

Deriving Sentiment from English Texts for the Creation of Digital Collages

Tianen (Benjamin) Liu
Advisor: Dr. Natalia Khuri



WAKE FOREST
UNIVERSITY

- [illegible]

The diagram illustrates a sentiment analysis pipeline. It starts with input data (SMS messages) on the left, which is processed by a central 'Sentiment Analysis Algorithm' (represented by a purple diamond). The output on the right shows three rows of sentiment classification results, each with an emoji and a corresponding JSON object:

- Positive Sentiment:** Represented by a happy emoji. The JSON object is:


```
{...
    "language": "english",
    "docSentiment": {
      "mixed": "0",
      "score": "0.72882",
      "type": "positive"
    }
  }
```
- Neutral Sentiment:** Represented by a neutral emoji. The JSON object is:


```
{...
    "language": "english",
    "docSentiment": {
      "mixed": "1",
      "score": "-0.15592",
      "type": "negative"
    }
  }
```
- Negative Sentiment:** Represented by an angry emoji. The JSON object is:


```
{...
    "language": "english",
    "docSentiment": {
      "mixed": "0",
      "score": "-0.31119",
      "type": "negative"
    }
  }
```

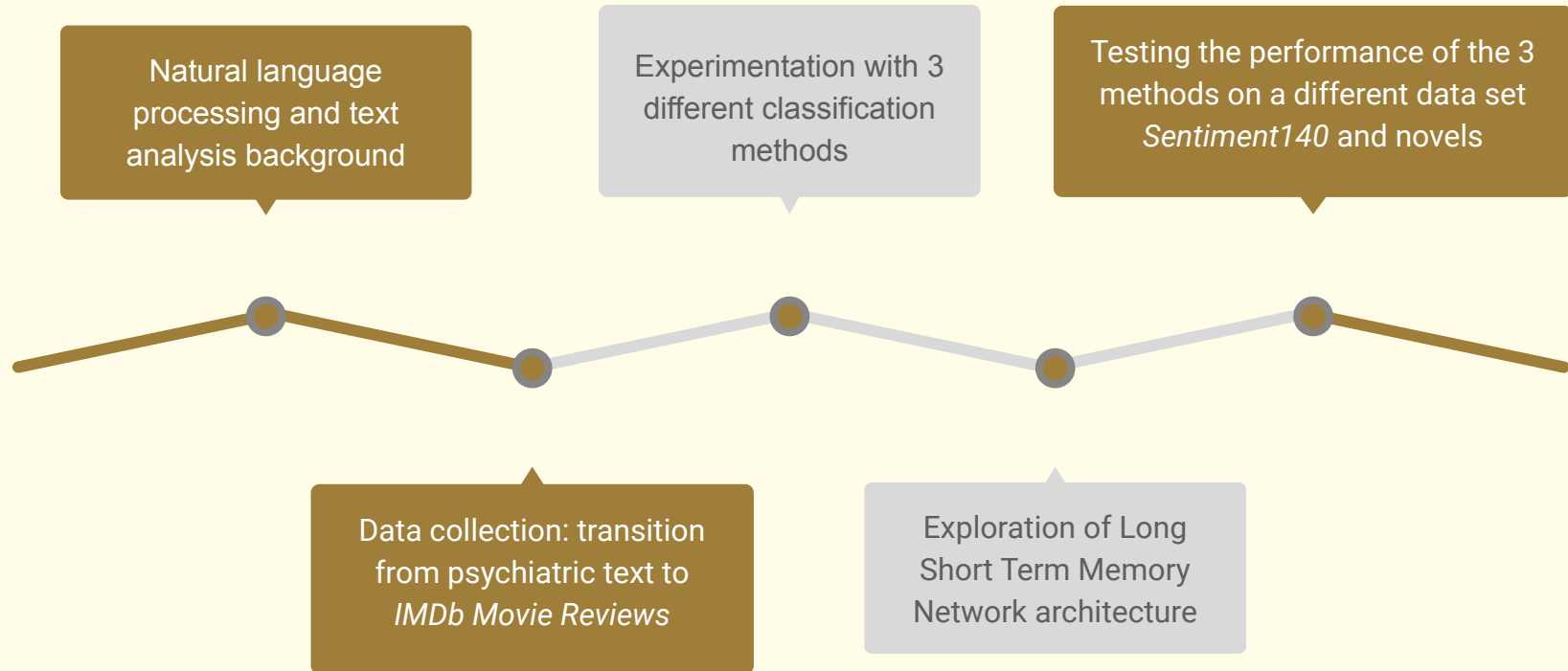
[illegible]

Figure 3



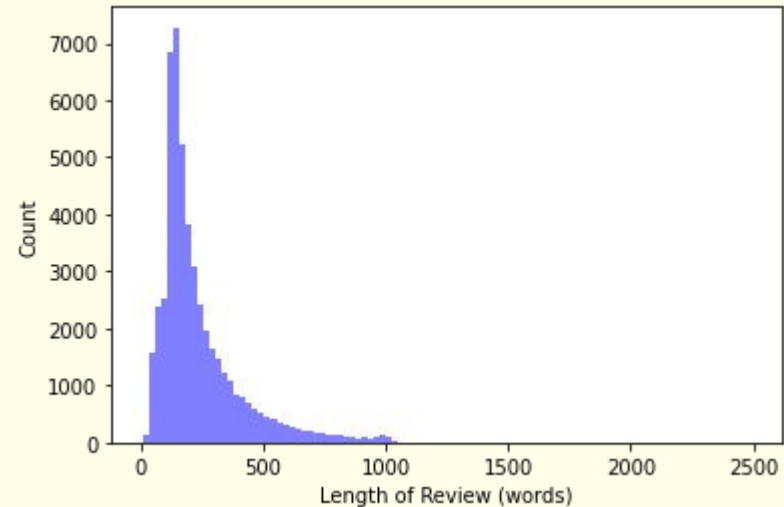
Explore machine learning methods to predict sentiments in different contexts

1. Train 3 machine learning algorithms with a data set of movie reviews
 2. Compare the performance of trained models
 3. Test their accuracy with the data sets from different domains
-





- **Source:** Internet Movie Database (IMDb)
Movie Reviews
- **Pre-processing:**
 - remove stopwords and punctuation
- **Data Overview:**
 - training size: 25000 reviews
(50% positive, 50% negative)
 - median length: 78 words
 - minimum length: 2 words
 - maximum length: 1079 words
 - length distribution





Algorithms used in this work

- Naïve Bayes Classifier
 - Turney Algorithm
 - Long Short Term Memory (LSTM) Network
-



Naïve Bayes Classifier



• Bayes Rule & Naïve Bayes

$$\Rightarrow P(c|d) = \frac{P(d|c)P(c)}{P(d)}$$

$$\Rightarrow c_{MAP} = \underset{c \in C}{\operatorname{argmax}} P(c|d) = \underset{c \in C}{\operatorname{argmax}} \frac{P(d|c)P(c)}{P(d)} = \underset{c \in C}{\operatorname{argmax}} \underbrace{P(d|c)}_{\text{likelihood}} \underbrace{P(c)}_{\text{prior}}$$

• Likelihood

$$\hat{P}(w_i|c_j) = \frac{\operatorname{count}(w_i, c_j) + 1}{\sum_{w \in V} (\operatorname{count}(w, c_j) + 1)} = \frac{\operatorname{count}(w_i, c_j) + 1}{(\sum_{w \in V} \operatorname{count}(w, c_j)) + |V|}$$

Notation

c : a class

C : all classes

c_{MAP} : max a posteriori (MAP) estimation

d : data (a set of texts)

w : a word in text

argmax_c : the class c that makes the following value maximum

V : vocabulary, distinctive words

$|V|$: size of vocabulary



Example: Naïve Bayes Classifier

- Vocabulary = 11
- Prior: probabilities of classes
 - $P(c_positive) = 2/4 = 1/2$
 - $P(c_negative) = 2/4 = 1/2$
- Likelihood: Conditional probabilities of individual words

Training	Data 1	i'm happy	Positive
	Data 2	avengers movie so good	Positive
	Data 3	hate movie very boring	Negative
	Data 4	college so stressful	Negative
Test	Data 5	college life so happy	Positive

$$\hat{P}(w_i|c_j) = \frac{\text{count}(w_i, c_j) + 1}{\sum_{w \in V} (\text{count}(w, c_j) + 1)} = \frac{\text{count}(w_i, c_j) + 1}{(\sum_{w \in V} \text{count}(w, c_j)) + |V|}$$

$$P(\text{college} | \text{positive}) = (0+1) / (6+11)$$

$$P(\text{life} | \text{positive}) = (0+1) / (6+11)$$

$$P(\text{so} | \text{positive}) = (1+1) / (6+11)$$

$$P(\text{happy} | \text{positive}) = (1+1) / (6+11)$$

$$P(\text{college} | \text{negative}) = (1+1) / (7+11)$$

$$P(\text{life} | \text{negative}) = (0+1) / (7+11)$$

$$P(\text{so} | \text{negative}) = (1+1) / (7+11)$$

$$P(\text{happy} | \text{negative}) = (0+1) / (7+11)$$

$$P(\text{positive} | \text{data}) = 0.0000239$$

$$P(\text{negative} | \text{data}) = 0.0000191$$



Turney Algorithm

- Part-of-speech tagging
- Bigram Extraction
“a beautiful day at Wake Forest”
- Learning the polarity of each phrase:
Pointwise Mutual Information (PMI)



1st word	2nd word	3rd word (not extracted)
JJ (<i>adjectives</i>)	NN, NNS (<i>nouns</i>)	anything
RB, RBR, RBS (<i>adverbs</i>)	JJ	not NN nor NNS
JJ	JJ	not NN nor NNS
NN, NNS	JJ	not NN nor NNS
RBm RBR, RBS	VB, VBD, VBN, VBG (<i>verbs</i>)	anything

$$Polarity(phrase) = PMI(phrase, "excellent") - PMI(phrase, "poor")$$

$$PMI(w1, w2) = \log_2 \frac{hits(w1 \text{ NEAR } w2)}{hits(w1)hits(w2)}$$





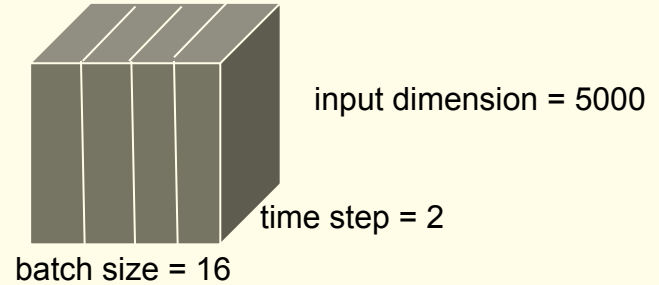
Long Short Term Memory Network

- **Pre-training**

- Padding: maximum input word length = 160

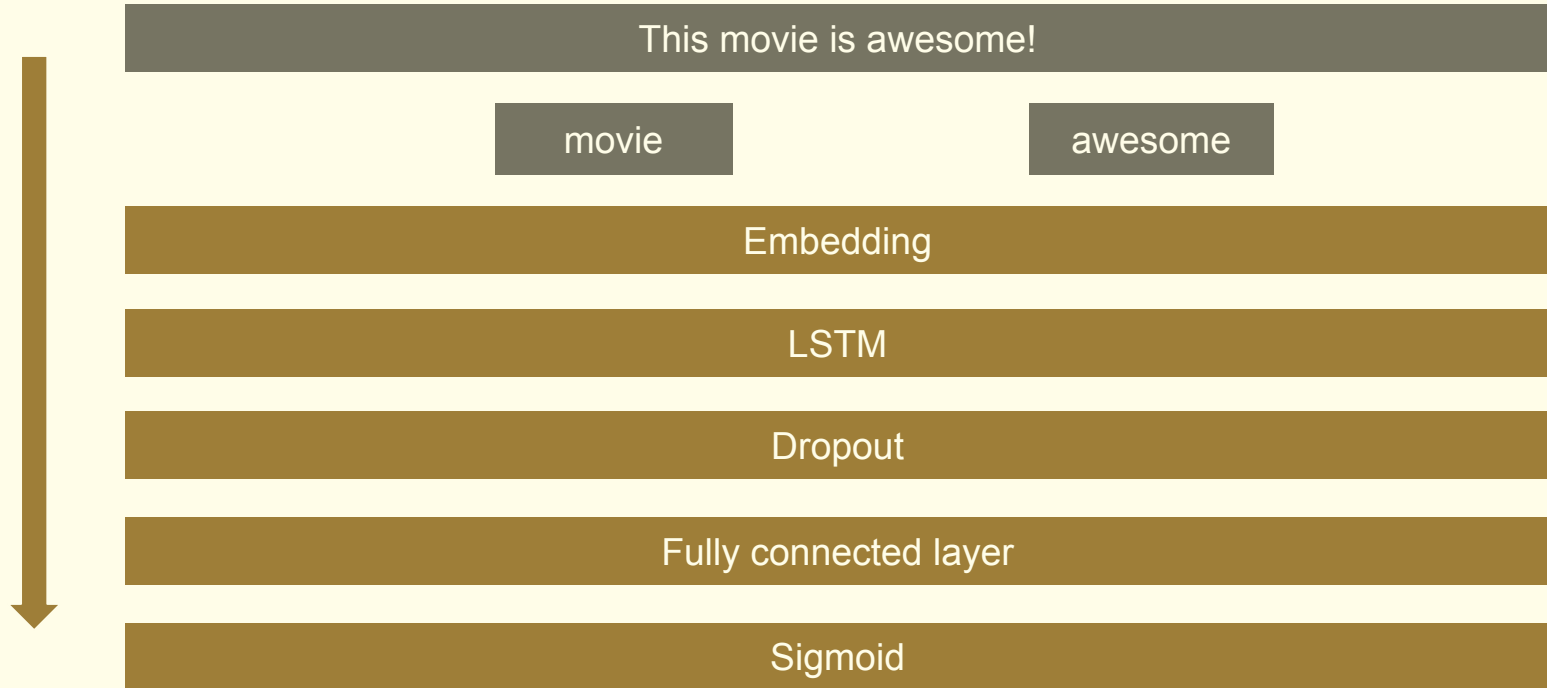
- **Parameters**

- embedding size = 64
- time step = 2
- input dimension = vocabulary size = 5000
- input unit = maximum input word length = 160
- batch size = 16
- LSTM input unit = 100
- dropout = 0.5



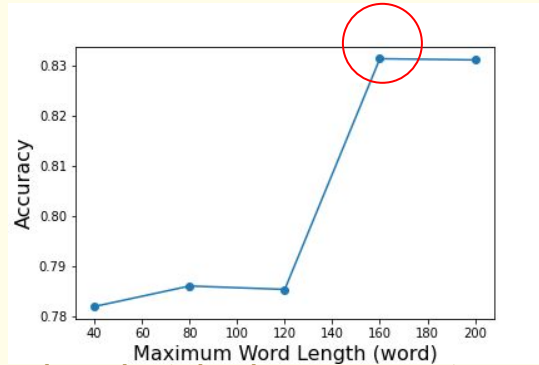


Architecture: LSTM Network

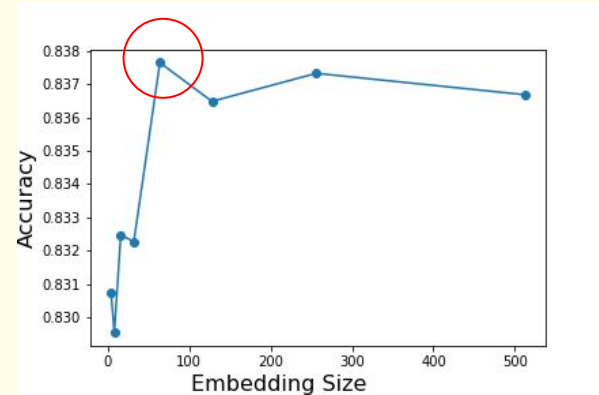


Parameter Tuning: LSTM Network

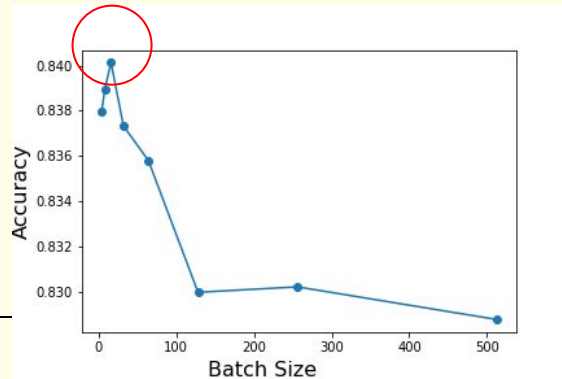
1. Explore maximum word length parameter



2. Explore embedding size parameter



3. Explore batch size parameter



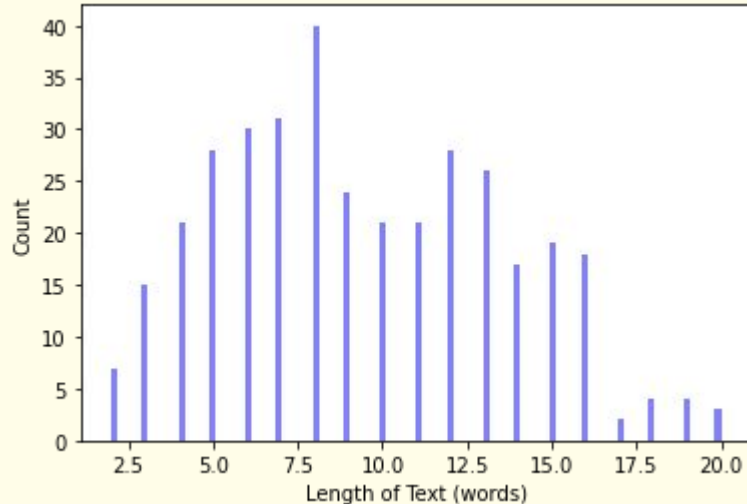
Comparison of the 3 methods

	Naïve Bayes	Turney Algorithm	LSTM Network
5-fold CV accuracy	83.73%	50.87%	86.17%
Standard deviation	0.0029	0.011	0.016

Aim 3: Application to Different Domain

Sample testing data of size 359 from *sentiment140*, a Twitter sentiment dataset

- Text length distribution



- Data format: csv

- 1st column values of 0 (neg) and 4 (pos)
- last column of texts

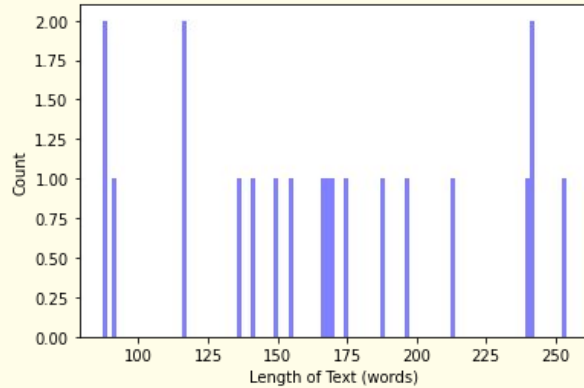
- Accuracies

- Naïve Bayes accuracy: 53.48%
- Turney Algorithm accuracy: 37.68%
- LSTM accuracy: 52.09%

4	26	Mon May 11	iphone app	CocoSavanna	downloading apps for my iphone! So much fun :-)	There literally is an app for just about anything.			
4	33	Mon May 11	visa	DreambigRa	good news, just had a call from the Visa office, saying everything is fine.....what a relief!	I am sick of scams out there! Stealing!			
4	34	Mon May 11	fredwilson	andrewwats	http://twurl.nl/epkr4b - awesome come back from @biz (via @fredwilson)				
4	35	Mon May 11	fredwilson	fredwilson	In montreal for a long weekend of R&R. Much needed.				
4	46	Thu May 14	("booz allen"	JoeSchueller	Booz Allen Hamilton has a bad ass homegrown social collaboration platform. Way cool!	#ttiv			
4	47	Thu May 14	("booz allen"	scottabel	[#MLUC09] Customer Innovation Award Winner: Booz Allen Hamilton -- http://ping.fm/c2hPP				

Sample testing data of size 10 from self-created novel excerpts

- Text length distribution

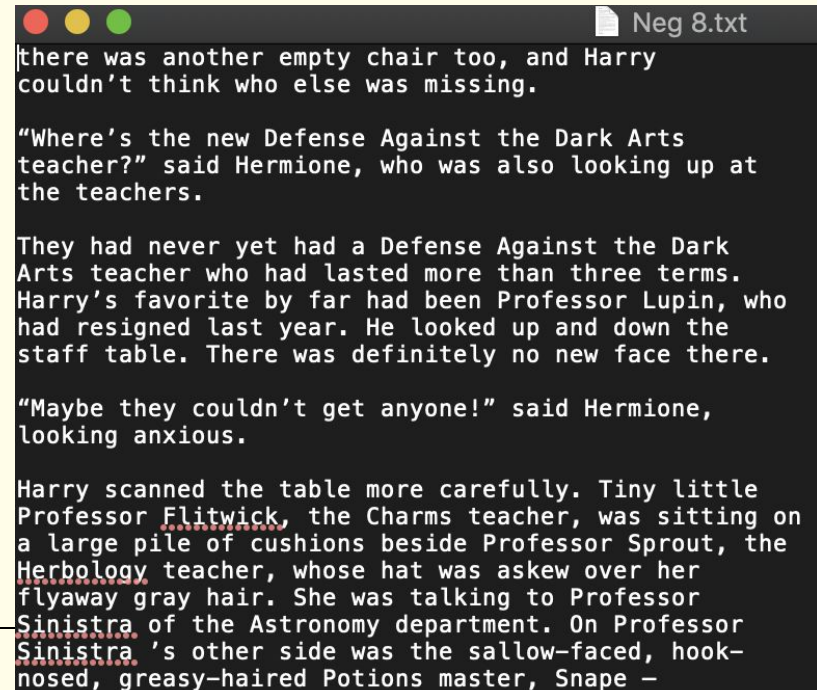


- Accuracies

- Naïve Bayes accuracy: 40.00%
- Turney Algorithm accuracy: 40.00%
- LSTM accuracy: 55.00%

- Data format: txt

Excerpted from *Harry Potter* and *Ulysses*



```
Neg 8.txt
there was another empty chair too, and Harry
couldn't think who else was missing.

"Where's the new Defense Against the Dark Arts
teacher?" said Hermione, who was also looking up at
the teachers.

They had never yet had a Defense Against the Dark
Arts teacher who had lasted more than three terms.
Harry's favorite by far had been Professor Lupin, who
had resigned last year. He looked up and down the
staff table. There was definitely no new face there.

"Maybe they couldn't get anyone!" said Hermione,
looking anxious.

Harry scanned the table more carefully. Tiny little
Professor Flitwick, the Charms teacher, was sitting on
a large pile of cushions beside Professor Sprout, the
Herbology teacher, whose hat was askew over her
flyaway gray hair. She was talking to Professor
Sinistra of the Astronomy department. On Professor
Sinistra's other side was the sallow-faced, hook-
nosed, greasy-haired Potions master, Snape -
```



- In 5 fold cross validation, LSTM has the best performance, followed by Naïve Bayes and Turney Algorithm.
 - Domain transfer is a challenge. LSTM and Naïve Bayes trained with movie reviews and Turney Algorithm fail to predict sentiments in other contexts due to overfitting and difference in phrasing words in different contexts.
 - Application of *text2collage* is a challenge.
-