# Cancer tweet sentiment analysis

Maggie and Ricardo
May 2nd, 2023

Big Data Project: Twitter Dashboard

# Content

- Objectives
- Data
- Challenges
- Sentiment predictions
- Relevance of findings

# Objectives

- Apply Spark ML for a sentiment classification model
- Use Athena to create tables and Quicksight to build dashboards
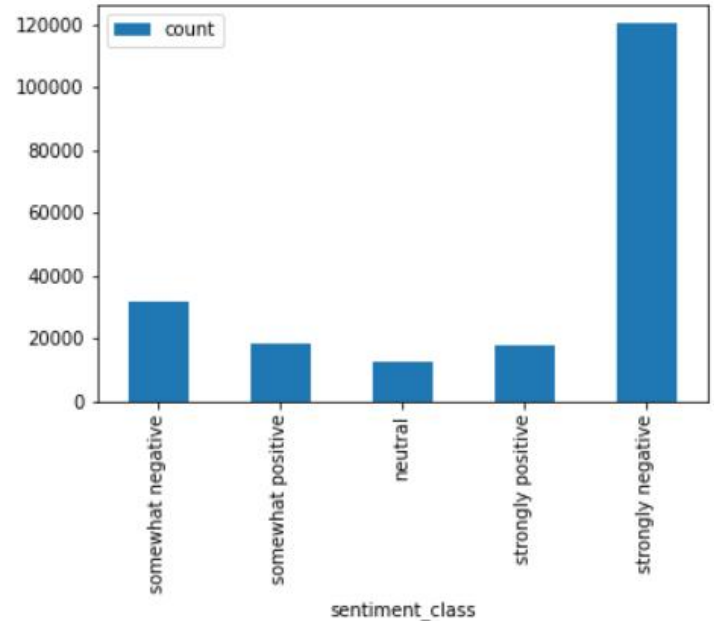
# Data

Columns: id, user_name, screen_name, tweet, followers, geo, coordinates, time

Rows: 201401 initial, 201067 after remove nulls

Processing:

Clean, tokenization, remove stopwords

count vectorization and TF-IDF

# Problems and challenges

- Understanding Py-spark and how it works with AWS
- Free version is too slow takes a long time to predict models.
- Understanding datafiles as parquet and the advantages of each one.
- Slight difference in SQL queries with athena.
- Quicksight is less intuitive than Tableau and more limited in their visuals.

# Sentiment Predictions

Logistic Regression: 1 Gram, Text Blob and binary classification.

```python
from pyspark.ml.feature import NGram, VectorAssembler, StopWordsRemover, HashingTF, IDF, Tokenizer, StringIndexer, CountVectorizer
from pyspark.ml.classification import LogisticRegression
from pyspark.ml import Pipeline
from pyspark.ml.evaluation import BinaryClassificationEvaluator
from pyspark.ml import PipelineModel

# Use 70% cases for training, 30% cases for testing
train, test = df_for_model.randomSplit([0.7, 0.3], seed=42)

# Create transformers for the ML pipeline
tokenizer = Tokenizer(inputCol="tweet_text", outputCol="tokens")
stopword_remover = StopWordsRemover(inputCol="tokens", outputCol="filtered")
cv = CountVectorizer(vocabSize=2**16, inputCol="filtered", outputCol='cv')
idf = IDF(inputCol='cv', outputCol="1gram_idf", minDocFreq=5) #minDocFreq: remove sparse terms
assembler = VectorAssembler(inputCols=["1gram_idf"], outputCol="features")
# assembler convert several columns to one call features so it can be fed to the model
label_encoder = StringIndexer(inputCol = "sentiment_label", outputCol = "label")
# we always need a column call features and one call label, if not you need to go inside the model and change the default names
lr = LogisticRegression(maxIter=100)
lr_pipeline = Pipeline(stages=[tokenizer, stopword_remover, cv, idf, assembler,label_encoder,lr])

lr_pipeline_model = lr_pipeline.fit(train)
lr_predictions = lr_pipeline_model.transform(test)

lr_evaluator = BinaryClassificationEvaluator(rawPredictionCol="rawPrediction")
lr_accuracy = lr_predictions.filter(lr_predictions.label == lr_predictions.prediction).count() / float(test.count())
lr_roc_auc = lr_evaluator.evaluate(predictions)

print("Accuracy Score: {0:.4f}".format(accuracy))
print("ROC-AUC: {0:.4f}".format(roc_auc))

lr_pipeline_model.save('dbfs:/mnt/s3_bucket/Final_models/LR')
```

Accuracy: 0.86
ROC-AUC: 0.82

Logistic Regression: 1 & 2 Gram, Vader and multiclass classification.

```python
# LogisticRegression Model

from pyspark.ml.feature import NGram, VectorAssembler, StopWordsRemover, HashingTF, IDF, Tokenizer, StringIndexer, NGram, ChiSqSelector, VectorAssembler
from pyspark.ml import Pipeline

# label encoder
from pyspark.ml.feature import StringIndexer

# label
label_encoder= StringIndexer(inputCol = "sentiment_class", outputCol = "label")

# Create transformers for the ML pipeline
tokenizer = Tokenizer(inputCol="tweet", outputCol="tokens")
stopword_remover = StopWordsRemover(inputCol="tokens", outputCol="filtered")
cv = CountVectorizer(vocabSize=2**16, inputCol="filtered", outputCol='cv')
idf = IDF(inputCol='cv', outputCol="1gram_idf", minDocFreq=5) #minDocFreq: remove sparse terms
ngram = NGram(n=2, inputCol="filtered", outputCol="2gram")
ngram_hashingtf = HashingTF(inputCol="2gram", outputCol="2gram_tf", numFeatures=20000)
ngram_idf = IDF(inputCol="2gram_tf", outputCol="2gram_tf", minDocFreq=5)

# assemble multiple input columns into a vector column, and then perform feature selection on the resulting vector column using the chi-squared test
# Assemble all text features
assembler = VectorAssembler(inputCols=["1gram_idf", "2gram_tf"], outputCol="rawFeatures")

# Chi-square variable selection
selector = ChiSqSelector(numTopFeatures=2**14,featuresCol='rawFeatures', outputCol="features")

# Regression model estimator
lr = LogisticRegression(maxIter=100)

# Build the pipeline
pipeline = Pipeline(stages=[label_encoder, tokenizer, stopword_remover, cv, idf, ngram, ngram_hashingtf, ngram_idf, assembler, selector, lr])

# Pipeline model fitting
pipeline_model = pipeline.fit(trainDF)
predictions = pipeline_model.transform(testDF)

evaluator = MulticlassClassificationEvaluator(predictionCol="prediction")
accuracy = predictions.filter(predictions.label == predictions.prediction).count() / float(testDF.count())
roc_auc = evaluator.evaluate(predictions)

print("Accuracy Score: {0:.4f}".format(accuracy))
print("ROC-AUC: {0:.4f}".format(roc_auc))
```
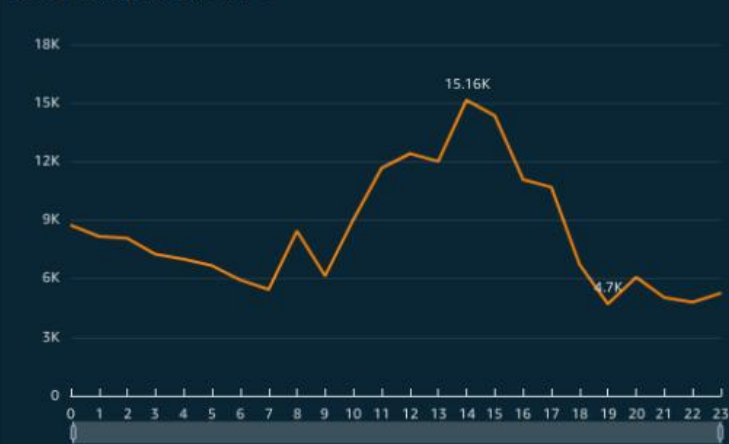
Accuracy: 0.8724
ROC-AUC: 0.8728

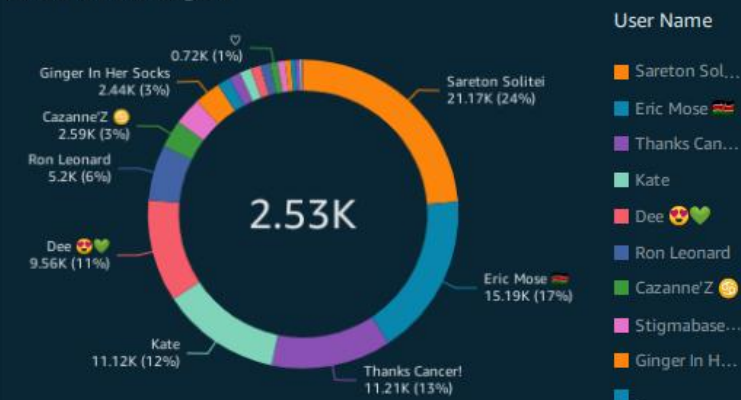# Raw data

## Tweet Count per hour GMT+0



Peak value labeled 15.16K (around hour 14), with a point labeled 4.7K (around hour 19). Y-axis: 0, 3K, 6K, 9K, 12K, 15K, 18K. X-axis hours: 0–23.

## WordCloud before cleaning data

SHOWING TOP 100 IN TWEET_TEXT



Prominent words: "Talk sensibly. cancer that's always ther..."
RT @fesshole: An old guy at my golf club...
RT @RepMTG: I witnessed in person human...
RT @blurayangel: Chadwick Boseman doing...
RT @POTUS: Too many men and women in uni...
RT @The_Law_Boy: American Lawyer gives ...

## Average followets by user

SHOWING TOP 20 IN USER_NAME



Center value: 2.53K

Segments:
- 0.72K (1%)
- Ginger In Her Socks 2.44K (3%)
- Cazanne'Z 2.59K (3%)
- Ron Leonard 5.2K (6%)
- Dee 9.56K (11%)
- Kate 11.12K (12%)
- Thanks Cancer! 11.21K (13%)
- Eric Mose 15.19K (17%)
- Sareton Solitei 21.17K (24%)

User Name legend:
- Sareton Sol...
- Eric Mose
- Thanks Can...
- Kate
- Dee 😍💚
- Ron Leonard
- Cazanne'Z 🌝
- Stigmabase...
- Ginger In H...
- .

## Tweets created from Nov 22 to Nov 24 2022 by user



- Md Alim Quemar — 489
- liam smith — 406
- Bobby6740 — 360
- . — 265
- Nicholas Tompanis — 179
- Thanks Cancer! — 96
- Lady Nique — 84
- Dee 😍💚 — 78
- Mike — 72
- Elizabeth Hampton... — 71
- J — 64

X-axis: 0, 100, 200, 300, 400, 500

May 2, 2023 1:52 AM (GMT)

# Prediction with Text Blob & 1 gram

Prediction with VADER 1 & 2 grams.

# Conclusion

- The sentiment labels from TextBlob and VADER are very different, and neither of them provides a fully fitted sentiment analysis.
- Logistic Regression was the best model in both cases and performed better than Random Forest.
- Experiment with a dataset that has the manually processed labels
- Takes emoticons into consideration
- Further analysis of N-grams and difference in multiclass and binary labels.