# Capstone 3

Mark Perez

# Overview

- IMDB Dataset
- 50K reviews
- Use NLP for prediction

# Why is it important?

- Sentiment analysis is an important technique for natural language processing
- Models focus on polarity, emotions, or intentions
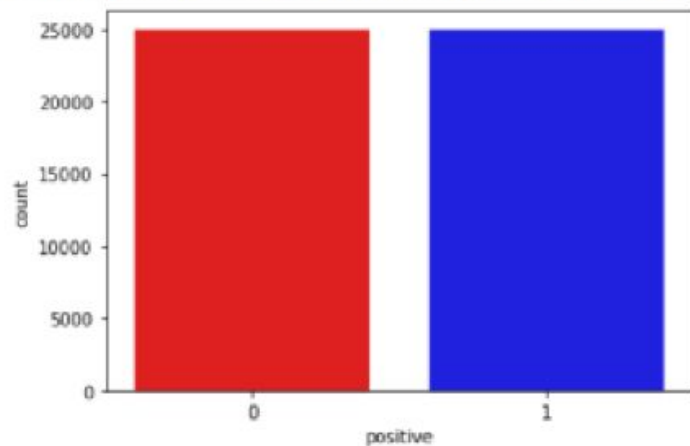- Helps businesses make decisions

# EDA

```
data.head()
```

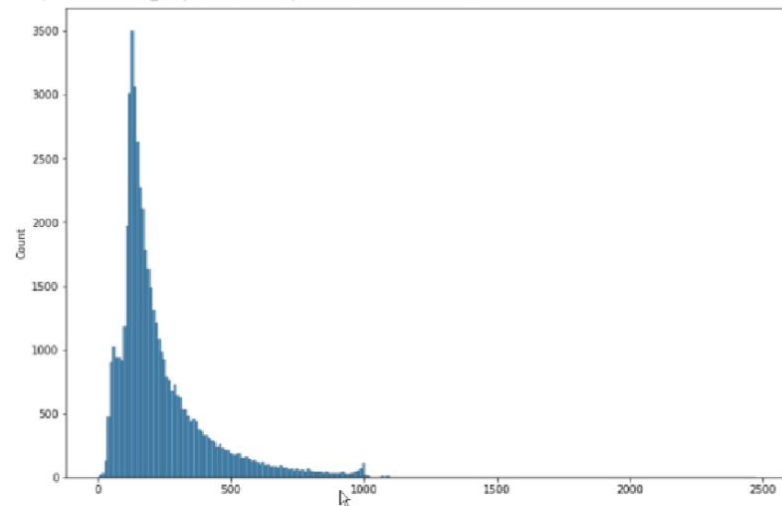|   | review | sentiment |
|---|--------|-----------|
| 0 | One of the other reviewers has mentioned that ... | positive |
| 1 | A wonderful little production. <br /><br />The... | positive |
| 2 | I thought this was a wonderful way to spend ti... | positive |
| 3 | Basically there's a family where a little boy ... | negative |
| 4 | Petter Mattei's "Love in the Time of Money" is... | positive |

```python
plt.figure(figsize = (15,12))
wc_pos = WordCloud(max_words = 1500, width = 1600 , height = 800).generate(" ".join(data[data.sentiment == 'positive'].clean_review))
plt.tight_layout(pad=0) #shrink the size of the border
plt.imshow(wc_pos, interpolation = 'bilinear')
#displayed image appear more smoothly
```

<matplotlib.image.AxesImage at 0x7f32441559d0>

```
[27] plt.figure(figsize = (15,12)) # Negative Review Text
     wc_neg = WordCloud(max_words = 2000, width = 1600 , height = 800).generate(" ".join(data[data.sentiment == 'negative'].clean_review))
     plt.imshow(wc_neg, interpolation = 'bilinear')
```
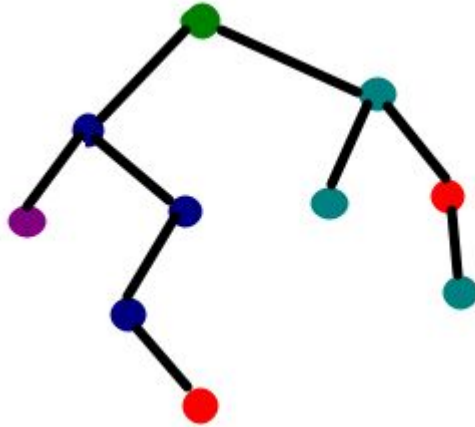
<matplotlib.image.AxesImage at 0x7f322783a050>

# Modeling

- Two models for Binary Classification:

  Decision Tree Classifier & BERT

# Decision Tree Classifier

- Simple to prepare and understand

|  | Precision | Recall | F1-score | Accuracy |
|---|---|---|---|---|
| 0 | .82 | .6 | .69 | |
| 1 | .47 | .73 | .57 | |
|  |  |  |  | .64 |

# What is BERT?

- Bidirectional Encoder Representations from Transformers
- Each word is contextualized based on the other words in the sentence.
- Different versions of BERT

Text Summarization

Text Encoding Similarity Retrieval

# How BERT works

- BERT relies on a Transformer
- The input for the encoder are three embedings: Token, Segment, and Positional
- BERT's language modeling task (MLM) masks 15% of words in the input and asks the model to predict the missing word, as well as predict whether a following sentence (after a SEP token) is random or not.

# Visual representation of BERT tokens

# Beauty of Tensorflow

- TensorFlow Hub is a repository for trained machine learning models
- We load the preprocessing model (to prepare the text) and the model (small BERT) from TF Hub.

```
[ ] bert_model_name = 'small_bert/bert_en_uncased_L-4_H-512_A-8'

    tfhub_handle_encoder = 'https://tfhub.dev/tensorflow/small_bert/bert_en_uncased_L-4_H-512_A-8/1'
    tfhub_handle_preprocess = 'https://tfhub.dev/tensorflow/bert_en_uncased_preprocess/3'

    print(f'BERT model selected           : {tfhub_handle_encoder}')
    print(f'Preprocess model auto-selected: {tfhub_handle_preprocess}')

    BERT model selected           : https://tfhub.dev/tensorflow/small_bert/bert_en_uncased_L-4_H-512_A-8/1
    Preprocess model auto-selected: https://tfhub.dev/tensorflow/bert_en_uncased_preprocess/3
```
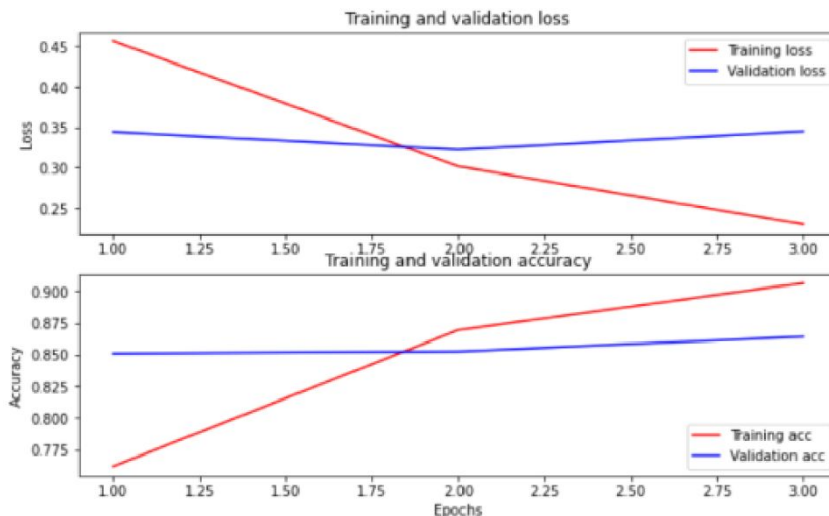
# Results



```
[40] history = classifier_model.fit(X_train, y_train,
                                     validation_data=(X_val, y_val),
                                     epochs=epochs)

     Epoch 1/3
     438/438 [==============================] - 122s 259ms/step - loss: 0.4804 - binary_accuracy: 0.7571 - val_loss: 0.4279 - val_binary_accuracy: 0.8048
     Epoch 2/3
     438/438 [==============================] - 112s 255ms/step - loss: 0.3137 - binary_accuracy: 0.8639 - val_loss: 0.3274 - val_binary_accuracy: 0.8577
     Epoch 3/3
     438/438 [==============================] - 112s 256ms/step - loss: 0.2291 - binary_accuracy: 0.9072 - val_loss: 0.3566 - val_binary_accuracy: 0.8660


[41] loss, accuracy = classifier_model.evaluate(X_test, y_test)

     print(f'Loss: {loss}')
     print(f'Accuracy: {accuracy}')

     157/157 [==============================] - 16s 104ms/step - loss: 0.3527 - binary_accuracy: 0.8590
     Loss: 0.3526787757873535
     Accuracy: 0.859000027179718
```

# Conclusion

- BERT Model is great at Sentiment analysis
- Accuracy increased by about 22%
- Accuracy is not always the best metric to measure

# Further Analysis and Constraints

- More data
- Compare other models (Random Forest and LSTM)
- Fine tuning