

NLP => Spam detection

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

## \* Sms spam detection :-

\* 1. App.py (python run file)

Pickle :- serializing and de-serializing (converts obj into

render-template request :- used for generate output from a template file based on jinja2 engine  
-> ".rb" => Read binary.

\* 2. Python Pandas (for Analysis and data structure, array)  
Pickle

(1) Importing libraries, nltk, re (regular expression)

(2) load dataset

(3) download (stop words)

(4) Porter Stemmer (reducing suffix and prefix, also common words)

(5) Cleaning the message.

-> (i) Clean special characters (.,,!,@,)

(ii) Convert message into lower case

(iii) Tokenizing (split full message into units)

(iv) Remove stop words

(v) stemming the words (grew, growing = grow)

(vi) Joining the word-stemmed words

(vii) build a Corpus of message

(taking format 2-17-18)

(6) Create the model.

Learn, feature extraction (get only row data, get only useful data)

Count Vectorizer (Convert text into Vector)

max features (used for build

Vocabulary, only consider top

features according to corpus

(v. fit - transform (corpus) to array (convert text into vector)



(7) Extracting dependent variable from data set.

→ get\_dummies (give 0 or 1 to label according to labels)

→ iloc (If you want to use specific rows/columns from data, use loc and iloc)

(8) Creating a pickle file for the Count Vectorizer  
pickle.dump(cv, open('cv.pkl', 'wb'))

(9) model building :-

from sklearn.model\_selection import train\_test\_split  
test\_size = 0.20 (train = 80%, test = 20%)

(10) Naive Bayes to training set.

Import multinomial NB

Classifier = multinomialNB(alpha=0.3) naive bayes classifier

.fit <- fit model in training (Laplace smoothing) used for Predict the Probabilities of various Attributes.

(11) Create Pickle file for multinomial model.

filename = 'spbm.pkl', pickle.dump(Classifier, open(filename, 'wb'))