

POKEMON GO ON TWITTER ANALYSIS

WHY POKEMON GO?



- Anit Gupta
- Anil Ravuru
- Shivani Bhoite
- Alejandro Benitez
- Umaraj Potla

GEO-LOCATION ANALYSIS

LDA

POKEMON GO FACTS
&
PREDICTIONS

SENTIMENT ANALYSIS



WHY POKEMON GO?

DATA PREPROCESSING
AND
PERFORMANCE

POKEMON GO
Vs
OTHER APPLICATIONS

WHY POKEMON GO?

Major events in 2016

- The U.S. Presidential Election
- The Zika Virus
- The continued threat of ISIS
- The 2016 Summer Olympics in Rio de Janeiro, Brazil

DATA PREPROCESSING
AND
PERFORMANCE

POKEMON GO
VS
OTHER APPLICATIONS

DATA PREPROCESSING AND PERFORMANCE

POKEMON GO
Vs
OTHER APPLICATIONS

WHY POKEMON GO?

Major events in 2016

- The U.S. Presidential Election
- The Zika Virus
- The continued threat of ISIS
- The 2016 Summer Olympics in Rio de Janeiro, Brazil

At the same time Nintendo launched a gaming app which revolutionized mobile gaming platform.

DATA PREPROCESSING AND PERFORMANCE

POKEMON GO
VS
OTHER APPLICATIONS

WHY POKEMON GO?

Major events in 2016

- The U.S. Presidential Election
- The Zika Virus
- The continued threat of ISIS
- The 2016 Summer Olympics in Rio de Janeiro, Brazil

At the same time Nintendo launched a gaming app which revolutionized mobile gaming platform.

Pokemon Go!

WHY POKEMON GO?

Major events in 2016

- The U.S. Presidential Election
- The Zika Virus
- The continued threat of ISIS
- The 2016 Summer Olympics in Rio de Janeiro, Brazil

At the same time Nintendo launched a gaming app which revolutionized mobile gaming platform.

Pokemon Go!

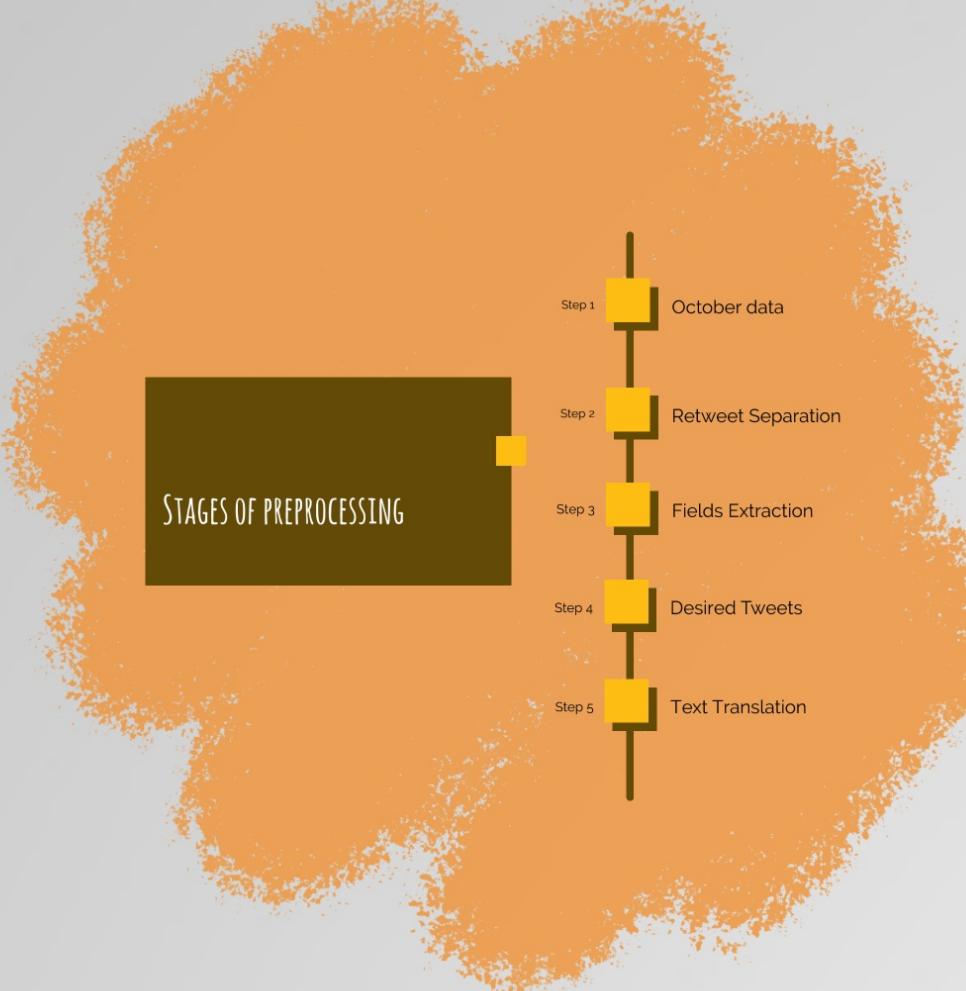
DATA PREPROCESSING
AND
PERFORMANCE

POKEMON GO
VS
OTHER APPLICATIONS

POKEMON GO
VS
OTHER GAMING APPS

CLUSTERING
ENVIRONMENT

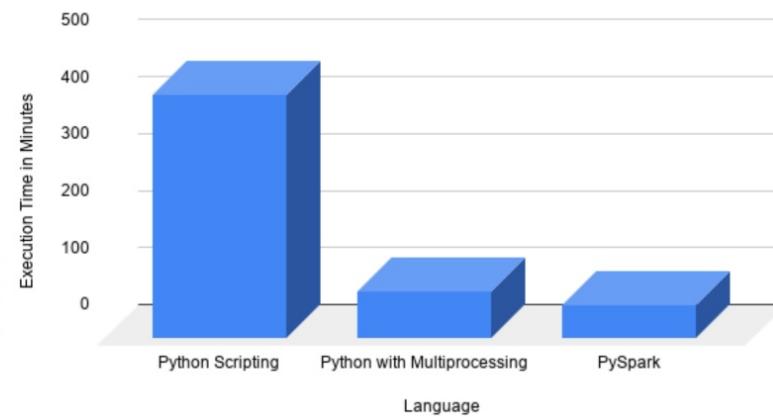
STAGES OF PREPROCESSING

- 
- Step 1 October data
 - Step 2 Retweet Separation
 - Step 3 Fields Extraction
 - Step 4 Desired Tweets
 - Step 5 Text Translation

PERFORMANCES

EXECUTION TIMES

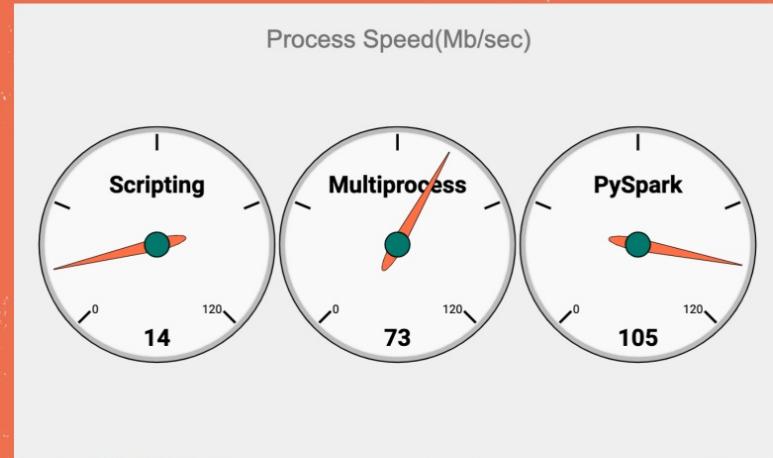
Execution Time in Minutes vs. Language



DESCRIPTIONS

SPEEDS

SPEED COMPARISONS





SYSTEM SPECIFICATIONS

We used 8 core, 16GB RAM Linux machine for running the python script.

For Python Multiprocessing we used 8 core, 16GB RAM Linux Machine.

For PySpark, We used 3 worker nodes with 48 Cores in Azure HDInsights.

AZTK - AZURE TOOLKIT

```
-bash: ztk: command not found
(bdp) (base) Anits-MacBook-Pro:bin agupta$ aztk spark cluster get --id mycluster
Cluster      mycluster
-----
State:      steady
Node Size:  standard_d12_v2
Created:    2019-11-08 04:16:47
Nodes:      5
| Dedicated: 0
| Low priority: 5

|       Nodes       |     State     |      IP:Port      | Dedicated |   Master   |
|tvmgs_09ea51d69f8f049d29b636d44ff07dc168bf297d33efc80c987fa414e61aa88_p| waitingforstarttask | 52.188.221.144:50004 |
|tvmgs_35f5bd615f39dd92df84c188bc7bd2a3c6ad2d06696a4275d2c3411fbdbab3d0b_p| waitingforstarttask | 52.188.221.144:50001 |
|tvmgs_b41c5f09cb78de645aa7961285507ed1fec7f37793e7de972030350ba45d8b78_p| waitingforstarttask | 52.188.221.144:50002 |
|tvmgs_ba30c2b02c0e8bc902b3ba0d379cd08f4add8fc7d328b08bd3fef70db86be8c2_p| waitingforstarttask | 52.188.221.144:50000 |
|tvmgs_bd53a926034cc5e6fd60376260d0471ea6c36e16b10bc5926ae0f8b24bb5b64d_p| waitingforstarttask | 52.188.221.144:50003 |
```

AZURE-HDINSIGHT

Apache Spark 2.3.0 **Spark Master at spark://10.0.0.5:7077**

URL: spark://10.0.0.5:7077
REST URL: spark://10.0.0.5:6066 (cluster mode)

Alive Workers: 5

Cores in use: 20 Total: 0 Used
Memory in use: 132.3 GB Total, 0.0 B Used
Applications: 0 Running, 0 Completed
Drivers: 0 Running, 0 Completed
Status: ALIVE

Workers (5)

Worker Id	Address	State	Cores	Memory
worker-20191108051128-10.0.0.5-45225	10.0.0.5:45225	ALIVE	4 (0 Used)	26.5 GB (0.0 B Used)
worker-20191108051138-10.0.0.6-46661	10.0.0.6:46661	ALIVE	4 (0 Used)	26.5 GB (0.0 B Used)
worker-20191108051148-10.0.0.4-35951	10.0.0.4:35951	ALIVE	4 (0 Used)	26.5 GB (0.0 B Used)
worker-20191108051152-10.0.0.7-44897	10.0.0.7:44897	ALIVE	4 (0 Used)	26.5 GB (0.0 B Used)
worker-20191108051223-10.0.0.8-35165	10.0.0.8:35165	ALIVE	4 (0 Used)	26.5 GB (0.0 B Used)

Running Applications (0)

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration

Completed Applications (0)

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration

HDINSIGHT

AZURE STORAGE EXPLORER

Why HDInsight?



Easy

Quickly spin up open source projects and clusters, with no hardware to install or infrastructure to manage.



Cost-effective

Reduce costs by creating big data clusters on demand, easily scaling them up or down, and paying only for what you use.



Enterprise-grade

Get enterprise-grade security and industry-leading compliance, with more than 30 certifications.



Open

Create optimized components for Hadoop, Spark, and more. Keep up to date with the latest versions.

The screenshot shows the Azure Storage Explorer interface. On the left, the 'EXPLORER' sidebar lists storage accounts, blob containers, and tables. The main area displays a list of blobs under the 'new' container, including 'anil_data', 'pokemonCharsCount', 'Scripts', 'twitter_Data', 'PokemonCharacters.csv', 'pokemonCharsCount', 'SparkAllRetweets.csv', 'SparkAllRetweets1.csv', 'SparkAllRetweets', and 'SparkAllRetweets1'. The 'Activities' section shows several completed transfers from local paths like '/Users/agupta/Documents/Semester 1/BDP/Project/Final/final data/bdp_all.py' to the 'new' container.

```

def write_to_AzureStorage(df,file_name):
    from azure.storage.blob import BlockBlobService
    import pandas as pd
    import io

    output = io.StringIO()
    output = df.to_csv(encoding = "utf-8")

    accountName = "bdpproj1"
    accountKey = "*****"
    containerName = "new"
    blobName = "ParallelRetweets.csv"

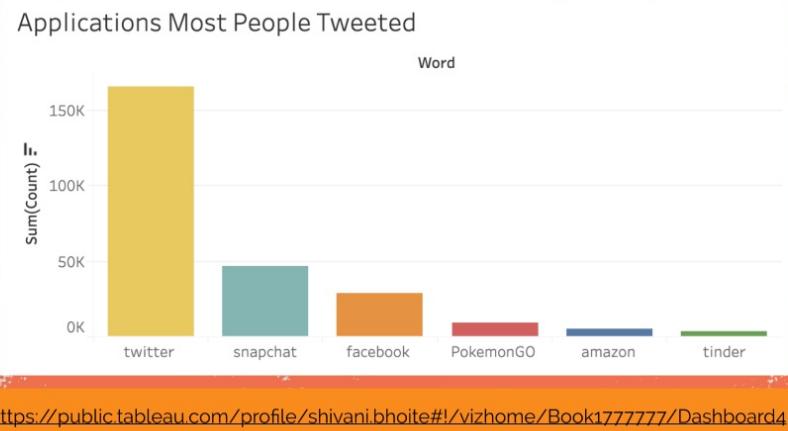
    blobService = BlockBlobService(account_name=accountName, account_key=accountKey)

    blobService.create_blob_from_text('new/anit',file_name, output)

def is_pokemon_tweet(text):

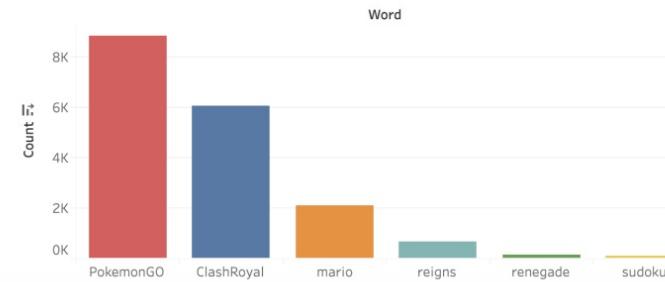
```

POKEMON GO VS OTHER APPLICATIONS



POKEMON GO VS GAMING APPLICATIONS

Mobile Games Most People Tweeted



<https://public.tableau.com/profile/shivani.bhoite#!/vizhome/Book1777777/Dashboard4>

POKEMON GO ON TWITTER ANALYSIS

WHY POKEMON GO?



- Anit Gupta
- Anil Ravuru
- Shivani Bhoite
- Alejandro Benitez
- Umaraj Potla

GEO-LOCATION ANALYSIS

LDA

POKEMON GO FACTS
&
PREDICTIONS

SENTIMENT ANALYSIS

GEO - LOCATION ANALYSIS

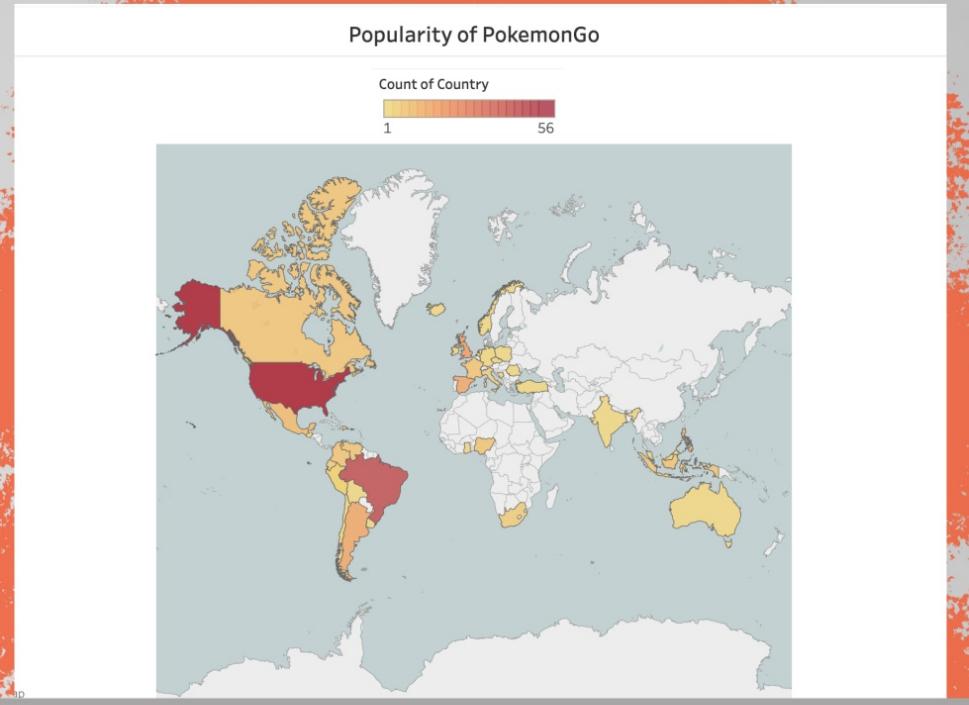
We wanted to analyze the places where people have been using the pokemongo App.

COUNTRIES PLAYING
POKEMON GO

COUNTRIES TWEETING
MOST
ABOUT POKEMONGO



<https://public.tableau.com/profile/shivani.bhoite#!/vizhome/Book1777777/3-CountriesPlayingPokemonGo>



<https://public.tableau.com/profile/shivani.bhoite#!/vizhome/Book1777777/3-PopularityofPokemonGo>

POKEMON GO ON TWITTER ANALYSIS

WHY POKEMON GO?



- Anit Gupta
- Anil Ravuru
- Shivani Bhoite
- Alejandro Benitez
- Umaraj Potla

GEO-LOCATION ANALYSIS

LDA

POKEMON GO FACTS
&
PREDICTIONS

SENTIMENT ANALYSIS

SENTIMENT ANALYSIS

We analyzed how people are reacting towards PokemonGo!!

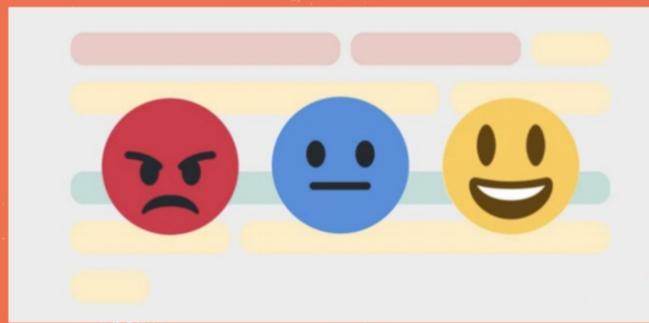


WHAT'S SENTIMENT
ANALYSIS?

POKEMON GO SENTIMENTS

SENTIMENT ANALYSIS DEFINITION

process of determining if a piece of writing is positive / negative / neutral



BUT WHY SENTIMENT ANALYSIS?

WHY?? SENTIMENT ANALYSIS??

- To understand the scope of PokemonGo better
- Twitter popular microblogging site
- Variety of twitter audiences – common man to celebrities



WHAT HAVE WE USED ?

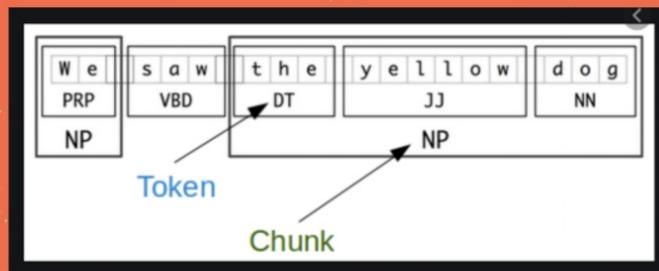
WHAT IS IT?

WHAT HAVE WE USED ?

- Natural Language Toolkit (NLTK) - open source library of tools for natural language processing
- VADER (Valence Aware Dictionary and sEntiment Reasoner)
- Chunking

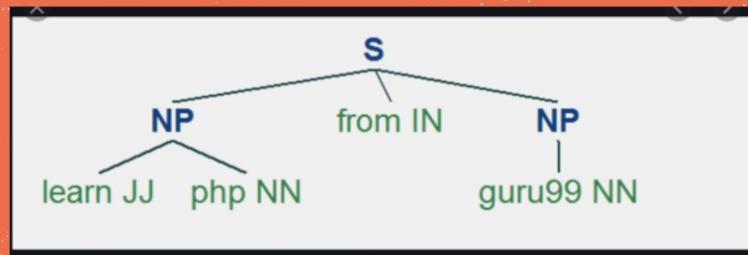
WHAT IS CHUNKING ?

Chunking breaks a text up into user-defined units ('chunks') that contain certain types of words (nouns, adjectives, verbs) or phrases (noun phrases, verb phrases, prepositional phrases).



GRAMMER

GRAMMAR



Reference : <https://www.guru99.com/pos-tagging-chunking-nltk.html>

WHY?

WHY??

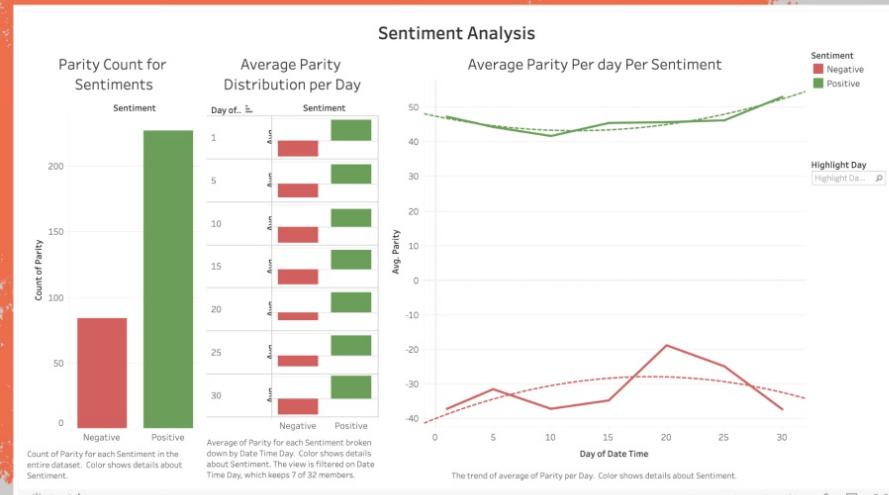
Why VADER ?

- It works exceedingly well on social media type text
- It doesn't require any training data but is constructed from a generalizable, valence-based, human-curated gold standard sentiment lexicon
- It is fast enough

Why Chunking ?

- ability to analyze the text and tag each word with its part of speech.
- allows us to pull out groups of words with set characteristics rather than selecting text by frequency

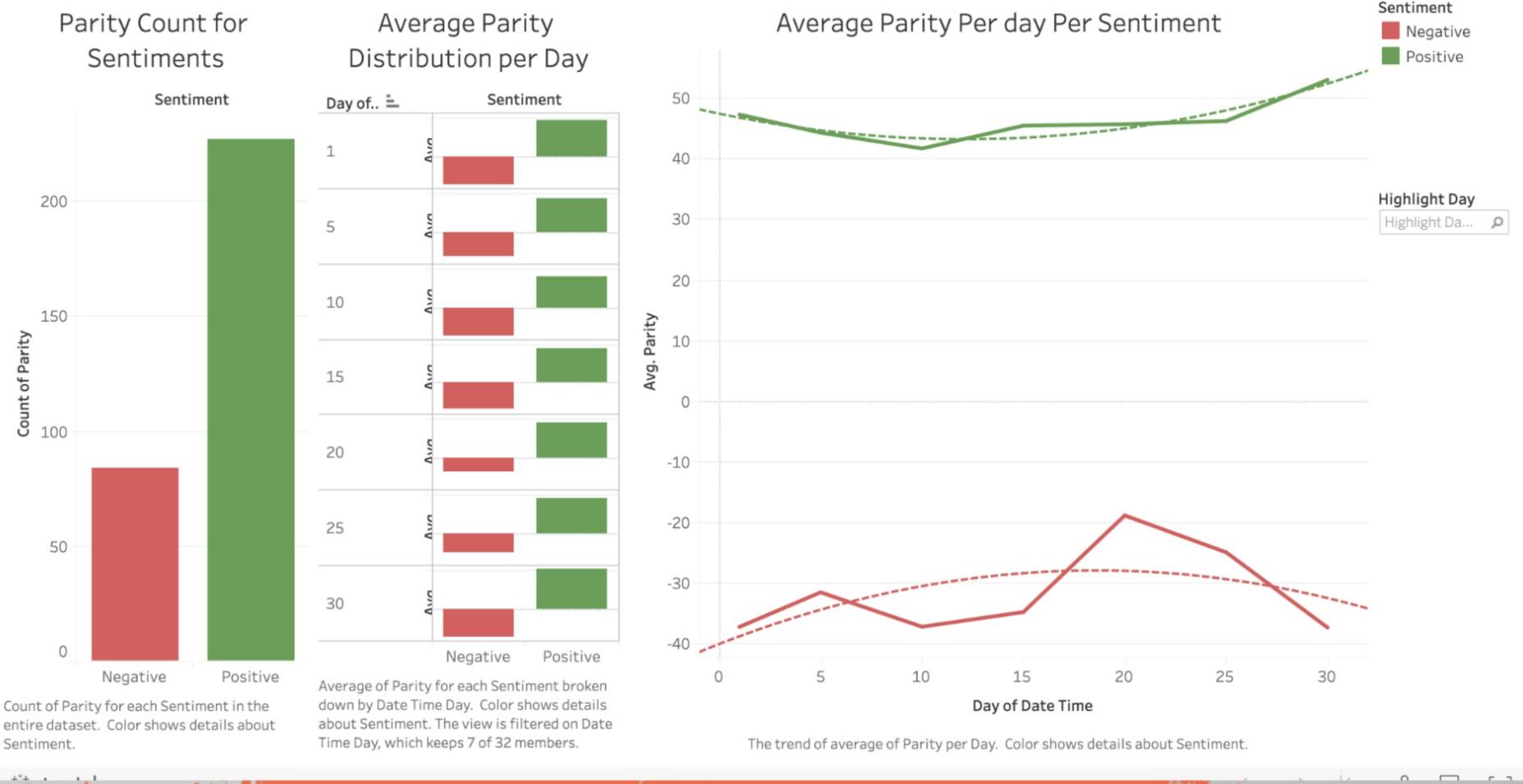
POKEMON TWEETS SENTIMENTS



Lets head to Tableau as its better viewable there

<https://public.tableau.com/profile/shivani.bhoite#!/vizhome/Book1777777/Dashboard1>

Sentiment Analysis



POKEMON GO ON TWITTER ANALYSIS

WHY POKEMON GO?



- Anit Gupta
- Anil Ravuru
- Shivani Bhoite
- Alejandro Benitez
- Umaraj Potla

GEO-LOCATION ANALYSIS

LDA

POKEMON GO FACTS
&
PREDICTIONS

SENTIMENT ANALYSIS

LATENT DIRICHLET ALLOCATION(LDA)

- Popular algorithm for topic modeling
- Topic modeling is a method for unsupervised clustering of a collection of documents to find natural groupings among them
- Goal for LDA: Better understand what is being said about PokemonGo without shifting through individuals tweets

TOPICS

FEATURE EXTRACTION

RESULTS

APPENDIX

LDA TOPICS

- LDA model treats each document as a mixture of topics
- Each topic is a mixture of words
- In this way, topics are allowed to overlap instead of discrete groups
- Topics are represented by keywords frequent in each topic
- Disadvantage is topics may be difficult interpret into coherent simple ideas



Documents

Topic proportions and assignments

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

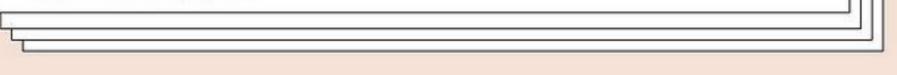
Although the numbers don't match precisely, those predictions

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson, at Umeå University in Sweden. "We arrived at the 800 number, but coming up with a consensus answer may be more than just a genetic numbers game." Particularly, more and more genomes are being sequenced and analyzed. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing all

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

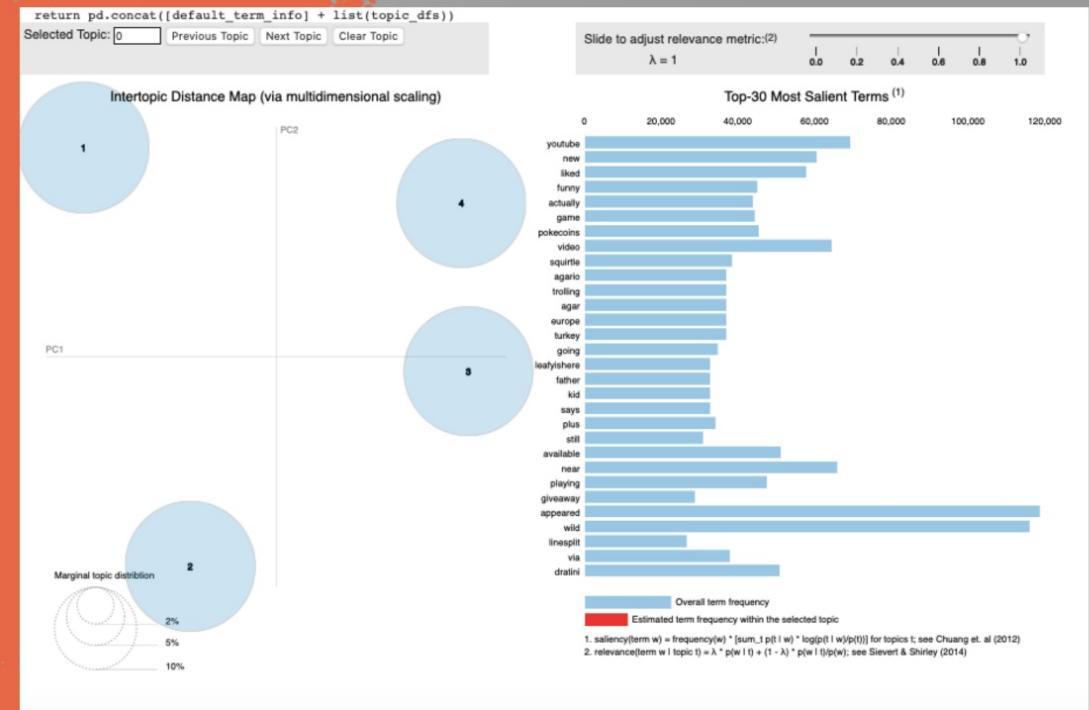
SCIENCE • VOL. 272 • 24 MAY 1996



TF-IDF

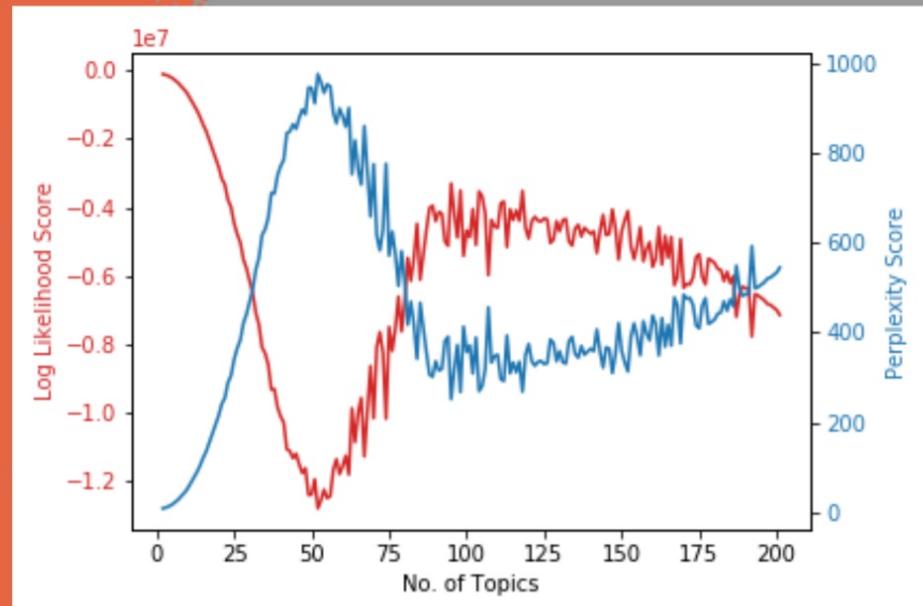
- Term Frequency-Inverse Document Frequency
- Numeric feature reflecting how important a word is in a document and penalizing it for how common the word is in many documents

Go to notebook



CHOOSING # OF TOPICS

- Use perplexity and log-likelihood metrics to evaluate separation of topics on subset of data
- Optimal topics is 2
- Increased Topics to 4 to break down topics further into more understandable topic-term distribution.
- More than 5 topics began to significantly overlap in topics



POKEMON GO ON TWITTER ANALYSIS

WHY POKEMON GO?



- Anit Gupta
- Anil Ravuru
- Shivani Bhoite
- Alejandro Benitez
- Umaraj Potla

GEO-LOCATION ANALYSIS

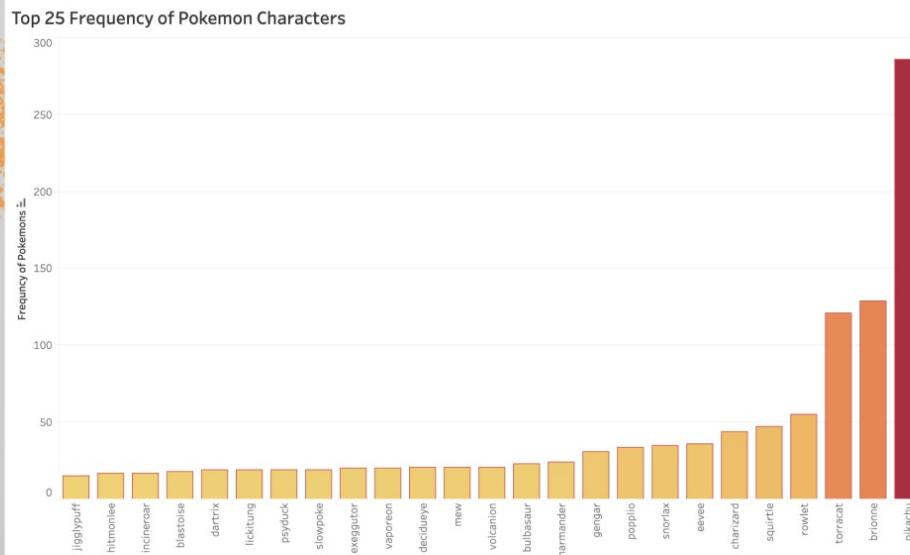
LDA

POKEMON GO FACTS
&
PREDICTIONS

SENTIMENT ANALYSIS

TOP POKEMON CHARACTERS TWEETED

Well most tweeted pokemon character is no more a mystery!!



FUTURE PREDICTIONS?

HOW MANY TWEETS SHOULD WE EXPECT IN FUTURE?

We predicted the popularity of the game using tweets counts on hourly-daily basis.

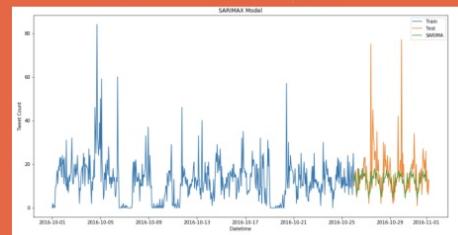
HOW DID WE DO?

TIME SERIES ANALYSIS
GRAPH

HOW DID WE DO?

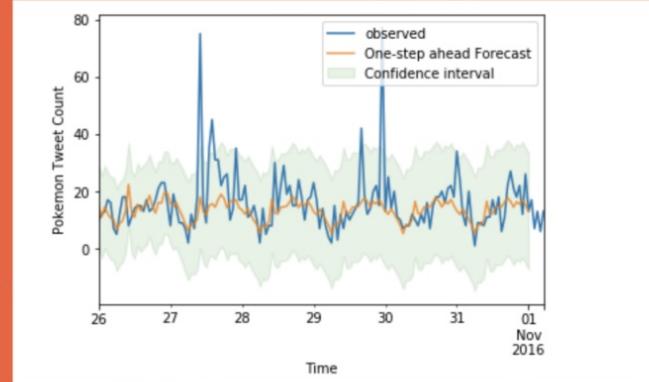
- Time Series Analysis by State Space Methods
'statespace' library
- **SARIMAX** - Seasonal Autoregressive Integrated Moving-Average
- SARIMAX works on a wider range of models by adding the estimation of additive and multiplicative seasonal effects, as well as arbitrary trend polynomials
- One-step-ahead prediction uses the true values of the training sample at each step to predict the next in-sample value.
- Dynamic predictions use one-step-ahead prediction up to some point in the dataset (specified by the dynamic argument); after that, the previous predicted values are used in place of the true values for each new predicted element.

Time Series Analysis using SARIMAX Model



MORE PREDICTIONS

LET'S ZOOM IN



Thank You!

POKEMON GO ON TWITTER ANALYSIS

WHY POKEMON GO?



- Anit Gupta
- Anil Ravuru
- Shivani Bhoite
- Alejandro Benitez
- Umaraj Potla

GEO-LOCATION ANALYSIS

LDA

POKEMON GO FACTS
&
PREDICTIONS

SENTIMENT ANALYSIS