# Language Detection using Machine Learning ¶

In [1]:
```python
import pandas as pd
import numpy as np
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
```

In [2]:
```python
data = pd.read_csv("dataset.csv")
print(data.head())
```

```
                                                Text   language
0  klement gottwaldi surnukeha palsameeriti ning ...   Estonian
1  sebes joseph pereira thomas  på eng the jesuit...    Swedish
2  ถนนเจริญกรุง อักษรโรมัน thanon charoen krung ...       Thai
3  விசாகப்பட்டினம் தமிழ்ச்சங்கத்தை இந்துப் பத்திர...      Tamil
4  de spons behoort tot het geslacht haliclona en...      Dutch
```

In [3]:
```python
data.isnull().sum()
```

Out[3]:
```
Text        0
language    0
dtype: int64
```

In [4]:
```python
data["language"].value_counts()
```

Out[4]:
```
Estonian      1000
Swedish       1000
English       1000
Russian       1000
Romanian      1000
Persian       1000
Pushto        1000
Spanish       1000
Hindi         1000
Korean        1000
Chinese       1000
French        1000
Portugese     1000
Indonesian    1000
Urdu          1000
Latin         1000
Turkish       1000
Japanese      1000
Dutch         1000
Tamil         1000
Thai          1000
Arabic        1000
Name: language, dtype: int64
```

## Language Detection Model

```
In [5]: x = np.array(data["Text"])
        y = np.array(data["language"])

        cv = CountVectorizer()
        X = cv.fit_transform(x)
        X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.33,
                                                    random_state=42)
```

As this is a problem of multiclass classification, so I will be using the Multinomial Naïve Bayes algorithm to train the language detection model as this algorithm always performs very well on the problems based on multiclass classification:

```
In [6]: model = MultinomialNB()
        model.fit(X_train,y_train)
        model.score(X_test,y_test)
```

```
Out[6]: 0.953168044077135
```

```
In [17]: user = input("Enter a Text: ")
         data = cv.transform([user]).toarray()
         output = model.predict(data)
         print(output)
```

Enter a Text: صبح بخير
['Urdu']

```
In [ ]:
```