

Effective Text Data Cleaning using Python

Benefits of mining for a brand?

You can do sentiment analysis to discover customer's sentiments for a brand

You can measure brand popularity using the actively engaged tweeters

It is used to identify the pain points of customers i.e. customer relationship management

It is widely used for predictions and forecasting



The Business Problem

Let's say, we want to find the features of an Apple Phone which are most popular amongst the fans on Twitter.

What to do next?

We've extracted all the tweets related to consumer opinions of iPhone. Here's a sample tweet on which we'll perform data cleaning.

TWEET
"I lov my :3 Iphone & you're awsm apple. Display is Awesome, sooo haaaaaappy :3 http://www.apple.com"

Steps for Data Cleaning

Escaping HTML characters

Code

```
import HTMLParser
html_parser = HTMLParser.HTMLParser()
tweet = html_parser.unescape(original_tweet)
```

Output

```
>> "I lov my :3 Iphone & you're awsm apple. Display is Awesome, sooo haaaaaappy :3 http://www.apple.com"
```

Decoding data

Code

```
tweet = original_tweet.decode("utf8").encode("ascii",ignore')
```

Output

```
>> "I lov my :3 Iphone & you're awsm apple. Display is Awesome, sooo haaaaaappy :3 http://www.apple.com"
```

Apostrophe Lookup

Code

```
APFPOSTOPHES = ["'a": " is", "'e": " are", ...] ## Need a huge dictionary
words = tweet.split()
reformed = APFPOSTOPHES[word] if word in APFPOSTOPHES else word for word in words
reformed = " ".join(reformed)
```

Outcome

```
>> "I lov my :3 Iphone & you are awsm apple. Display is Awesome, sooo haaaaaappy :3 http://www.apple.com"
```

Removal of Stop-Words

When data analysis needs to be data driven at the word level, the commonly occurring words (stop-words) should be removed. One can either create a long list of stop-words or one can use predefined language specific libraries.

Removal of Punctuations

All the punctuation marks according to the priorities should be dealt with. For example: ".,,," are important punctuations that should be retained while others need to be removed.

Removal of Expressions

Textual data (usually speech transcripts) may contain human expressions like laughing, crying, audience paused. These expressions are usually non relevant to content of the speech and hence need to be removed.

Split Attached Words

Code

```
cleaned = " ".join(re.findall('[A-Z]([A-Z])', original_tweet))
```

Outcome

```
>> "I lov my :3 Iphone & you are awsm apple. Display is Awesome, sooo haaaaaappy :3 http://www.apple.com"
```

Slangs lookup

Code

```
tweet = _slang_lookup(tweet)
```

Outcome

```
>> "I love my :3 Iphone & you are awesome apple. Display is Awesome, sooo haaaaaappy :3 http://www.apple.com"
```

Standardizing word

Code

```
tweet = " ".join(" ".join(s) for _, s in lertools.grouphy(tweet))
```

Outcome

```
>> "I love my :3 Iphone & you are awesome apple. Display is Awesome, so happy :3 http://www.apple.com"
```

Removal of URLs

URLs and hyperlinks in text data like comments, reviews, and tweets should be removed.

Final cleaned tweet:

```
>> "I love my Iphone & you are awesome apple. Display is Awesome, so happy", :3, :
```

Advanced Data Cleaning

Grammar checking

Grammar checking is majority learning based, huge amount of proper text data is learned and models are created. Many online tools are available for grammar correction purposes.

Spelling correction

In natural language, misspelled errors are encountered. One can use algorithms like the Levenshtein Distances, Dictionary Lookup etc. other modules and packages to fix these errors.

Your Next Steps...

Now that the data (tweet) is cleaned, you are ready to practice and learn the following techniques (in no order) of Text Mining:

1. Framework to build a niche dictionary for text mining
<http://bit.ly/text4web>
2. Step by Step guide to extract insights from free texts
<http://bit.ly/1j4Yv>
3. 2014 FIFA World Cup Prediction using Twitter Mining
<http://bit.ly/1kLxY5k>
4. Text Mining Hack using Google API
<http://bit.ly/1LDFF6c>



For more resources on analytics/data science, visit

www.analyticsvidhya.com

