

Same Mickey Different Cultures

Twitter Sentiment Analysis of
Disney Parks

Justin Huang





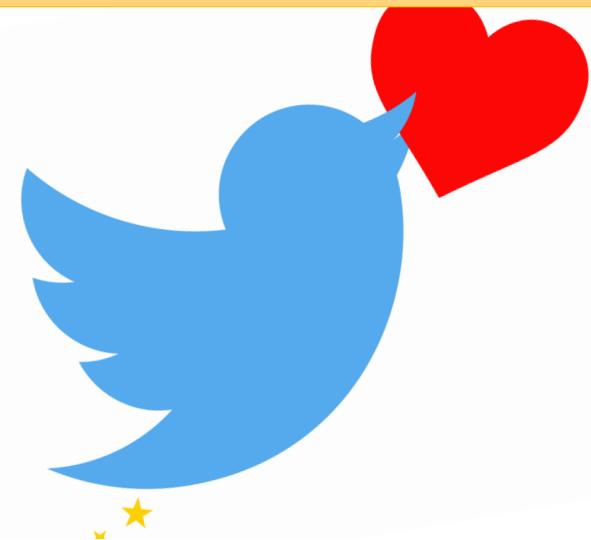
Intro

- Why Sentiment Analysis
 - By seeing trends of what people liked and disliked it can help relay information to Disney's business partners of what worked and didn't work.
- Imagineers
 - Imagineers are in charge of dreaming, designing and building Disney theme parks, attractions, cruise ships, resorts etc...
- Same name different Ownership
 - Oriental Lands
 - A Japanese leisure and tourism corporation headquartered in Urayasu, Chiba, Japan.
 - Tokyo Disneyland
 - Tokyo Disney Sea
 - Disney
 - Owners of several renown and loved character IPs
 - Disneyland
 - Disney California Adventure

Twitter API

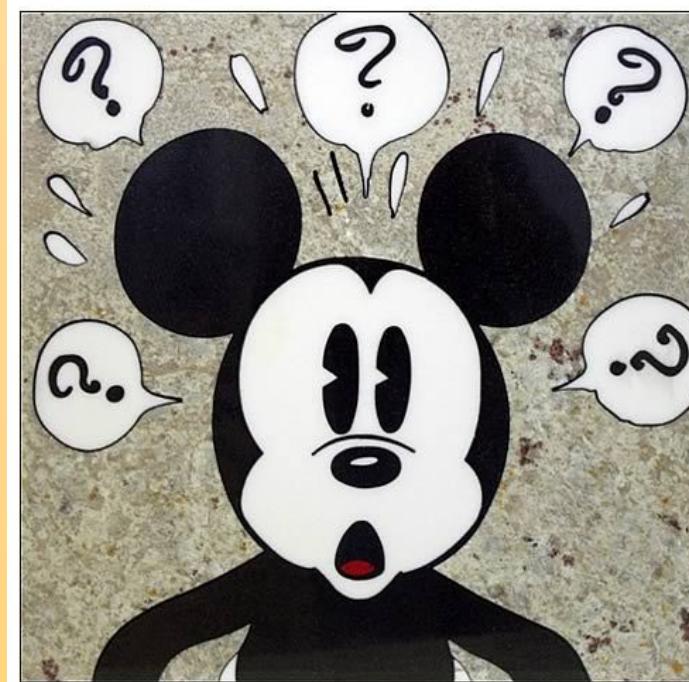
{"created_at": "Tue Sep 24 03:24:13 +0000 2019", "id": "1176336509575712768", "id_str": "1176336509575712768", "text": "RT @tdrblognote: \u2035\u203c\u3034c\u65b0\u3057\u2035"}, {"created_at": "Tue Sep 24 03:24:16 +0000 2019", "id": "1176336519969345536", "id_str": "1176336519969345536", "text": "@OrangeGrove55 Galaxy\u2019s Edge at Walt Disney World"}, {"created_at": "Tue Sep 24 03:24:19 +0000 2019", "id": "117633653139919360", "id_str": "117633653139919360", "text": "RT @hp_yec: \u206628\u2065e\u3068\u303a\u3093\u304b\u308a"}, {"created_at": "Tue Sep 24 03:24:25 +0000 2019", "id": "11763365796567744", "id_str": "11763365796567744", "text": "RT @gourmedty: New Halloween themed cups! #Shan"}, {"created_at": "Tue Sep 24 03:24:28 +0000 2019", "id": "1176336572305874944", "id_str": "1176336572305874944", "text": "Why would Disney need to promote #StarWars: #Galaxy"}, {"created_at": "Tue Sep 24 03:24:44 +0000 2019", "id": "11763366377992681472", "id_str": "11763366377992681472", "text": "RT @tdrblognote: \u2035\u203c\u3034c\u65b0\u3057"}, {"created_at": "Tue Sep 24 03:24:47 +0000 2019", "id": "1176336652878303232", "id_str": "1176336652878303232", "text": "RT @MezzoMikiD: \u206771\u2064eac\u303c\u2037\u303a\u303b\u303a"}, {"created_at": "Tue Sep 24 03:24:49 +0000 2019", "id": "11763366580829002752", "id_str": "11763366580829002752", "text": "You could win a vacation for four to a galaxy far"}, {"created_at": "Tue Sep 24 03:24:52 +0000 2019", "id": "1176336672872517633", "id_str": "1176336672872517633", "text": "RT @_erikagirl00: a bitch just wants to go to Disney"}, {"created_at": "Tue Sep 24 03:24:54 +0000 2019", "id": "1176336681152208896", "id_str": "1176336681152208896", "text": "RT @EricMProzek: Why would Disney need to promote"}, {"created_at": "Tue Sep 24 03:24:58 +0000 2019", "id": "1176336698705186816", "id_str": "1176336698705186816", "text": "RT @pinku_tai: \u20e0\u0e01\u20e2\u20e3\u20e35\u20e0\u0e35\u20e0\u0e14\u20e0\u0e35"}, {"created_at": "Tue Sep 24 03:25:00 +0000 2019", "id": "1176336706124963840", "id_str": "1176336706124963840", "text": "To Popcorn Bucket or not to Popcorn Bucket? That"}, {"created_at": "Tue Sep 24 03:25:04 +0000 2019", "id": "1176336721442566146", "id_str": "1176336721442566146", "text": "RT @kelseyr1101: My diet consists of Starbucks and"}, {"created_at": "Tue Sep 24 03:25:06 +0000 2019", "id": "117633673262876264", "id_str": "117633673262876264", "text": "RT @gourmedty: New Halloween popcorn bucket+ \n"}, {"created_at": "Tue Sep 24 03:25:09 +0000 2019", "id": "1176336741344526336", "id_str": "1176336741344526336", "text": "RT @MezzoMikiD: \u206771\u2064eac\u303c\u2037\u303a\u303b\u303a"}, {"created_at": "Tue Sep 24 03:25:17 +0000 2019", "id": "1176336777348411392", "id_str": "1176336777348411392", "text": "RT @MezzoMikiD: \u206771\u2064eac\u303c\u2037\u303a\u303b\u303a"}, {"created_at": "Tue Sep 24 03:25:20 +0000 2019", "id": "1176336790166355968", "id_str": "1176336790166355968", "text": "Late night sugar rush with my princess!! \n\n\n#tre"}, {"created_at": "Tue Sep 24 03:25:21 +0000 2019", "id": "1176336792372596737", "id_str": "1176336792372596737", "text": "RT @KatyFBrand: Passport checks at St Pancras after"}, {"created_at": "Tue Sep 24 03:25:22 +0000 2019", "id": "1176336797741117440", "id_str": "1176336797741117440", "text": "RT @MezzoMikiD: \u206771\u2064eac\u303c\u2037\u303a\u303b\u303a"}, {"created_at": "Tue Sep 24 03:25:24 +0000 2019", "id": "1176336805685161984", "id_str": "1176336805685161984", "text": "RT @MezzoMikiD: \u206771\u2064eac\u303c\u2037\u303a\u303b\u303a"}, {"created_at": "Tue Sep 24 03:25:32 +0000 2019", "id": "1176336838992089088", "id_str": "1176336838992089088", "text": "i want to go to disneyland", "source": "\u2003ca href"}, {"created_at": "Tue Sep 24 03:25:34 +0000 2019", "id": "1176336848899076098", "id_str": "1176336848899076098", "text": "RT @DisneyParks: Disney\u2019s Grand Californian"}, {"created_at": "Tue Sep 24 03:25:39 +0000 2019", "id": "1176336867235238624", "id_str": "1176336867235238624", "text": "RT @fairejk: im ready to have my heart broken http://t.co/2X"}, {"created_at": "Tue Sep 24 03:25:40 +0000 2019", "id": "117633687509904128", "id_str": "117633687509904128", "text": "RT @ochamz_dinsey: \u2030df\u203c\u303d\u2030d\u303f\u203c\u203b"}, {"created_at": "Tue Sep 24 03:25:49 +0000 2019", "id": "1176336909762727937", "id_str": "1176336909762727937", "text": "RT @isneylandToday: .@TheGoldbergsABC are on vac"}, {"created_at": "Tue Sep 24 03:25:54 +0000 2019", "id": "1176336930838958081", "id_str": "1176336930838958081", "text": "\u2030d\u203c\u303d\u2030d\u303f\u203c\u203b"}, {"created_at": "Tue Sep 24 03:26:00 +0000 2019", "id": "1176336959705020800", "id_str": "1176336959705020800", "text": "RT @cartarsauce: Since you asked, \nTHIS is what"}, {"created_at": "Tue Sep 24 03:26:03 +0000 2019", "id": "117633697794987817", "id_str": "117633697794987817", "text": "RT @lorrreanaa: I NEED to go to Disneyland! :)", "source": "\u2003ca href"}, {"created_at": "Tue Sep 24 03:26:04 +0000 2019", "id": "1176336972224311298", "id_str": "1176336972224311298", "text": "RT @KatyFBrand: Passport checks at St Pancras after"}]

- Why Twitter
 - Twitter is popular in Japan and the US
 - Twitter Data collected
 - Japanese and English Tweets
 - 270000 tweets for Halloween Event
 - 26000 for Chinese New Year event
 - Live Stream Tweets collected



Data Wrangling and cleaning

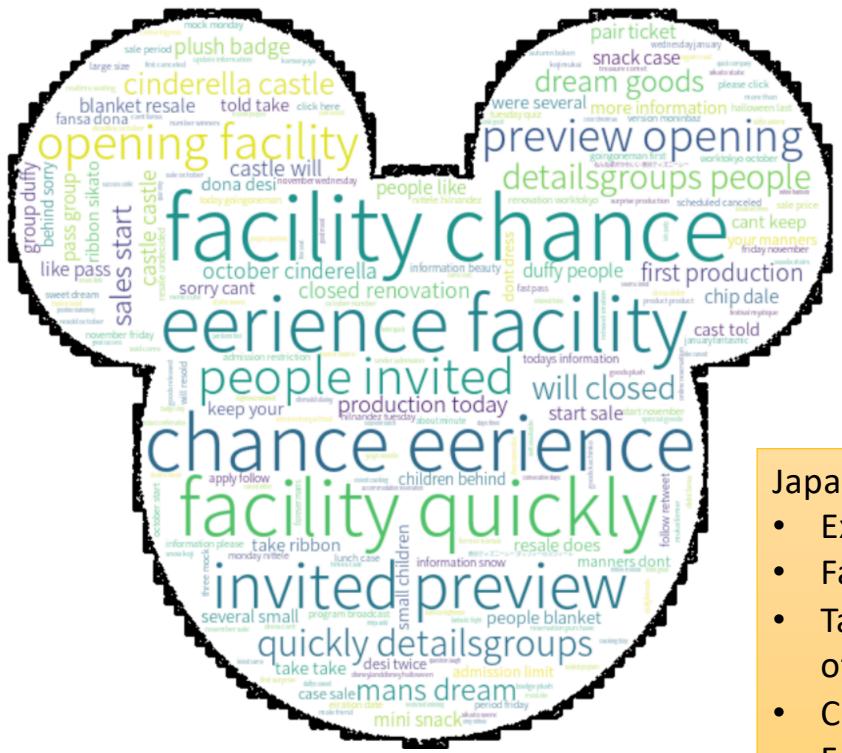
- Collected Disney tweets based off Anaheim and Tokyo Disney researched words
- Json file was structured using tweepy.
- Had to find Text through 4 deeply nested dictionaries
 - Extended Tweets
 - Quoted Status
 - Retweeted Status
 - Extended Entities
- Columns were formatted for Datetime and categorical
- Longitude and Latitude added back in using Geocoder
- Used office 365 translation to translate Japanese Tweets to English
- Then preprocessed the text
- Got the hashtags, mentions, links and emojis out



Exploratory Data Analysis

Most retweeted during Halloween collection period

Exploratory Data Analysis



English Tweets

- Fake Annual Passes
 - Quick
 - First Tasting
 - Stopping Lovely

Japanese Translated Tweets

- Experience
 - Facility
 - Talking about the last shows of One Mans Dream
 - Cinderella's Castle
 - Food

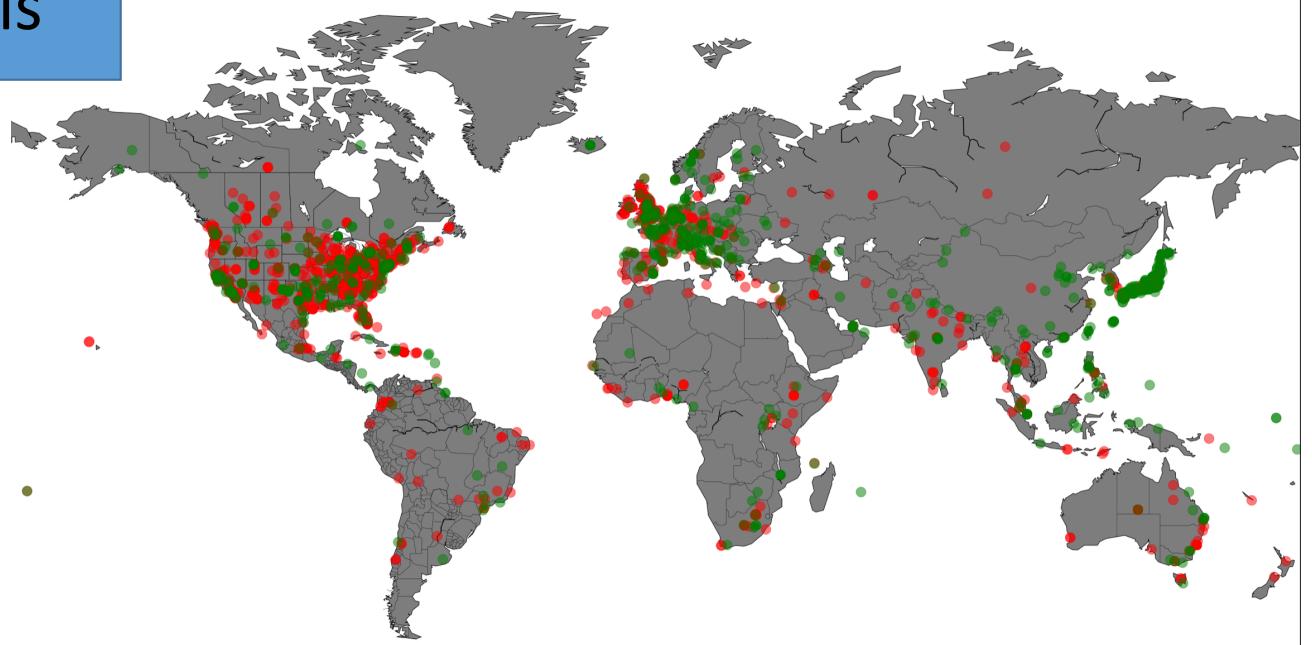


Exploratory Data Analysis

```
# Geocoder api
from opencage.geocoder import OpenCageGeocode
from pprint import pprint
key = # get api key from: https://opencagedata.com
geocoder = OpenCageGeocode(key)

en_lat = [] # create empty lists
en_long = []

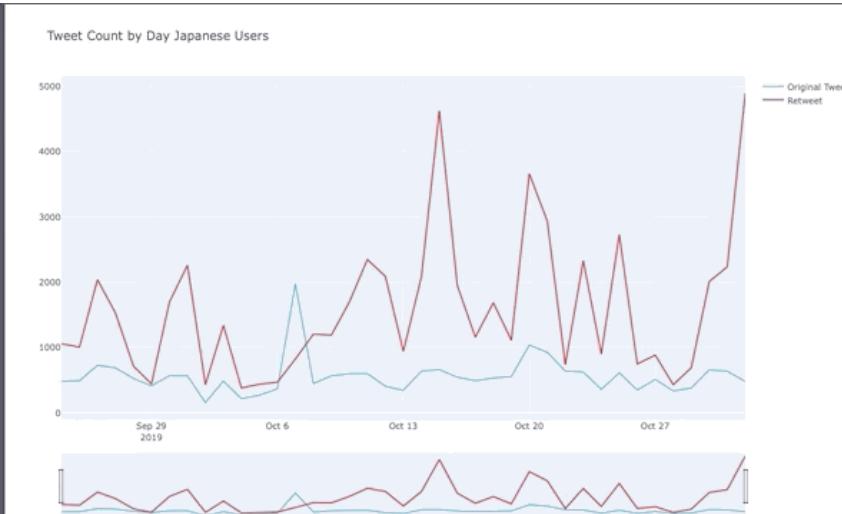
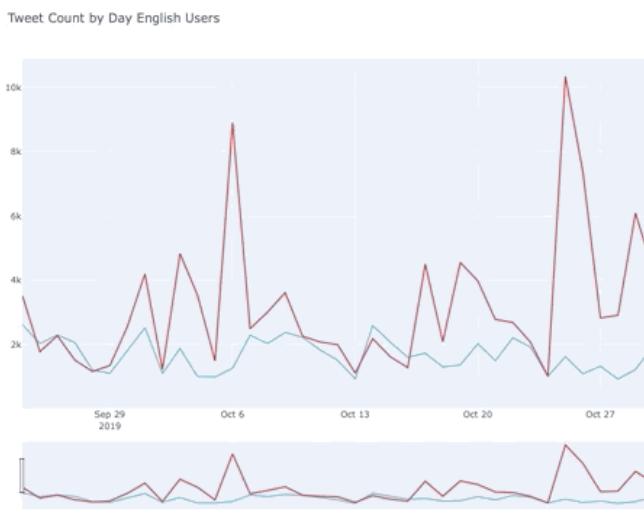
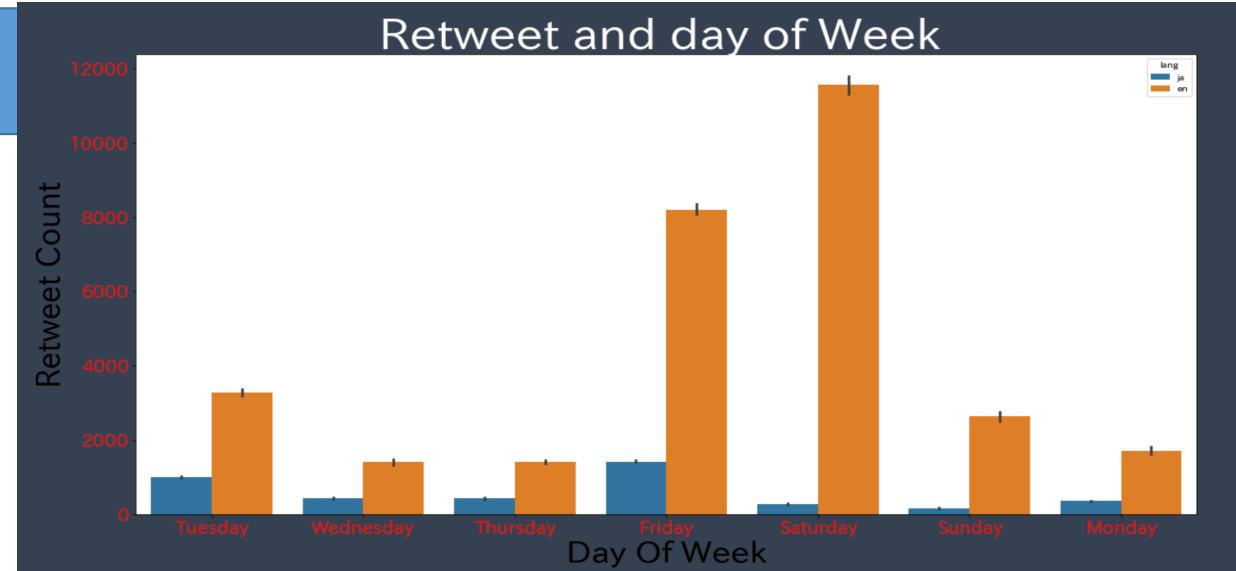
for x in en_loc3000.user_location: # iterate over rows in dataframe
    try:
        query = x
        results = geocoder.geocode(query)
        lat = results[0]['geometry']['lat']
        long = results[0]['geometry']['lng']
        en_lat.append(lat)
        en_long.append(long)
    except IndexError:
        en_lat.append(np.nan)
        en_long.append(np.nan)
```



- Geocoder to find longitude and latitude for user location
- Most US influencers were in the US marked in red
- Most Japanese influencers marked in green were in Japan and also in Europe

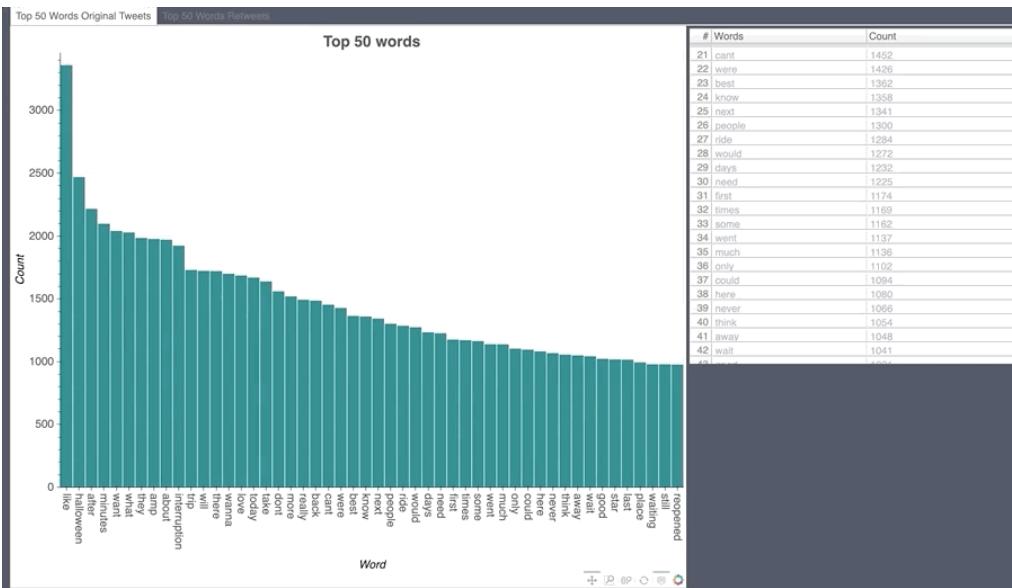
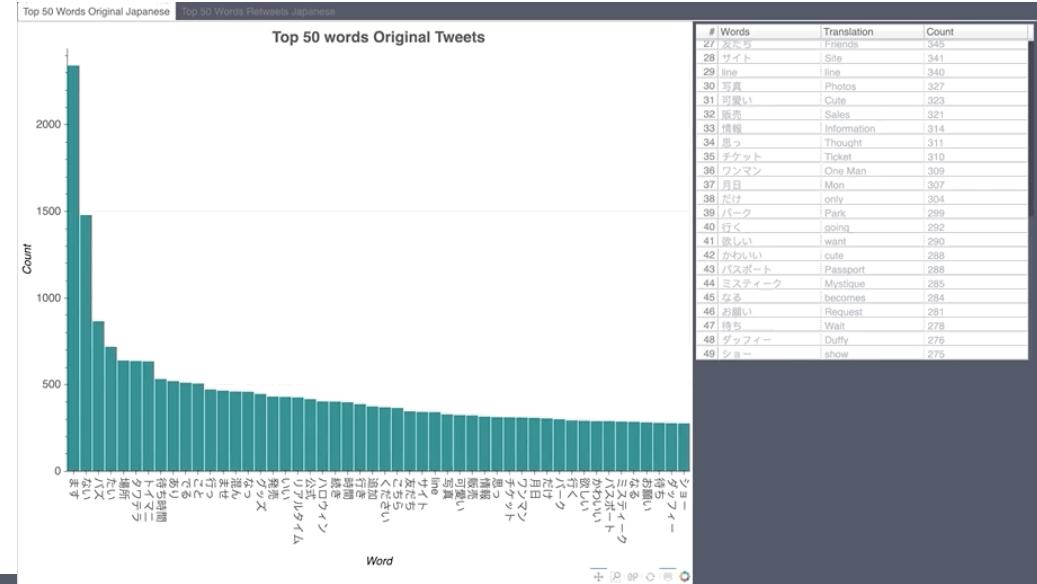
Exploratory Data Analysis

- For both English and Japanese users retweets and tweets happened mainly on Friday and Saturday.
- For peak periods the biggest peak was the weekend before Halloween for English users
- Japanese users was during the JAL event and Halloween day.



Exploratory Data Analysis

- Top 2 words original tweet English Users
 - Like and Halloween
- Top 2 words retweets
 - Fuck and annual

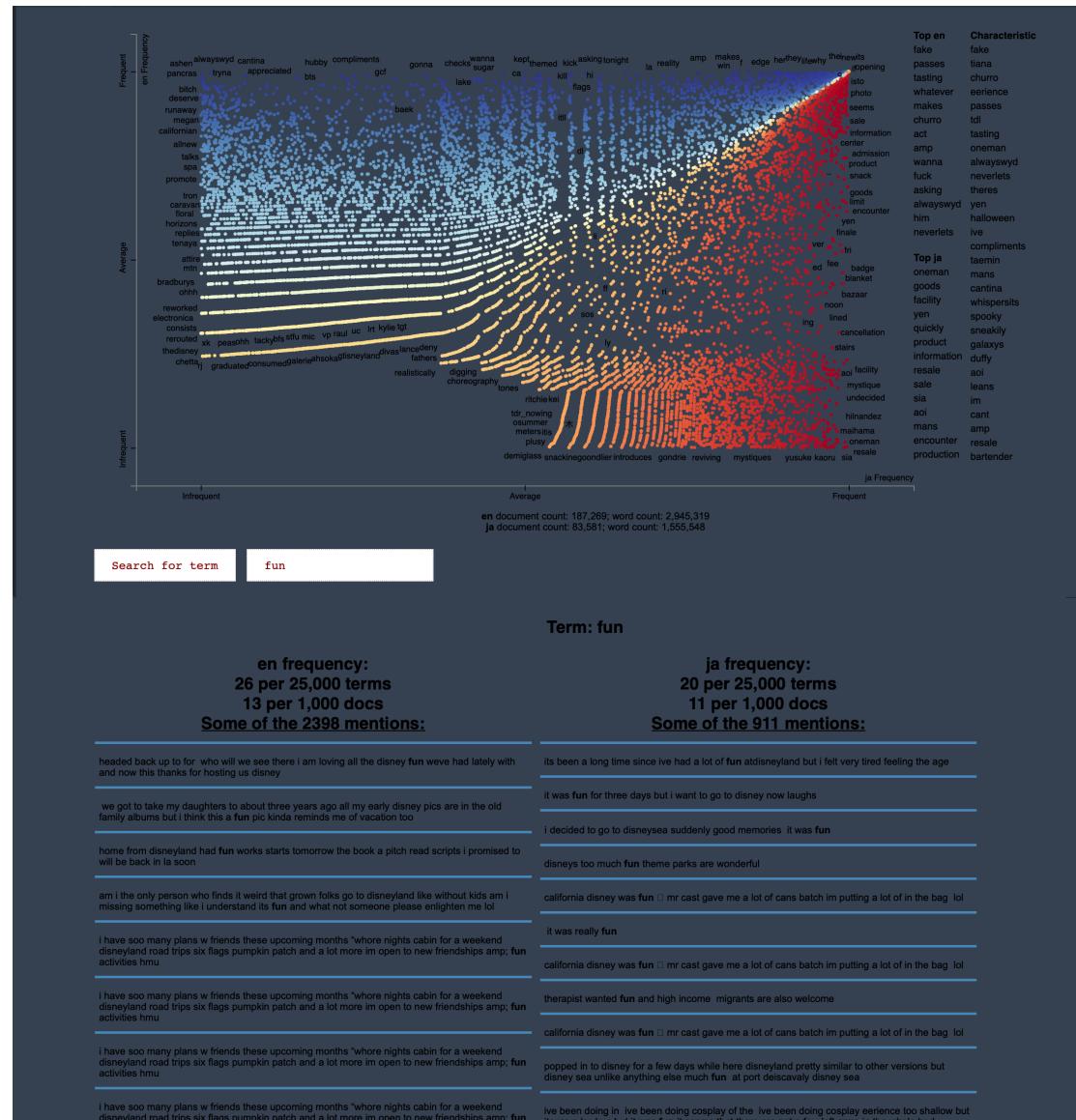


- Top 2 words original tweet Japanese Users
 - Monday and Sale
- Top 2 words retweets
 - Friend and Site

Exploratory Data Analysis

Scatter Text part of the Spacy Library

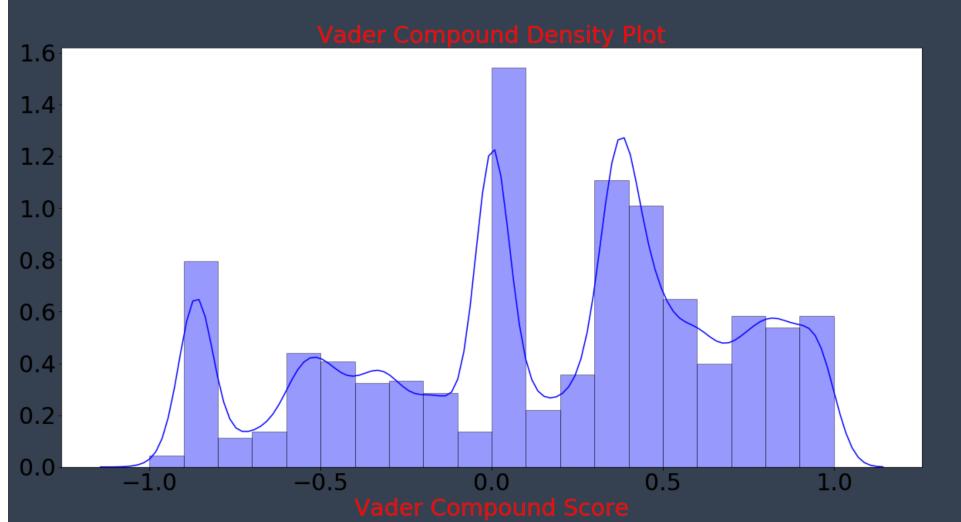
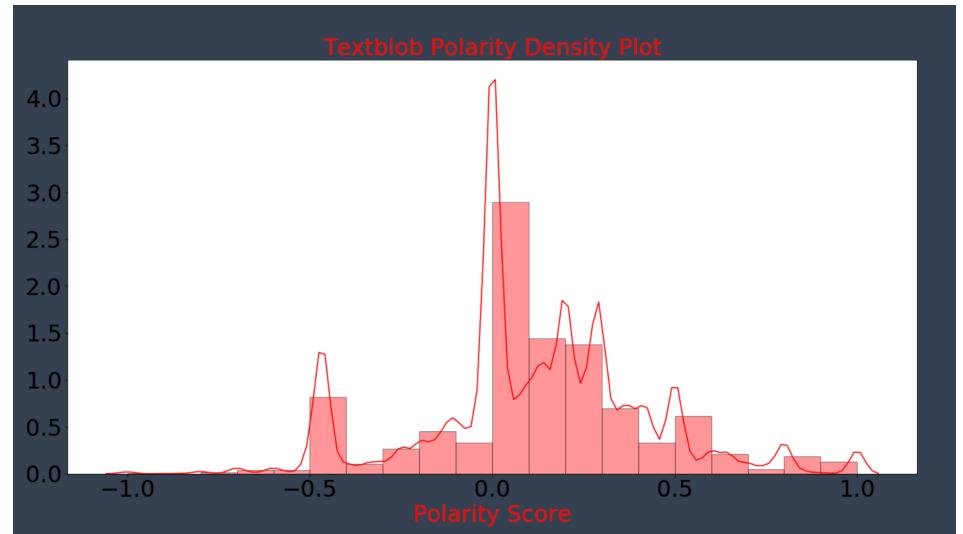
- Able to see word frequency across documents
- Able to do a search to see how many times it shows up per document or how many mentions
- Did Japanese documents that were translated compared to English documents



Stats

TextBlob and Vader sentiment compound score were used to compare:

- Japanese vs English user Sentiment towards the Disney Parks
- Tokyo Disney Resort vs Anaheim Disney Resort to see which one was more favorably received.



First had to test normality. Both were not normal. However should have large enough population for CLM to hold.

- Mean difference test showed that there was a difference between Japanese and English twitter users as well as the parks.
- Confirmed with non-parametric test, Mann- Whitney U Test. Both extremely low p-value at alpha 0.05

Stats

	user_favorite	user_follower	retweet	retweet_count	retweet_fav_count	text_length	polarity_tb	subjectivity_tb	vader_pos	vader_score
user_favorite	1	-0.0019	0.14	-0.023	-0.0093	0.06	0.0098	0.009	-0.0021	0.014
user_follower	-0.0019	1	-0.02	-0.0056	-0.0054	0.0062	0.0017	0.00089	-0.0005	0.0016
retweet	0.14	-0.02	1	0.22	0.22	0.18	-0.013	0.12	-0.0027	-0.019
retweet_count	-0.023	-0.0056	0.22	1	0.92	-0.16	-0.34	0.32	-0.092	-0.34
retweet_fav_count	-0.0093	-0.0054	0.22	0.92	1	-0.13	-0.26	0.28	-0.071	-0.27
text_length	0.06	0.0062	0.18	-0.16	-0.13	1	0.24	0.15	0.088	0.26
polarity_tb	0.0098	0.0017	-0.013	-0.34	-0.26	0.24	1	0.16	0.49	0.62
subjectivity_tb	0.009	0.00089	0.12	0.32	0.28	0.15	0.16	1	0.25	0.023
vader_pos	-0.0021	-0.0005	-0.0027	-0.092	-0.071	0.088	0.49	0.25	1	0.72
vader_score	0.014	0.0016	-0.019	-0.34	-0.27	0.26	0.62	0.023	0.72	1

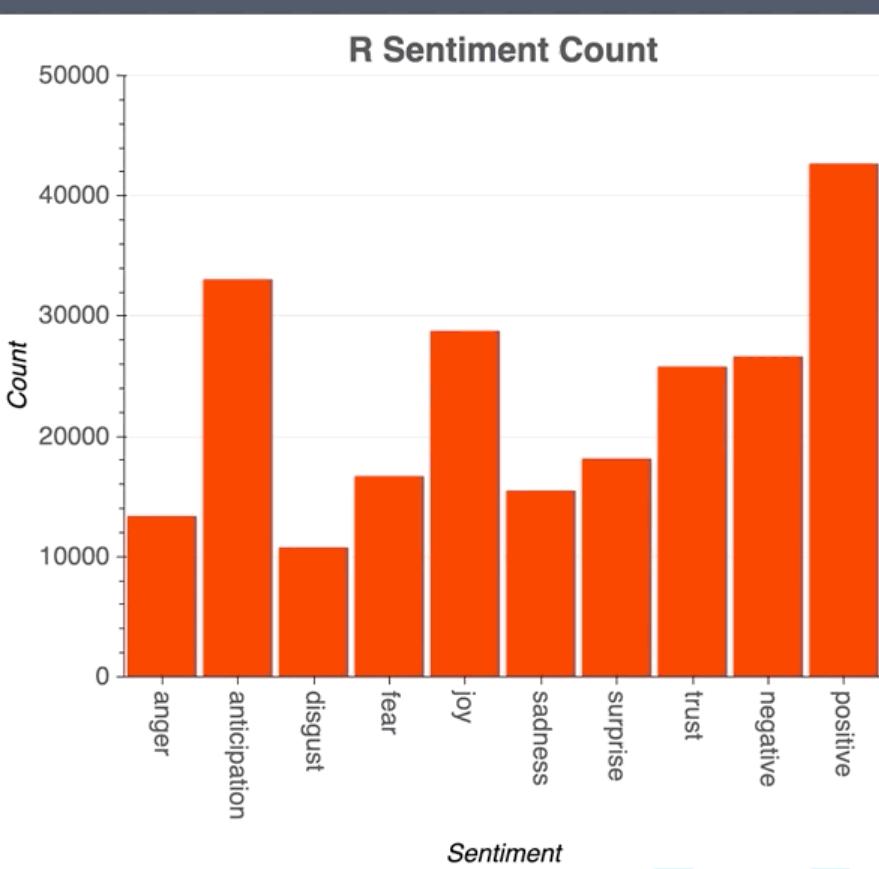
- Not too correlated.
- Only thing that is really correlated is retweet count and favorite count with a correlation of 0.92
- Vader pos and compound score of course would have a high correlation at 0.72

Machine Learning

- Text preprocess steps
 - Remove amp,&, /n, lxml
 - Remove href
 - Remove @ mentions
 - Remove RT
 - Lemmatize
 - Remove punctuation
 - Remove short words less than 2 characters
 - Remove special Japanese Punctuation
 - Tokenize



Machine Learning



#	Sentiment	Count
9	positive	42646
1	anticipation	33016
4	joy	28734
8	negative	26603
7	trust	25758
6	surprise	18106
3	fear	16645
5	sadness	15442
0	anger	13331
2	disgust	10719

Feature engineering and selection

- R Sentiment Emotion
 - Text Length
- Number of Capital letters
- Number of Punctuations
 - Number of Emojis
 - Number of Hashtags

Machine Learning

- Document term matrix selection
 - CountVectorizer best results overall
 - TFIDF Naïve Bayes TFIDF produced better results
 - Word2Vec worst results however can show word similarity
- Hyper Tuning
 - Random Forest best
 - Naïve Bayes second best
 - Logistic Regression third best
 - KNN worst

Important Features from Random Forest

```
[ (0.07264998189604532, 'tweet_len'),  
  (0.04641853216744782, 'cap_count'),  
  (0.04272377734677494, 'punc_count'),  
  (0.038893919132389765, 'hash_count'),  
  (0.02146655519317955, 'sadness'),  
  (0.018894953953547726, 'joy'),  
  (0.017807906473101077, 'fear'),  
  (0.016860613435313247, 'emoji_count'),  
  (0.014306119628743563, 'anger'),  
  (0.014305397515462942, 'anticipation'),  
  (0.012611158225566781, 'trust'),  
  (0.012546505297928749, 'disgust'),  
  (0.009523934533368706, 172),  
  (0.008760837448836236, 'surprise'),  
  (0.008609675090072008, 470),  
  (0.0050045463033912265, 547),  
  (0.004917858903756468, 937),  
  (0.004756201692522895, 496),  
  (0.004426082061033796, 318),  
  (0.004410321722824947, 182)]
```

Machine Learning

- Naïve Bayes vs Random Forest
 - Naïve Bayes worked with 15000 max features and (1,3) grams best
 - Random Forest worked with unigrams at 5000 features best



Machine Learning LDA



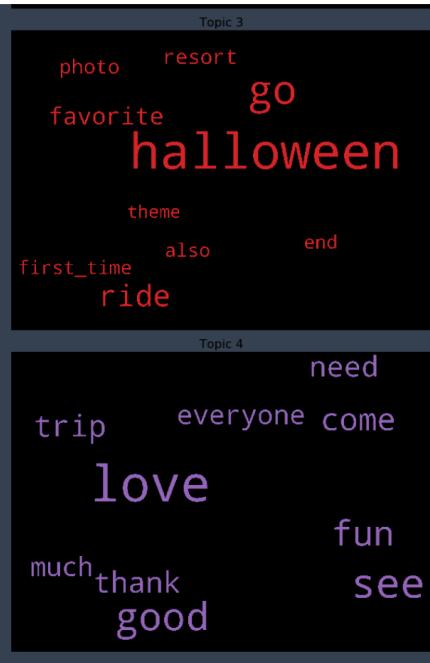
Japanese Positive

- One mans dream with Donald and friends
- Food and show with Mickey
- Good Happy and Halloween

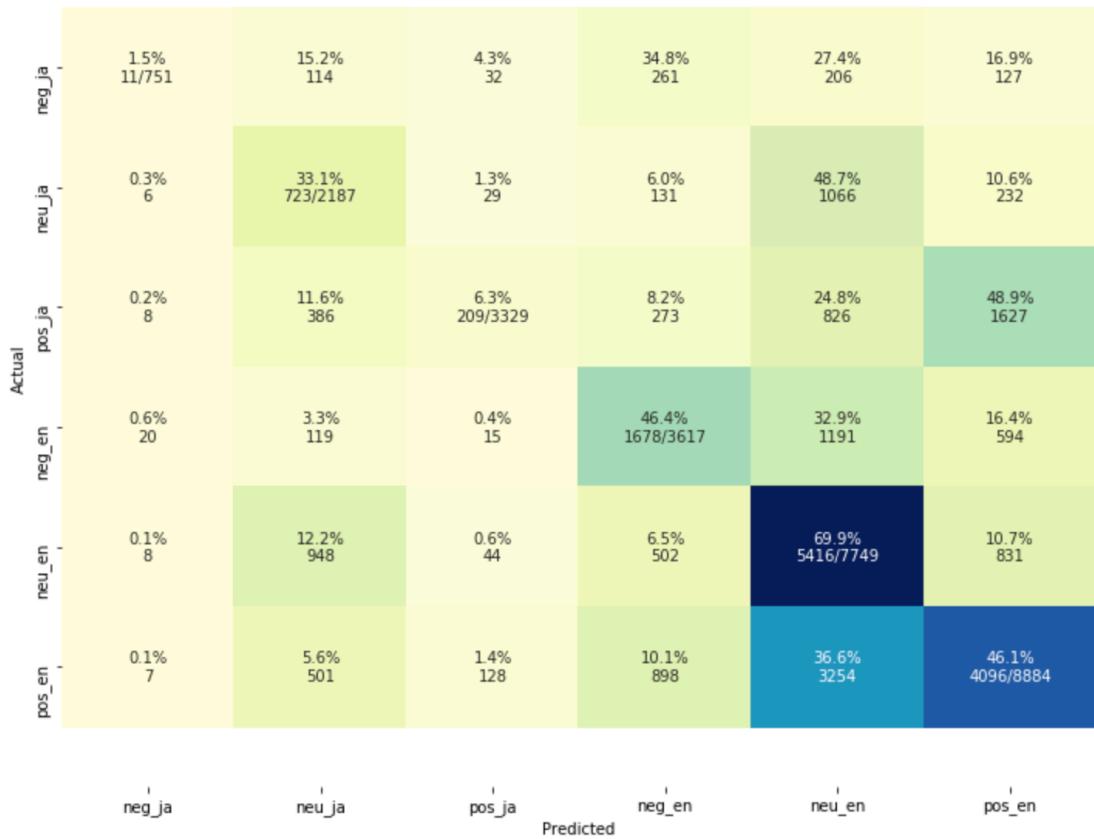


English Positive

- Halloween and resort first time for the ride
- Cute Amazing friends
- Love good thank and fun



Chinese New Year Model Validation



- Model performed worse
- Only 46 percent accuracy using RF single n-gram model.
- 33 percent 1 to 3-gram Naïve Bayes model