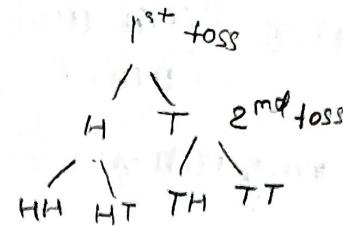
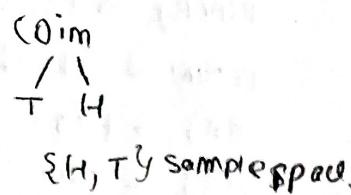


## NAIVE BAYES CLASSIFIER.

Based on  
bayes theorem.



$$P(\text{at least } 1H) = 3/4$$

$$P(2H) = 1/4$$

$$\text{conditional Prob } P(H/T) = 1/4$$

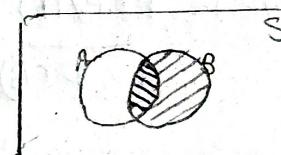
Prob of H if T is  
already occurred.

Independent events →

A & B are independent

$$\text{if } P(A \cap B) = P(A) * P(B) \quad \text{--- (1)}$$

$$P(A/B) = \frac{m(A \cap B)}{m(B)}$$



$$P(A) = \frac{P(A \cap B)}{P(B)} \quad \{ \text{from (1)} \}$$

$$= \frac{\frac{m(A \cap B)}{m(S)}}{\frac{m(B)}{m(S)}} = \frac{m(A \cap B)}{m(B)}$$

$\therefore P(A/B) = P(A)$   
as occurrence of B does  
not affect the occurrence  
of A.

Thus A & B are independent.

Mutually exclusive events →

A & B are mutually exclusive

$$\text{if } P(A \cap B) = 0$$

$$P(A/B) = \frac{P(A \cap B)}{P(B)} = 0$$

$$P(A/B) = 0$$

BAYE'S THEOREM  $\rightarrow$

If there are two events A & B

$$P(A/B) = \frac{P(B/A) P(A)}{P(B)}$$

given,  $P(B) \neq 0$

$P(A/B) \rightarrow$  Posterior  
 $P(B/A) \rightarrow$  Likelihood  
 $P(A)$  → Prior  
 $P(B)$  → Evidence

Proof -

From conditional P

$$P(A/B) = \frac{P(A \cap B)}{P(B)} \quad \text{--- ①}$$

$$A \cap B = B \cap A$$



$$P(B/A) = \frac{P(B \cap A)}{P(A)} = \frac{P(A \cap B)}{P(A)} \quad \text{--- ②}$$

$$\boxed{P(A/B) = \frac{P(B/A) P(A)}{P(B)}}$$

Ex  $\rightarrow$   
if factory manufacture 100 Markers

20      30      50  
made by      M<sub>1</sub>      M<sub>2</sub>      M<sub>3</sub>

M<sub>1</sub>  
produce      3%      1%.

5%  
defective

Given -

$$P(M_1) = \frac{1}{5} \quad P(M_2) = \frac{3}{10} \quad P(M_3) = \frac{1}{2}$$

$$P(D/M_1) = \frac{1}{20} \quad P(D/M_2) = \frac{3}{100} \quad P(D/M_3) = \frac{1}{100}$$

$$\text{To find: } P(M_3/D) = \frac{P(D/M_3) P(M_3)}{P(D)}$$

Prob that  
Marker is made  
by M<sub>3</sub> if it is  
defective.

Prob of  
Marker to  
be defective

$$P(D) = P(D \cap M_1) + P(D \cap M_2) + P(D \cap M_3)$$

$$\boxed{P(D) = P(D/M_1) P(M_1) + P(D/M_2) P(M_2) + P(D/M_3) P(M_3)}$$

## Naive Baye's theorem -

(SK. {lost, Mumbai, sunny})

toss	venue	outlook	Result
won	Mumbai	overcast	won
lost	Chennai	sunny	won
won	Kolkata	sunny	won
won	Chennai	sunny	won
lost	Mumbai	sunny	lost
won	Chennai	overcast	lost
won	Kolkata	overcast	lost
won	Mumbai	sunny	won

$$P(W \mid \text{lost} \cap \text{Mumbai} \cap \text{sunny})$$

$$= P(\text{Lost}, \text{Mumbai}, \text{sunny} \mid W) P(W)$$

$$P(\text{lost, Mumbai, sunny})$$

$$P(\text{lost} \mid \text{lost, Mumbai, sunny})$$

$$= P(\text{lost, Mumbai, sunny} \mid L) P(L)$$

$$P(\text{lost, Mumbai, sunny})$$

denominator  
is same  
so we remove  
the denominator

$$P(W) = \frac{5}{8}$$

$$P(L) = \frac{3}{8}$$

$$P(W \mid \text{lost} \cap \text{Mumbai} \cap \text{sunny}) = P(\text{lost, Mumbai, sunny} \mid W) P(W)$$

↓ & replace

$$P(\text{lost} \mid W) P(\text{Mumbai} \mid W) P(\text{sunny} \mid W) P(W)$$

$$\frac{1}{8} * \frac{2}{8} * \frac{4}{8} * \frac{5}{8}$$

## Numerical Data in Naive Baye's -

$$H = 185, W = 170, G = ?$$

height	weight	Gender
172	150	M
180	170	M
165	140	M
190	200	M
139	100	F
145	120	F
160	140	F
172	150	F

$$P(M \mid H=185, W=170) =$$

$$P(H=185 \mid M) P(W=170 \mid M) P(M)$$

If we assume that  
the height is gaussian distributed  
random variable

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

or we can assume

Binomial Distribution

Multinomial Distribution

Poisson Distribution

IMDB Dataset

50,000

## Sentiment Analysis

Review	Sentiment
Great movie	Positive
Boring	Positive
Movie.	Positive
.	.
.	.

① Text Preprocessing

② Vectorize (Bag of words)

③ Putting data into ALGO

↳ Train

④ Deployment

Bag of words →

For example we have only 4 unique words in review

"great" "awesom" "poor" "boring" label

no of occurrences for each review upto 50,000	3	1	0	0	1
	0	1	0	4	0

→ In this way we use the most common words

→ This is done by SKLearn library

↳ CountVectorizer

→ How Naive Bayes work -

Given movie was great, but slightly boring → {0, 1, 1} label = ?

awesom great boring label

3	1	0	Positive
0	0	1	Positive
1	2	0	Positive
.	.	.	
.	.	.	

$$P(\text{label} = \text{+ve} / f_1 = 0, f_2 = 1, f_3 = 1)$$

$$= \frac{P(f_1 = 0 / \text{+ve}) P(f_2 = 1 / \text{+ve}) P(f_3 = 1 / \text{+ve})}{P(\text{+ve})}$$

$$P(\text{label} = \text{-ve} / f_1 = 0, f_2 = 1, f_3 = 0)$$

$$= \frac{P(f_1 = 0 / \text{-ve}) P(f_2 = 1 / \text{-ve}) P(f_3 = 0 / \text{-ve})}{P(\text{-ve})}$$

DATASET

## Text cleaning →

- ① Sample 1000 rows
- ② Remove HTML tags
- ③ Remove special tags
- ④ Converting every string to lower case
- ⑤ Removing stop words { is, are, the... }
- ⑥ Stemming { play  
playing ⇒ play }  
played

## Advantages of NB classifiers:

- easy to implement and computationally efficient
- effective in cases with a large number of features
- performs well even with limited data.
- performs well in the presence of categorical features.

## Disadvantages of NB classifiers:

- Assumes that features are independent, which may not always hold in real world data.
- can be influenced by irrelevant attributes
- May assign zero probability to unseen events, leading to poor generalization.

## Applications of NB classifier:

- spam email filtering
- Text classification ← sentiment analysis
- Weather Prediction
- Credit scoring
- Medical Diagnosis

# Sentiment Analysis code for IMDB Dataset

```
import numpy as np } importing libraries  
import pandas as pd
```

```
df = pd.read_csv('...') } loading data
```

```
df.head()
```

```
df['review'][0]
```

```
df.info()
```

```
df['sentiment'].replace({'positive':  
    1, 'negative': 0}, inplace=True)
```

## TEXT CLEANING

- Sample 10000 rows
- Remove html tags
- Remove special characters
- Converting everything to lower
- Removing stop words
- Stemming

```
import re  
clear = re.compile('<.*?>')  
clear = re.compile('<.*?>')  
df['review'] = df.iloc[2].review
```

Removing html tags

```
def clear_html(text):  
    clear = re.compile('<.*?>')  
    return re.sub(clear, '', text)
```

```
df['review'] = df['review'].apply(clear_html)
```

```
def convert_lower(text):  
    return text.lower()
```

Converting text to lower

```
def remove_special(text):  
    x = ''  
    for i in text:  
        if i.isalnum():  
            x = x + i  
        else:  
            x = x + ' '
```

Removing special characters.

```
df['review'] = df['review'].apply(remove_special)
```

```

import nltk
from nltk.corpus import stopwords
stopwords.words('english')

def remove_stopwords(text):
    x = []
    for i in text.split():
        if i not in stopwords.words('english'):
            x.append(i)

    y = x[:]
    x.clear()
    return y

```

Removing  
stop words.

```
df['review'] = df['review'].apply(remove_stopwords)
```

```

from nltk.stem.porter import PorterStemmer
ps = PorterStemmer()

```

```

y = []
def stem_words(text):
    for i in text:
        y.append(ps.stem(i))

z = y[:]
y.clear()
return z

```

Stemming.

```

def join_back(list_input):
    return " ".join(list_input)

```

```
df['review'] = df['review'].apply(join_back)
```

```

from sklearn.feature_extraction.text import CountVectorizer
cv = CountVectorizer(max_features=1000)

```

```
X = cv.fit_transform(df['review'])
```

vectorize  
Bag of  
words

```
y = df.iloc[:, -1].values
```

```
from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```
from sklearn.naive_bayes import GaussianNB, MultinomialNB, BernoulliNB  
clf1 = GaussianNB()  
clf2 = MultinomialNB()  
clf3 = BernoulliNB()
```

```
clf1.fit(X_train, y_train)  
clf2.fit(X_train, y_train)  
clf3.fit(X_train, y_train)
```

```
y_pred1 = clf1.predict(X_test)  
y_pred2 = clf2.predict(X_test)  
y_pred3 = clf3.predict(X_test)
```

} Predicting

```
from sklearn.metrics import accuracy_score
```

```
print("Gaussian", accuracy_score(y_test, y_pred1))  
print("Multinomial", accuracy_score(y_test, y_pred2))  
print("Bernoulli", accuracy_score(y_test, y_pred3))
```

} Accuracy

Score