

# REPORT

## 1. Background

Text sentiment analysis (Sentiment Analysis) refers to the process of analyzing, processing and extracting sentimental subjective text using natural language processing and text mining technology. At present, the research of text sentiment analysis covers many fields including natural language processing, text mining, information retrieval, information extraction, machine learning, and ontology. It has attracted the attention of many scholars and research institutions, and has continued to become natural language processing in recent years. And one of the hot issues in the field of text mining. Sentiment analysis tasks can be divided into chapter-level, sentence-level, word or phrase-level according to the granularity of their analysis; according to the type of text processed, they can be divided into sentiment analysis based on product reviews and sentiment analysis based on news reviews; according to their research tasks Types can be divided into sub-problems such as emotion classification, emotion retrieval and emotion extraction.

## 2. Purpose

The model is built through Python's NLTK library. The 10,000 emotive tweets in NLTK were used as the training set data, and the data were

processed through word segmentation and data normalization to construct the model data. Bayesian classification is used to train the model. Finally, the crawler is used to obtain all the reviews of a certain Amazon product as a test set and is brought into the trained model for training. After testing, the accuracy of the model can reach 78%.

### 3. load data

We use `negative_tweets.json` (5000 tweets with negative emotions) and `positive_tweets.json` under `twitter_samples`: (5000 tweets with positive emotions are used to train the model)

```
po_file_path = 'positive_tweets.json'
ne_file_path = 'negative_tweets.json'

positive_tweets = twitter_samples.strings(po_file_path)
negative_tweets = twitter_samples.strings(ne_file_path)
for i in range(6):
    print(positive_tweets[i])
    print(negative_tweets[i])
```

```
#FollowFriday @France_Inte @PKuchly57 @Milipol_Paris for being top engaged members in my community this week :)
hopeless for tmr :(
@Lamb2ja Hey James! How odd :/ Please call our Contact Centre on 02392441234 and we will be able to assist you :) Many tha
Everything in the kids section of IKEA is so cute. Shame I'm nearly 19 in 2 months :(
@DespiteOfficial we had a listen last night :) As You Bleed is an amazing track. When are you in Scotland?!
@Hegelbon That heart sliding into the waste basket. :(
@97sides CONGRATS :)
"@ketchBurning: I hate Japanese call him "bani" :(:("

Me too
yeaaaah yippppy!!! my acct verified rqst has succeed got a blue tick mark on my fb profile :) in 15 days
Dang starting next week I have "work" :(
@BhaktisBanter @PallaviRuhail This one is irresistible :)
#FlipkartFashionFriday http://t.co/EbZOL2VENM
oh god, my babies' faces :( https://t.co/9fcwGvaki0
```

## 4. Word segmentation

Word segmentation is to decompose long texts such as sentences, paragraphs, and articles into data structures in units of words to facilitate subsequent processing and analysis.

```
po_fenci_res = fenci(po_file_path)
be_fenci_res = fenci(ne_file_path)

print('Positive participle result: {}'.format(po_fenci_res))
print('Negative participle result: {}'.format(be_fenci_res))
```

```
Positive participle result: ['#FollowFriday', '@France_Inte', '@FKuchly57', '@Milipol_Paris', 'for', 'being', 'top', 'engaged', 'members', 'in', 'my', 'community', 'this', 'week', ':')].
Negative participle result: [['hopeless', 'for', 'tmr', ':('), ['Everything', 'in', 'the', 'kids', 'section', 'of', 'IKEA', 'is', 'so', 'cute', '.'], 'Shame', 'I'm', 'nearly', '19', 'in',
```

## 5. Data normalization

Data normalization includes the following steps:

Part-of-speech tagging

Junk data processing

Part of speech

```
def cleaned_list_func(evert_tweet):
    new_text = []
    cixing_list = pos_tag(evert_tweet)
    for word, cixing in cixing_list:
        word = re.sub('http[s]?://(?:[a-zA-Z]|[0-9]|[$-_@.&+#!*%\(\),]|(?:[0-9a-fA-F][0-9a-fA-F]))+', '', word)
        word = re.sub('@[A-Za-z0-9_]+', '', word)
        if cixing.startswith('NN'):
            pos = 'n'
        elif cixing.startswith('VB'):
            pos = 'v'
        else:
            pos = 'a'
        lemmatizer = WordNetLemmatizer()
        new_word = lemmatizer.lemmatize(word, pos)
        if len(new_word) > 0 and new_word not in string.punctuation and new_word.lower() not in stopwords.words(
            'english'):
            new_text.append(new_word.lower())
    return new_text
```

```
Positive tweet results after processing: ['#followfriday', 'top', 'engage', 'member', 'community', 'week', ':'), ['hey', 'james', 'odd', ':/', 'please', 'call', 'contact', 'centre', '023
original data: ['#FollowFriday @France_Inte @FKuchly57 @Milipol_Paris for being top engaged members in my community this week :'), '@Lamb2ja Hey James! How odd :/ Please call our Contact C
```

## 6. construct model data

```
def get_tweets_for_model(clean_tokens_list, tag):
    li = []
    for every_tweet in clean_tokens_list:
        data_dict = dict([token, True] for token in every_tweet)
        li.append((data_dict, tag))
    return li
```

positive data: [(('followfriday': True, 'top': True, 'engage': True, 'member': True, 'community': True, 'week': True, ':': True), 'Positive'), (('hey': True, 'james': True, 'odd': True, 'negative data: [(('hopeless': True, 'tmr': True, ':': True), 'Negative'), (('everything': True, 'kid': True, 'section': True, 'ikea': True, 'cute': True, 'shame': True, 'i'm': True, 'near

## 7. Use a crawler to get Amazon reviews

```
import emoji
url = 'https://www.amazon.com/Greenlights-Matthew-McConaughey-ebook/product-reviews/B086823SWK/ref=cm_cr_getr_d_paging_btm_next_?ie=UTF8&
headers = {'accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,image/avif,image/webp,image/apng,*/*;q=0.8,application/signed-ex
# 'accept-encoding': 'gzip, deflate, br',
'accept-language': 'en-US,en;q=0.9,zh-CN;q=0.8,zh;q=0.7',
'cache-control': 'max-age=0',
'cookie': 'ubid-main=135-3761994-2155044; skin=noskin; lc-main=en_US; csd-key=eyJ2IjoxLCJraWQiOiJmMGYzOTEiLCJrZXkiOiJmbHVkaEhsbWVhGmNBER2hoR
'downlink': '10',
'ect': '4g',
'referer': 'https://www.amazon.com/dp/B086823SWK',
'rtt': '0',
'sec-fetch-dest': 'document',
'sec-fetch-mode': 'navigate',
'sec-fetch-site': 'same-origin',
'sec-fetch-user': '?1',
'upgrade-insecure-requests': '1',
'user-agent': 'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_14_6) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/87.0.4280.88 Safari/537.36'}

if __name__ == '__main__':
    fp = open('climb.csv', mode='w', encoding='utf-8')
    fp.write('comment\\trate\\n')
    for i in range(1,79):
        page = url%(i, i)
        response = requests.get(page, headers=headers)
        response.encoding = 'utf-8'
        text = response.text
        html = etree.HTML(text)
        comments = html.xpath('//div[@id="cm_cr-review_list"]/div[@class="a-section review aok-relative"]')
        for comment in comments:
            c = comment.xpath('..//span[@class="a-size-base review-text review-text-content"]/span/text()')[0].strip()
            rate = comment.xpath('..//span[@class="a-icon-alt"]/text()')[0].strip()
            c = emoji.demojize(c)
            fp.write('%s\\t%s\\n' % (c, rate))
        print('Page ', i, ' saved!')
        time.sleep(1)
    fp.close()
```

	A	B	C
1	comment	rate	
2	Greenlights is a remarkable first book from an already renowned artist. Kind of a mashup of	4.0 out of 5 stars	
3	Worst book ever written. The only redeeming value here is that it gives you a glimpse into th	1.0 out of 5 stars	
4	Amazing book, clever and funny. Not your typical memoir, by design. As a huge fan, I enjoye	5.0 out of 5 stars	
5	Already received and enjoying it immensely. Beautifully written with candor, humor and refle	5.0 out of 5 stars	
6	Excellent book !! Matthew is far more complex than I ever realized. " Y" all quit bitching ab	5.0 out of 5 stars	
7	My book arrived in great shape, no damage, tears or nothing. My 14 year old son eyed it so I	5.0 out of 5 stars	
8	All races matter. We' ve all experienced discrimination. But, this statement takes the cake. W	1.0 out of 5 stars	
9	Question! Why is the dust cover at least an inch shorter than the book . Looks weird. Since it'	1.0 out of 5 stars	
10	I have always liked Matthew McConaughey' s dramatic roles and was anxious to read his bo	1.0 out of 5 stars	
11	Absolutely amazing! I have no been able to put it down since I' ve opened the box.	5.0 out of 5 stars	
12	The writing is entertaining, but the message is disturbing. Written by an abused, resilient, sui	2.0 out of 5 stars	
13	As one who never engages in needless exaggerations, I can honestly say that this is absolute	1.0 out of 5 stars	
14	I am a UT Alum and a big fan of Matthew McConaughey, but if you were expecting some tho	2.0 out of 5 stars	
15	It just delivered and unfortunately this is how it came. :weary_face: whyyy?!	1.0 out of 5 stars	

## 8. Train and test the model

```
def train_model(train_data):
    from nltk import NaiveBayesClassifier
    model = NaiveBayesClassifier.train(train_data)
    return model

def test(test_text):
    from nltk.tokenize import word_tokenize
    custom_tokens = cleaned_list_func(word_tokenize(test_text))
    result = dict([token, True] for token in custom_tokens)
    return result
```

## 9. Result

```
positive data: [({'#followfriday': True, 'top': True, 'engage': True, 'member': Tru
negative data: [({'hopeless': True, 'tmr': True, ':(': True}, 'Negative'), ({'every
[1, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1
```

		0	1	2	3
0	greenlights is a remarkable first book from an...	4.0	1	1	
1	worst book ever written. the only redeeming v...	1.0	0	0	
2	amazing book, clever and funny. not your typic...	5.0	1	1	
3	already received and enjoying it immensely. be...	5.0	1	1	
4	excellent book matthew is far more complex th...	5.0	0	1	
..	...	...	...	...	...
770	arrived with marks on the cover. was expecting...	1.0	1	0	
771	we received our book today and it was damaged...	1.0	0	0	
772	i love this book	5.0	1	1	
773	un de mes livres préféré. juste sublime.	5.0	0	1	
774	not what i was expecting. but its kinda ok.	2.0	1	0	

```
[775 rows x 4 columns]
```

```
Accuray: 0.7858064516129032
```

```
Process finished with exit code 0
```

The accuracy is about 78.6%