

# Sentiment Classification on UGC data

## 1. Text Processing

Apply the following procedures on text data.

- All lower-case
- Remove url
- Expand short forms
- Remove punctuations
- Tokenization
- Remove numbers
- Remove stop words

## 2. Model Selection

### 2.1 Label Encoding

Encode all labels in train and test set with numbers.

Label	Encoded Label
Negative	0
Neutral	1
Positive	2

### 2.2 Feature Engineer

Use TF-IDF to create features, set min\_df to 5. The model generates 7015 features. Reduce feature dimension using SVD to 100.

### 2.3 Modeling and Evaluation

Experiment with 6 models and the results are showed below.

	model_name	accuracy_score	precision_score	recall_score	f1_score
2	Random Forest	0.609339	0.597	0.59635	0.590595
3	AdaBoost	0.566426	0.551304	0.552421	0.545379
5	K Nearest Neighbor	0.54949	0.552702	0.544515	0.540771
1	Decision Tree	0.479448	0.470986	0.471049	0.470777
4	Gaussian Naive Bayes	0.480434	0.507912	0.476016	0.463437
0	Dummy	0.349227	0.345812	0.346062	0.345819

The best model in terms of overall accuracy is random forest. Tune the hyperparameter using 5-fold-cross validation and refit the model using best hyperparameters. The result is showed below.

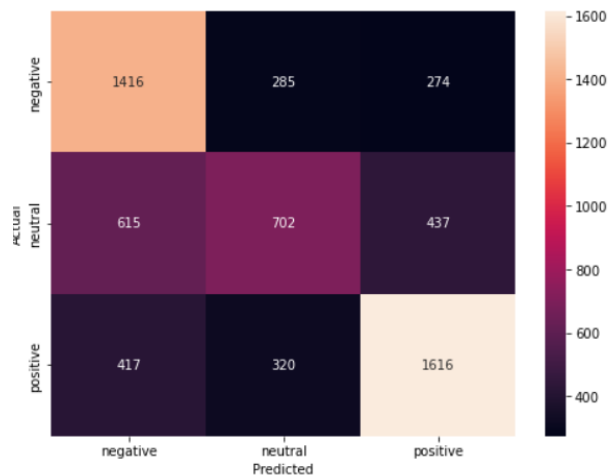
Metric	Score
Accuracy	0.6140
Precision	0.6033
Recall	0.6013
F1	0.5695

## 3. Result Evaluation

As showed below, positive sentiments are best classified, neutral sentiments are worst classified. There are two potential reasons.

- **Data Imbalance:** the number of positive and negative instances is much larger than the number of neutral instances. This can cause the model to favor the majority class and perform poorly on the minority class.
- Neutral sentiments often **lack clear positive or negative cues**, making them harder to classify.

Label	Precision	Recall
Negative	0.5784	0.7169
Neutral	0.5371	0.4002
Positive	0.6945	0.6868



## 4. Error Analysis

There are 6 types of errors.

### 4.1 Negative Classified as Neutral

cleaned_text	text
35 itchy miserable	I'm itchy and miserable!
41 cant sleep tooth aching	cant sleep ... my tooth is aching
48 started think cili really deep gon na survive turmoil gon na next aig	started to think that Citi is in really deep s & * i . Are they gonna survive the turmoil or are they gonna be the next AIG?
77 needs someone explain lambda calculus	needs someone to explain lambda calculus to him i :
89 burning cash chrysler gm stop financial tsunami bailout means taking handout	Are YOU burning more cash \$ than Chrysler and GM ? Stop the financial tsunami. "Where" "bailout" means taking a handout !
90 insects infected spinach plant	insects have infected my spinach plant :(
95 history exam studying ough	History exam studying ough
137 unfortunate stimulus plan put place twice help gm back american people led inevitable	It's unfortunate that after the Stimulus plan was put in place twice to help GM on the back of the American people has led to the inevitable
160 recovering surgerywishing julesrenner	Recovering from surgery .wishing @julesrenner was here :(
212 naive bayes using em text classification really frustrating	Naive Bayes using EM for Text Classification. Really Frustrating ...

- **Floppy spelling:** For example, misspell "can't" as "cant". So the negative "not" is not identified.

- Should not remove some **punctuations** like "!", '(', ':)' . "!" expresses excitement. ':( ' expresses negative feelings and ':)' expresses positive feelings. Should replace them with their meaning in words.
- The **TF-IDF model** considers the frequency and importance of each word in a corpus, so if the overall context and usage of the words in the corpus do not indicate a negative sentiment, the model may not classify the sentence as negative. This explanation applies to all errors. Data balance is important.

## 4.2 Negative Classified as Positive

	cleaned_text	text
10	firmlly believe obama pelosi zero desire civil charade slogan want destroy conservation.	I firmly believe that Obama / Pelosi have ZERO desire to be civil. It's a charade and a slogan, but they want to destroy conservation.
13	dear nike stop flywire shit waste science ugly love	dear nike , stop with the flywire. that shit is a waste of science, and ugly. love. @vincents24x
15	talking guy last night telling die hard spurs fan also told hates lebron james	I was talking to this guy last night and he was telling me that he is a die hard Spurs fan. He also told me that he hates LeBron James.
17	lebron beast still cheering all end	Lebron is a Beast. but I'm still cheering 4 the A..til the end.
32	played android google phone slide screen scares would break fucker fast still peddle phone	Played with an android google phone. The slide out screen scares me I would break that fucker so fast. Still peddle my iPhone.
34	omg bored tattoosos itchy help aha	omg so bored & my tattoosos are so itchy I I help I aha ->
36	no not itchy maybe later lol	no , I'm not itchy for now. Maybe later. lol
42	blah blah blah old old no plans today going back sleep guess	Blah. blah. blah same old same old. No plans today. going back to sleep I guess.
43	glad didnt bay breakers today freakin degrees san francisco uff	glad i didnt do Bay to Breakers today. it's 1000 freakin degrees in San Francisco uff
63	annoying new trend internets people picking apart michael lewis malcolm gladwell nobody wants read	annoying new trend on the internets. people picking apart michael lewis and malcolm gladwell. nobody wants to read that.

- Appearance of **positive words** in negative sentence: e.g. firmly believe, dear, cheering, glad
- Wrong label:** "Lebron is a Beast , but I'm still cheering 4 the A..til the end." It looks positive to me.

## 4.3 Neutral Classified as Negative

	cleaned_text	text
69	next time call nike	" Next time , I'll call myself Nike "
127	saw new night museum moviet vascoakayol	Just saw the new Night at the Museum movie. it was...okay. Jul 7   10
152	time warner cable pulls plug griffriend experience vawwtinyurlcom	Time Warner Cable Pulls the Plug on "The Griffriend Experience" - ( www.tinyurl.com / m596K )
165	climate focus turns beijing united nations us european governments called china cros	Climate focus turns to Beijing. The United Nations , the US and European governments have called on China to do so.
199	time buy gm car	is now the time to buy a GM car ?
162	dentist tomorrow brush well morning like make hair rice get cut why	Dentist tomorrow. Have to brush well in the morning. Like I make my hair all nice before I get it cut. Why ?
173	waiting line sawfaway	waiting in line at sawfaway.
179	jake going sawfaway	Jake's going to sawfaway I
180	found sawfaway picking staples	Found a sawfaway. Picking up a few staples.
184	normal weight get normal eating blog	Your Normal Weight ( and How to Get There ) ? Normal Eating Blog

- Using **punctuations without space** is a problem. After removing punctuations, multiple words are joined

## 4.4 Neutral Classified as Positive

	cleaned_text	text
9	check video president obama white house correspondents dinner	Check this video out -- President Obama at the White House Correspondents' Dinner
29	need suggestions good r filter canon got pls dm	need suggestions for a good IR filter for my canon 40D ... got some ? pls DM
30	checked google business btp shows second entry huh good	I just checked my google for my business - btp shows up as the second entry I Huh. is that a good or ba... ? -> identify
46	san francisco today suggestions	San Francisco today. Any suggestions ?
51	way see star trek esquire	On my way to see Star Trek @ The Esquire
62	going see star trek soon dad	Going to see star trek soon with my dad
60	playing curl twitter api	playing with cURL and the Twitter API
62	playing java twitter api	playing with Java and the Twitter API
68	nike owns nba playoffs ads w lebron kobe carmelo	Nike owns NBA Playoffs ads w / LeBron, Kobe, Carmelo ?
70	new blog post nike sb dunk low premium white gum	New blog post. Nike SB Dunk Low Premium " White Gum "

- Appearance of **positive words** in neutral sentence: e.g. good, play

## 4.5 Positive Classified as Negative

	cleaned_text	text
4	fair enough think perfect	Fair enough. But I have the Kindle2 and I think it's perfect. )
21	lebron beast nobody nba comes even close	Lebron is a beast. nobody in the NBA comes even close.
27	customer innovation award winner booz allen hamilton	[ #MILUC09 ] Customer Innovation Award Winner. Booz Allen Hamilton --
61	hello twitter api	Hello Twitter API.
84	work til 5pm. lets go lakere	work til 5pm. lets go lakere I I I
110	ever amazing psyop goodby silvenstein partners hp go play effects	The ever amazing Psyop and Goodby Silvenstein & Partners for HP I Have to go play with After Effects now I
161	wrist still hurts get looked hate dr dentist scary places time watch eagle eye want join bet	My wrist still hurts. I have to get it looked at. I HATE the dr / dentist / scary places. ( Time to watch Eagle eye. If you want to join, bet I
165	studing math tomorrow exam dentist	is studing math. ) tomorrow exam and dentist. )
197	right lol get high expectations warren buffet style	right I I I LOL we'll get there I I I have high expectations. Warren Buffet style.
200	warren buffet became time richest man united states not working investing big idea lead future	Warren Buffet became ( for a time ) the richest man in the United States. not by working but investing in 1 Big idea which lead to the future

- Multiple emotions** within one text
- Should not remove **emoji**
- Wrong label:** e.g. 165, 200 are hard to tell if they are positive or sarcasm.

## 4.6 Positive Classified as Neutral

	cleaned_text	text
2	ok first assessment fucking rocks	Ok, first assessment of the #Kindle2 ... it fucking rocks I I I
40	going sleep bike ride	is going to sleep then on a bike ride. )
65	scrapbooking nic	is scrapbooking with Nic :D)
73	way totally inspired freaky nike commercial	By the way, I'm totally inspired by this freaky Nike commercial
112	zong	zong I I I I I have a Q2 I I I I I I I I
115	guess retiring start using developer woot fogosies	Guess I'll be retiring my G1 and start using my developer G2 woot fogosies
286	highly recommend malcolm gladwell tipping point next audiobook probably one well	I highly recommend Malcolm Gladwell's "The Tipping Point." My next audiobook will probably be one of his as well.
313	learning lambda calculus	Learning about lambda calculus. )
316	atebits finished watching stanford phone class session really appreciate rock	atebits I just finished watching your Stanford iPhone Class session. I really appreciate it. You Rock I
321	mcdonalds dinner goosood big mac meal	Just had McDonalds for dinner. :D it was goosood. Big Mac Meal. :)

- Should not remove **emoji and exclamations**.

## 5. Limitation

Did not experiment any deep learning models.

## References:

<https://towardsdatascience.com/multi-class-text-classification-with-scikit-learn-12f1e60e0a9f>

<https://medium.com/@robert.salgado/multiclass-text-classification-from-start-to-finish-f616a8642538>

YouTube channel: NormalizedNerd