In [ ]:

```python
import numpy as np
import pandas as pd
```

In [ ]:

```python
df = pd.read_csv('../input/sms-spam-collection-data-set/SMSSpamCollection',sep='\t',names=[
```

In [ ]:

```python
df
```

# Data Pre-processing

In [ ]:

```python
df.shape
```

In [ ]:

```python
import nltk #!pip install nltk
```

In [ ]:

```python
nltk.download('stopwords')
```

In [ ]:

```python
sent = 'How are you friends?'
```

In [ ]:

```python
from nltk.tokenize import word_tokenize
word_tokenize(sent)
```

In [ ]:

```python
from nltk.corpus import stopwords
swords = stopwords.words('english')
```

In [ ]:

```python
clean = [word for word in word_tokenize(sent) if word not in swords]
```

In [ ]:

```python
clean
```

In [ ]:

```python
# Stemming words with NLTK
from nltk.stem import PorterStemmer
ps = PorterStemmer()
clean = [ps.stem(word) for word in word_tokenize(sent)
         if word not in swords]
clean
```

In [ ]:

```python
sent = 'Hello friends! How are you? We will learning python today'
```

In [ ]:

```python
def clean_text(sent):
    tokens = word_tokenize(sent)
    clean = [word for word in tokens if word.isdigit() or word.isalpha()]
    clean = [ps.stem(word) for word in clean
             if word not in swords]
    return clean
```

In [ ]:

```python
clean_text(sent)
```

In [ ]:

```python
# Pre-processing
from sklearn.feature_extraction.text import TfidfVectorizer
```

In [ ]:

```python
tfidf = TfidfVectorizer(analyzer=clean_text)
```

In [ ]:

```python
x = df['text']
y = df['label']
```

In [ ]:

```python
x_new = tfidf.fit_transform(x)
```

In [ ]:

```python
x.shape
```

In [ ]:

```python
x_new.shape
```

In [ ]:

```python
# tfidf.get_feature_names()
```

In [ ]:

```python
import seaborn as sns
sns.countplot(x=y)
```

In [ ]:

```python
#cross validation
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(x_new,y,test_size=0.25,random_state=1)
```

In [ ]:

```python
print(f"Size of splitted data")
print(f"x_train {x_train.shape}")
print(f"y_train {y_train.shape}")
print(f"y_test {x_test.shape}")
print(f"y_test {y_test.shape}")
```

In [ ]:

```python
from sklearn.naive_bayes import GaussianNB
```

In [ ]:

```python
nb = GaussianNB()
nb.fit(x_train.toarray(),y_train)
y_pred_nb = nb.predict(x_test.toarray())
```

In [ ]:

```python
y_test.value_counts()
```

In [ ]:

```python
from sklearn.metrics import ConfusionMatrixDisplay, accuracy_score
from sklearn.metrics import classification_report
import matplotlib.pyplot as plt
```

In [ ]:

```python
ConfusionMatrixDisplay.from_predictions(y_test,y_pred_nb)
plt.title('Naive bayes')
plt.show()
print(f" Accuracy is {accuracy_score(y_test,y_pred_nb)}")
print(classification_report(y_test,y_pred_nb))
```

In [ ]:

```python
from sklearn.ensemble import RandomForestClassifier
model_rf = RandomForestClassifier(random_state=1)
model_rf.fit(x_train,y_train)
```

In [ ]:

```python
y_pred_rf = model_rf.predict(x_test) #float
```

In [ ]:

```python
ConfusionMatrixDisplay.from_predictions(y_test,y_pred_rf)
plt.title('Random Forest')
plt.show()
print(f" Accuracy is {accuracy_score(y_test,y_pred_rf)}")
print(classification_report(y_test,y_pred_rf))
```

In [ ]:

```python
from sklearn.linear_model import LogisticRegression
model_lr = LogisticRegression(random_state=1)

model_lr.fit(x_train,y_train)
y_pred_lr = model_lr.predict(x_test)
```

In [ ]:

```python
ConfusionMatrixDisplay.from_predictions(y_test,y_pred_lr)
plt.title('Logistic regression')
plt.show()
print(f" Accuracy is {accuracy_score(y_test,y_pred_lr)}")
print(classification_report(y_test,y_pred_lr))
```

# Hyper parameter tunning

In [ ]:

```python
from sklearn.model_selection import GridSearchCV
```

In [ ]:

```python
para = {

    'criterion':['gini', 'entropy','log_loss'],
#     'max_features': ['sqrt','Log2'],
#     'random_state': [0,1,2,3,4],
    'class_weight':['balanced','balanced_subsample']
}
```

In [ ]:

```python
grid = GridSearchCV(model_rf, param_grid=para, cv=5, scoring='accuracy')
```

In [ ]:

```python
grid.fit(x_train,y_train)
```

In [ ]:

```python
rf = grid.best_estimator_
```

In [ ]:

```python
y_pred_grid = rf.predict(x_test)
```

In [ ]:

```python
ConfusionMatrixDisplay.from_predictions(y_test,y_pred_grid)
plt.title('Gride Search')
plt.show()
print(f" Accuracy is {accuracy_score(y_test,y_pred_grid)}")
print(classification_report(y_test,y_pred_grid))
```

In [ ]: