2. **Dataset insights & analysis**

3. **Feature Engineering & Word Embeddings**

4. **Experiments and Results**

# Problem Statement

## Quora Insincere Questions Classification

*"In this competition, Kagglers will develop models that identify and flag insincere questions. To date, Quora has employed both machine learning and manual review to address this problem. With your help, they can develop more scalable methods to detect toxic and misleading content."*

# Dataset insights & analysis

Dataset & Target

Wordclouds

Dominant N-grams

# Dataset - Target

*"Actual questions from Quora"*

**Train dataset**

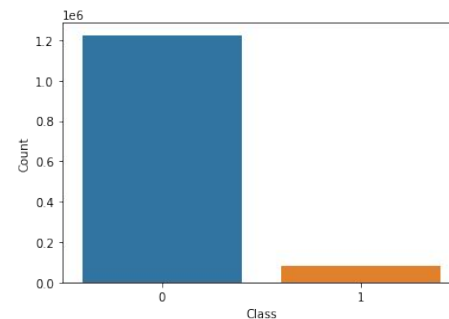1.306.122 tuples of 3 columns each: [Question identifier, question text, target]
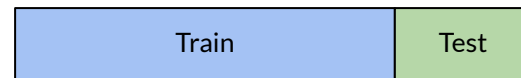
**Test dataset**

375.806 tuples, no target

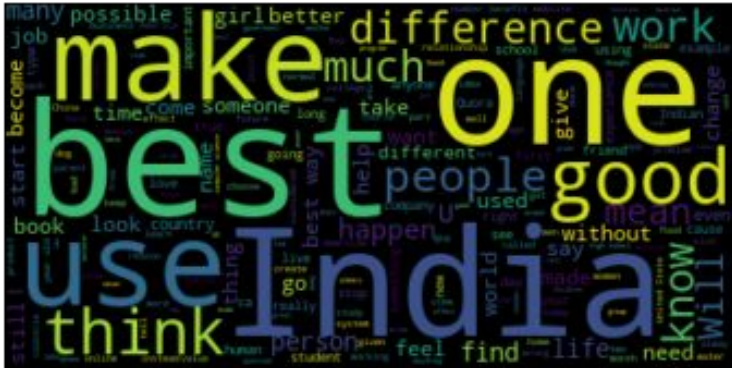No null values or empty strings

**Binary Classification**

- Target heavily imbalanced (94% sincere, 6% insincere)
- Due to imbalance, accuracy is useless: F1 Score used instead for performance evaluation



$$F1 = 2 * \frac{precision * recall}{precision + recall}$$

# Wordclouds



Sincere Question WordCloud



Insincere Question WordCloud
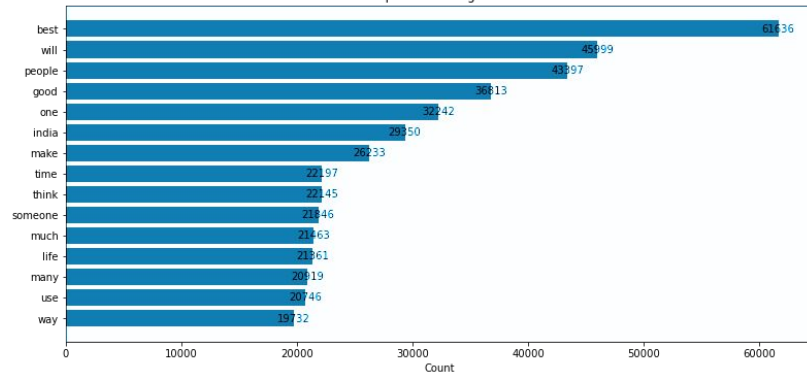
# Dominant n-grams
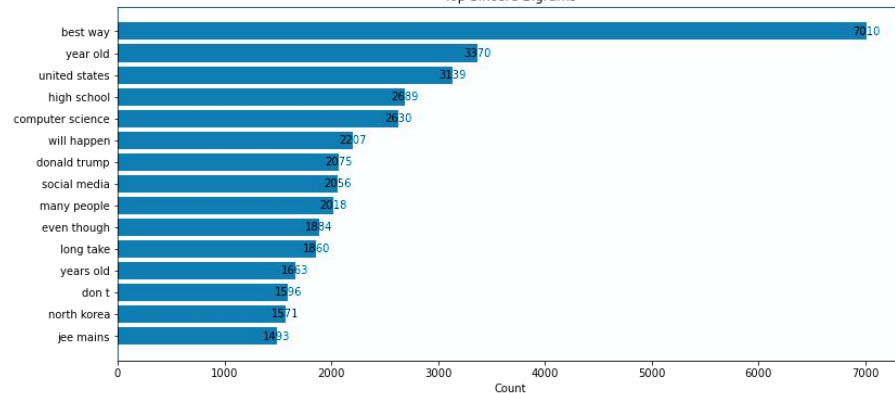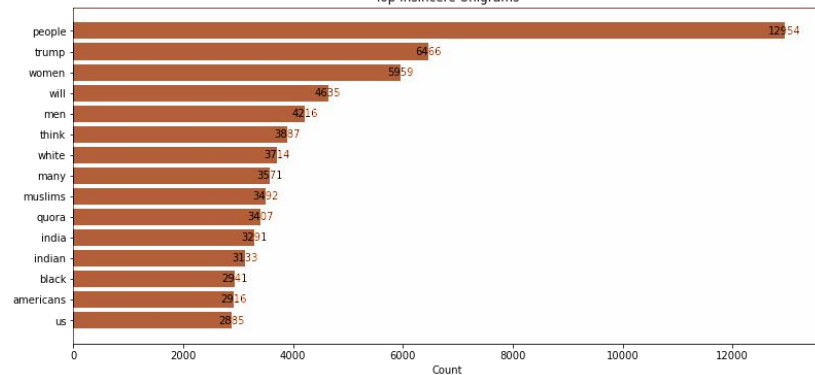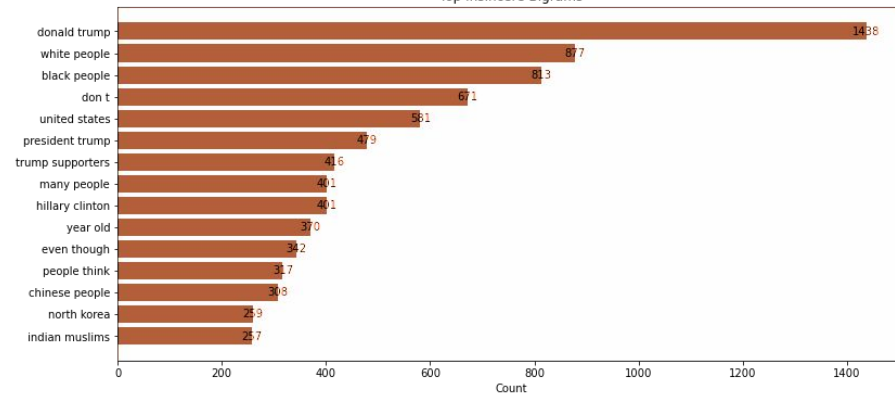


### Top Sincere Unigrams

| Word | Count |
|------|-------|
| best | 61636 |
| will | 45999 |
| people | 43397 |
| good | 36813 |
| one | 32242 |
| india | 29350 |
| make | 26233 |
| time | 22197 |
| think | 22145 |
| someone | 21846 |
| much | 21463 |
| life | 21361 |
| many | 20919 |
| use | 20746 |
| way | 19732 |

### Top Sincere Bigrams

| Bigram | Count |
|--------|-------|
| best way | 7010 |
| year old | 3370 |
| united states | 3139 |
| high school | 2689 |
| computer science | 2630 |
| will happen | 2207 |
| donald trump | 2075 |
| social media | 2056 |
| many people | 2018 |
| even though | 1884 |
| long take | 1860 |
| years old | 1663 |
| don t | 1596 |
| north korea | 1571 |
| jee mains | 1493 |

### Top Insincere Unigrams

| Word | Count |
|------|-------|
| people | 12954 |
| trump | 6466 |
| women | 5959 |
| will | 4635 |
| men | 4216 |
| think | 3837 |
| white | 3714 |
| many | 3571 |
| muslims | 3492 |
| quora | 3407 |
| india | 3291 |
| indian | 3133 |
| black | 2941 |
| americans | 2916 |
| us | 2835 |

### Top Insincere Bigrams

| Bigram | Count |
|--------|-------|
| donald trump | 1438 |
| white people | 877 |
| black people | 813 |
| don t | 671 |
| united states | 581 |
| president trump | 479 |
| trump supporters | 416 |
| many people | 401 |
| hillary clinton | 401 |
| year old | 370 |
| even though | 342 |
| people think | 317 |
| chinese people | 308 |
| north korea | 259 |
| indian muslims | 257 |

01   02   03   04

# Feature Engineering & Word Embeddings

Text preprocessing

Basic Feature Engineering

Sentiment & Readability

Spell Correction

Word Embeddings

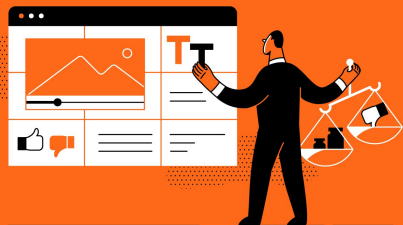# Text Preprocessing

NLTK

1. lower **Case**
2. nltk **Tokenize**
3. strip **Punctuation**
4. remove **Stopwords**
5. stem with **PorterStemmer**
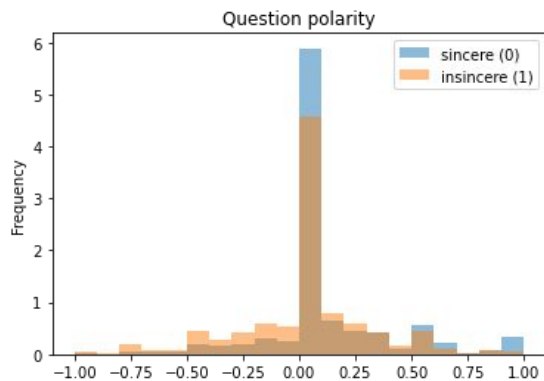6. spelling **Correction**\*

# (Basic) Feature Engineering

1. Number of words
2. Number of characters
3. Number of punctuation marks
4. Number of special characters (+ smilies)
5. Capitals proportional to length
6. Mean word length
7. Sentiment: Polarity
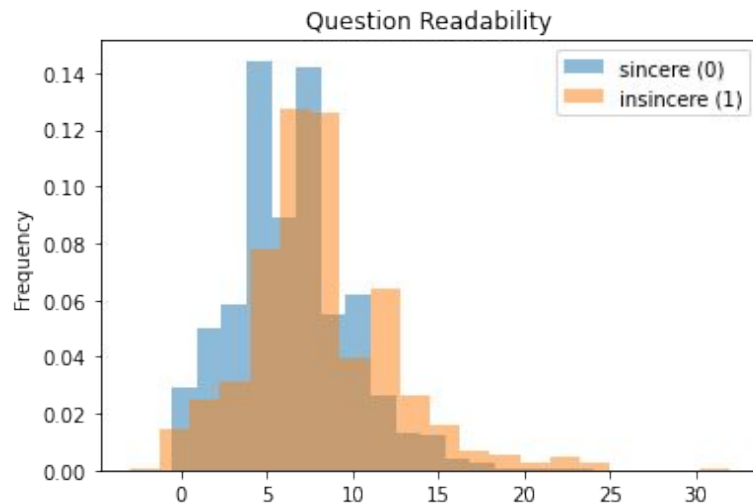8. Sentiment: Subjectivity
9. Readability

**TF-IDF:** Term Frequency - Inverse Document Frequency

# Feature Engineering: Sentiment & Readability



- Various readability indexes (i.e. SMOG)
- Values indicate required level to understand text

# What about unknown words? Spell Correction.

Use words from **fasttext** *(pre-trained word embeddings)* as vocabulary

↓

Find words outside the vocabulary

↓

Generate possible spelling corrections for the word

All edits that are **one edit away** from word

All edits that are **two edits away** from word

↓

Replace with the most probable spelling correction for word

# Word Embeddings



word2vec
Tool for computing continuous distributed representations of words.

*GoogleNews-vectors-negative300*

*wikinews_300d_1M*

**GloVe: Global Vectors for word representation**

*glove.840B.300d*

**Paragram embeddings**

*paragram_300_sl999*

# Experiments

**Available Hardware - Machine Specifications**

| | |
|---|---|
| Model name | Asus TUF Gaming A15 |
| OS | Ubuntu 18.04.5 LTS |
| Processor | AMD Ryzen 7 4800H CPU @4.20GHz (8 cores, 16 threads) |
| GPU | NVidia GeForce RTX 2060 6GB |
| Memory | 40 GB |

# Baseline Machine Learning models

- Logistic Regression

- Naive Bayes

- Random Forest

- XGBoost

**Feature pool:**

1. Handcrafted features
2. TF-IDF
3. Embeddings (*mean*)

# NN: Keras model Architecture

Input(shape=(max_length,))

Embedding(nb_words,embedding_size)

SpatialDropout1D(0.3)

x1 = Bidirectional(CuDNN LSTM(256, return_sequences=True))

x2 = Bidirectional(CuDNNGRU(128, return_sequences=True))(x1)

GlobalMaxPooling1D()(x1)

GlobalMaxPooling1D()(x2)

Concatenate()([max_pool1, max_pool2])

Dense(1)

**Hyper-parameters (I)**
- **optimizer**: *adam*
- **loss**: *binary_crossentropy*
- **metrics**: *accuracy*
- **activation**: *sigmoid*

**Hyper-parameters (II)**
- **max_length**: 55
- **embedding_size**: 600
- **learning_rate**: 0.001
- **batch_size**: 512
- **num_epoch**: 4

01   02   03   04

# Competition winning solution

*wikinews_300d_1M*

hello   $\begin{bmatrix} 1 & 10 & 1 & 2 & \ldots \\ 3 & 0.2 & 0.4 & 2 & \ldots \\ 1 & 0 & 0.2 & 2 & \ldots \end{bmatrix}$
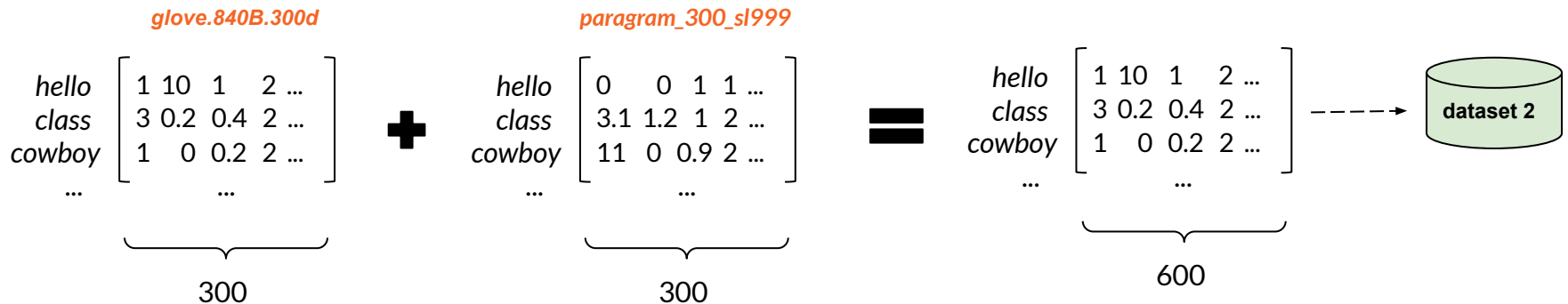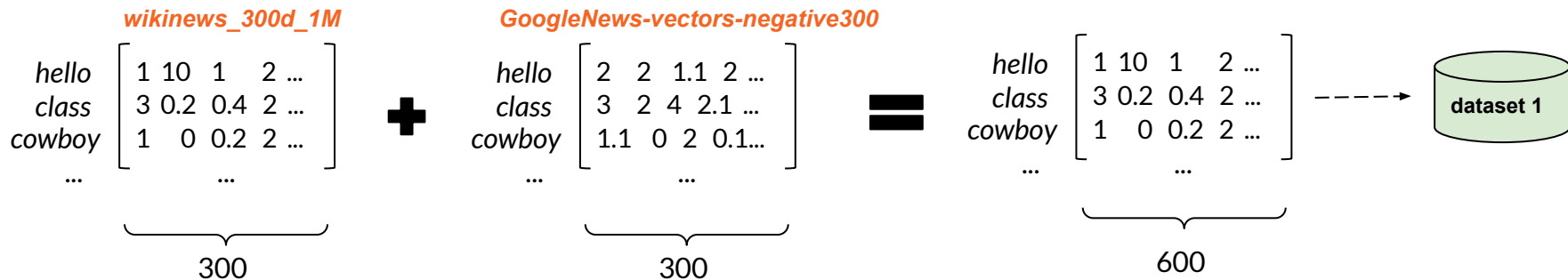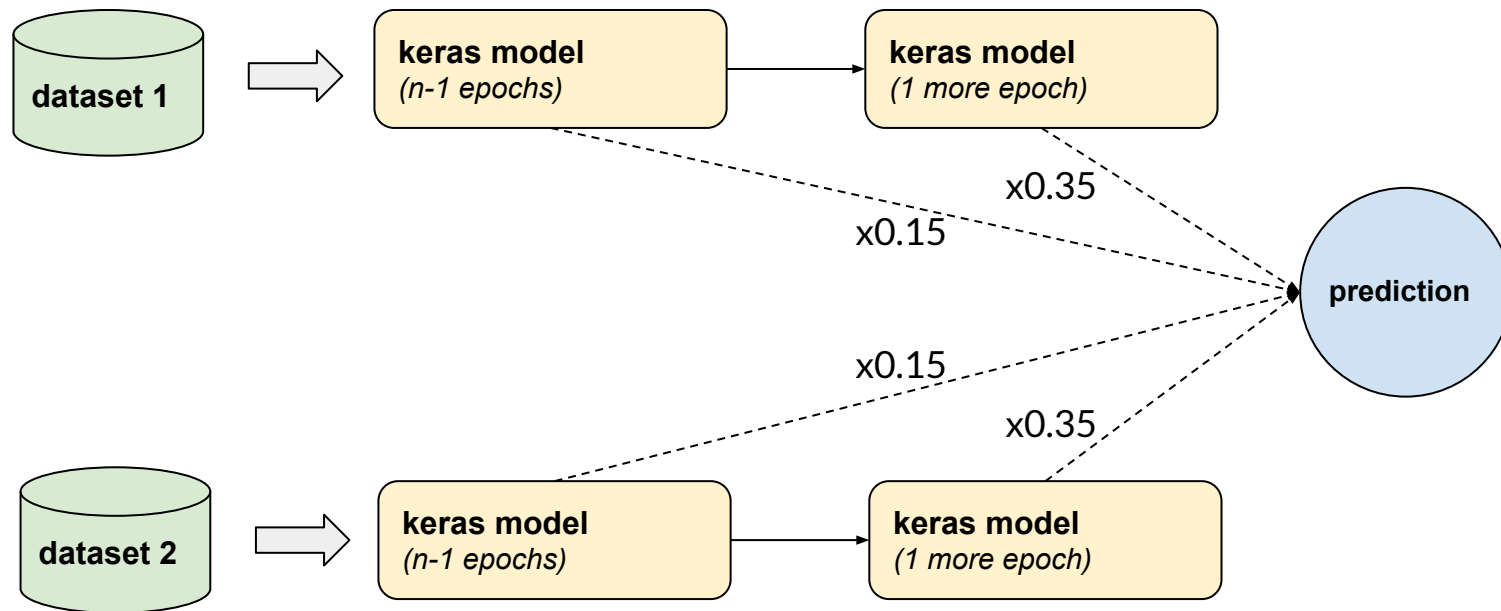class
cowboy
...

$\underbrace{\phantom{xxxxxxxxxxxx}}_{300}$

**+**

*GoogleNews-vectors-negative300*

hello   $\begin{bmatrix} 2 & 2 & 1.1 & 2 & \ldots \\ 3 & 2 & 4 & 2.1 & \ldots \\ 1.1 & 0 & 2 & 0.1 & \ldots \end{bmatrix}$
class
cowboy
...

$\underbrace{\phantom{xxxxxxxxxxxx}}_{300}$

**=**

hello   $\begin{bmatrix} 1 & 10 & 1 & 2 & \ldots \\ 3 & 0.2 & 0.4 & 2 & \ldots \\ 1 & 0 & 0.2 & 2 & \ldots \end{bmatrix}$
class
cowboy
...

$\underbrace{\phantom{xxxxxxxxxxxx}}_{600}$

- - - → **dataset 1**

*glove.840B.300d*

hello   $\begin{bmatrix} 1 & 10 & 1 & 2 & \ldots \\ 3 & 0.2 & 0.4 & 2 & \ldots \\ 1 & 0 & 0.2 & 2 & \ldots \end{bmatrix}$
class
cowboy
...

$\underbrace{\phantom{xxxxxxxxxxxx}}_{300}$

**+**

*paragram_300_sl999*

hello   $\begin{bmatrix} 0 & 0 & 1 & 1 & \ldots \\ 3.1 & 1.2 & 1 & 2 & \ldots \\ 11 & 0 & 0.9 & 2 & \ldots \end{bmatrix}$
class
cowboy
...

$\underbrace{\phantom{xxxxxxxxxxxx}}_{300}$

**=**

hello   $\begin{bmatrix} 1 & 10 & 1 & 2 & \ldots \\ 3 & 0.2 & 0.4 & 2 & \ldots \\ 1 & 0 & 0.2 & 2 & \ldots \end{bmatrix}$
class
cowboy
...

$\underbrace{\phantom{xxxxxxxxxxxx}}_{600}$

- - - → **dataset 2**

01   02   03   04

# Competition winning solution

**Typical Performance:** Evaluation Metric - f1 score

| | | |
|---|---|---|
| 1. | **Baseline ML models + Handcrafted Features** | ~ 0.20 |
| 2. | **Baseline ML models + Embeddings (mean)** | ~ 0.45 |
| 3. | **Baseline ML models + TF-IDF** | ~ 0.55 |
| 4. | **Simple keras architecture** | ~ 0.60 |
| 5. | **keras + all embeddings (Competition winning model)** | ~ 0.70 |

❏ TF-IDF superior to handcrafted features

❏ Baseline ML models doesn't work well with pre-trained word embeddings (mean across all words kills the information)

❏ The keras embedding, LSTM and GRU layers do the work!

❏ Large coverage between corpus and pre-trained word embeddings is extremely important.

# Results / Conclusions

# Спасибо большое!