```
X_train = ["This was really awesome an awesome movie",
           "Great movie! Ilikes it a lot",
           "Happy Ending! Awesome Acting by hero",
           "loved it!",
           "Bad not upto the mark",
           "Could have been better",
           "really Dissapointed by the movie"]
# X_test = "it was really awesome and really disspntd"

y_train = ["positive","positive","positive","positive","negative","negative","negative"] # 1- Positive class, 0- negative class
```

```
X_train # Reviews
```

```
['This was awesome an awesome movie',
 'Great movie! Ilikes it a lot',
 'Happy Ending! Awesome Acting by hero',
 'loved it!',
 'Bad not upto the mark',
 'Could have been better',
 'Dissapointed by the movie']
```

## ⌄ Cleaning of the data

```
# Tokenize
# "I am a python dev" -> ["I", "am", "a", "python", "dev"]


from nltk.tokenize import RegexpTokenizer
# NLTK -> Tokenize -> RegexpTokenizer


# Stemming
# "Playing" -> "Play"
# "Working" -> "Work"


from nltk.stem.porter import PorterStemmer
# NLTK -> Stem -> Porter -> PorterStemmer

from nltk.corpus import stopwords
# NLTK -> Corpus -> stopwords


# Downloading the stopwords
import nltk
nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
True
```

```
tokenizer = RegexpTokenizer(r"\w+")
en_stopwords = set(stopwords.words('english'))
ps = PorterStemmer()


def getCleanedText(text):
  text = text.lower()

  # tokenizing
  tokens = tokenizer.tokenize(text)
  new_tokens = [token for token in tokens if token not in en_stopwords]
  stemmed_tokens = [ps.stem(tokens) for tokens in new_tokens]
  clean_text = " ".join(stemmed_tokens)
  return clean_text
```

## ⌄ Input from the user

```python
X_test = ["it was bad"]
```

```python
X_clean = [getCleanedText(i) for i in X_train]
xt_clean = [getCleanedText(i) for i in X_test]
```

```python
X_clean
```

```
['awesom awesom movi',
 'great movi ilik lot',
 'happi end awesom act hero',
 'love',
 'bad upto mark',
 'could better',
 'dissapoint movi']
```

```python
# Data before cleaning
'''
X_train = ["This was awesome an awesome movie",
           "Great movie! Ilikes it a lot",
           "Happy Ending! Awesome Acting by hero",
           "loved it!",
           "Bad not upto the mark",
           "Could have been better",
           "Dissapointed by the movie"]
'''
```

```
'\nX_train = ["This was awesome an awesome movie",\n           "Great movie! Ilikes it
a lot",\n           "Happy Ending! Awesome Acting by hero",\n           "loved it!",\n
"Bad not upto the mark",\n           "Could have been better",\n           "Dissapointe
```

## ˅ Vectorize

```python
from sklearn.feature_extraction.text import CountVectorizer
```

```python
cv = CountVectorizer(ngram_range = (1,2))
# "I am PyDev" -> "i am", "am Pydev"
```

```python
X_vec = cv.fit_transform(X_clean).toarray()
```

```python
X_vec
```

```
array([[0, 0, 2, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 1, 0, 0, 0],
       [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 1,
        1, 0, 0, 1, 1, 0, 0],
       [1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0,
        0, 0, 0, 0, 0, 0, 0],
       [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 1, 0, 0, 0, 0, 0],
       [0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 1, 0, 0, 1, 1],
       [0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0],
       [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 1, 0, 0, 0]])
```

```python
print(cv.get_feature_names())
```

```
['act', 'act hero', 'awesom', 'awesom act', 'awesom awesom', 'awesom movi', 'bad', 'bad upto', 'better', 'could', 'could better', 'dissa
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function get_feature_names is deprecated; get_fea
  warnings.warn(msg, category=FutureWarning)
```

```python
Xt_vect = cv.transform(xt_clean).toarray()
```

```python
Xt_vect
```

```
array([[0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0]])
```

## ⌄ Multinomial Naive Bayes

```python
from sklearn.naive_bayes import MultinomialNB

mn = MultinomialNB()

mn.fit(X_vec, y_train)
```

```
MultinomialNB()
```

```python
y_pred = mn.predict(Xt_vect)

y_pred
```

```
array(['negative'], dtype='<U8')
```