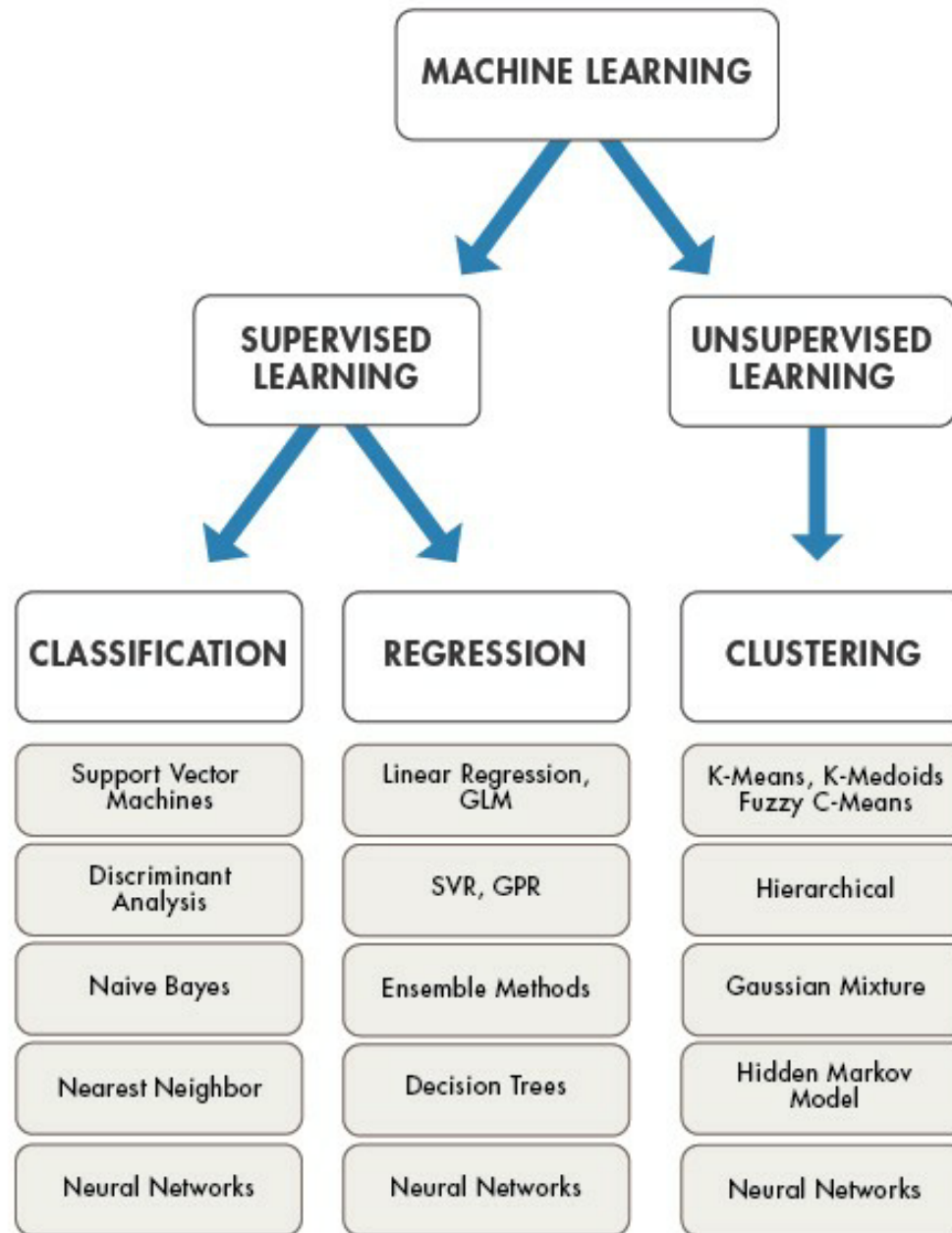


# Text Analysis

**AMAZON FINE  
FOOD REVIEWS**

amazon.com





# Attribute Information



- ID
- Product id
- User id
- Profile Name
- Helpfulness Numerator
- Helpfulness Denominator
- Score - 1 to 5
- Time
- summary
- Text

# CONTENTS



1.ABSTRACT



2.EXISTING  
SYSTEM



3.PROPOSED  
SYSTEM



4. SCOPE OF  
DEVELOPMENT



5.  
REQUIREMENTS

# ABSTRACT

- This is an Machine Learning project which can be used in E-commerce websites. The development of the Internet changed the way people eat and responding for food. Amazon is a biggest website where users can easily purchase all kind of food they need with only a mouse-click. Here our aim is to create a classifier that classifies the reviews in to either positive or negative based on the data given such as the time at which the review was written, rating given by the customer. We can analyze the using various models like Naive Bayes, KNN, Logistic Regression etc.

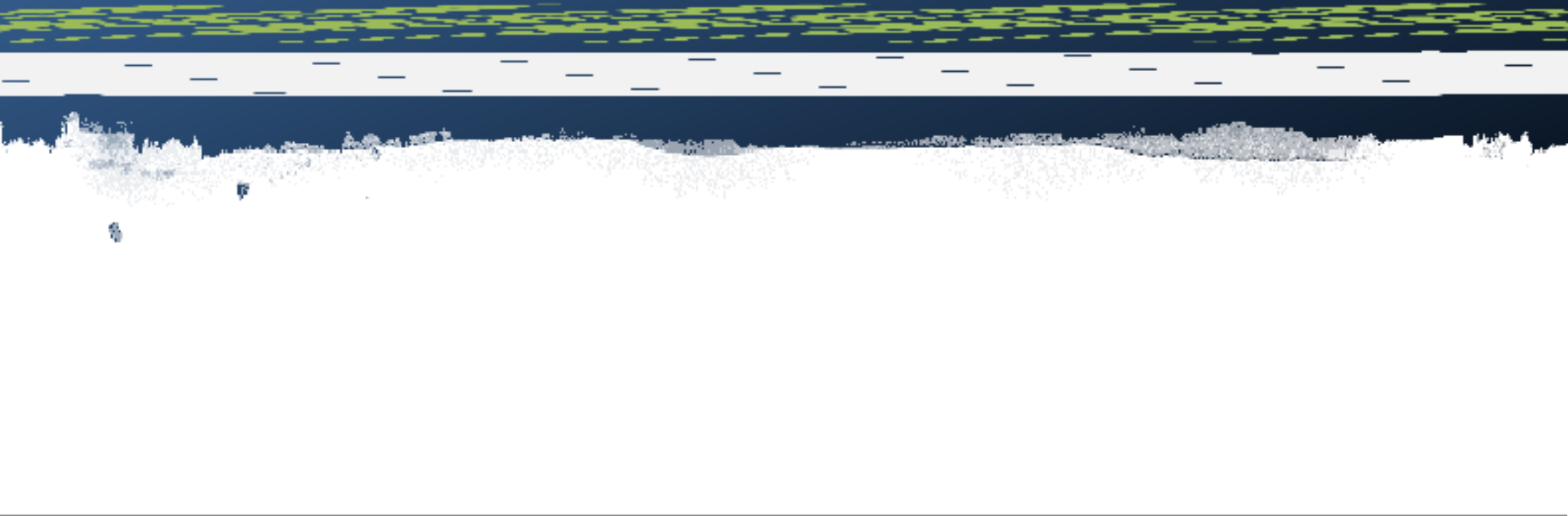


Eventually we want to build an acceptable model which helps us better understand how customers rate and review the food they purchased.

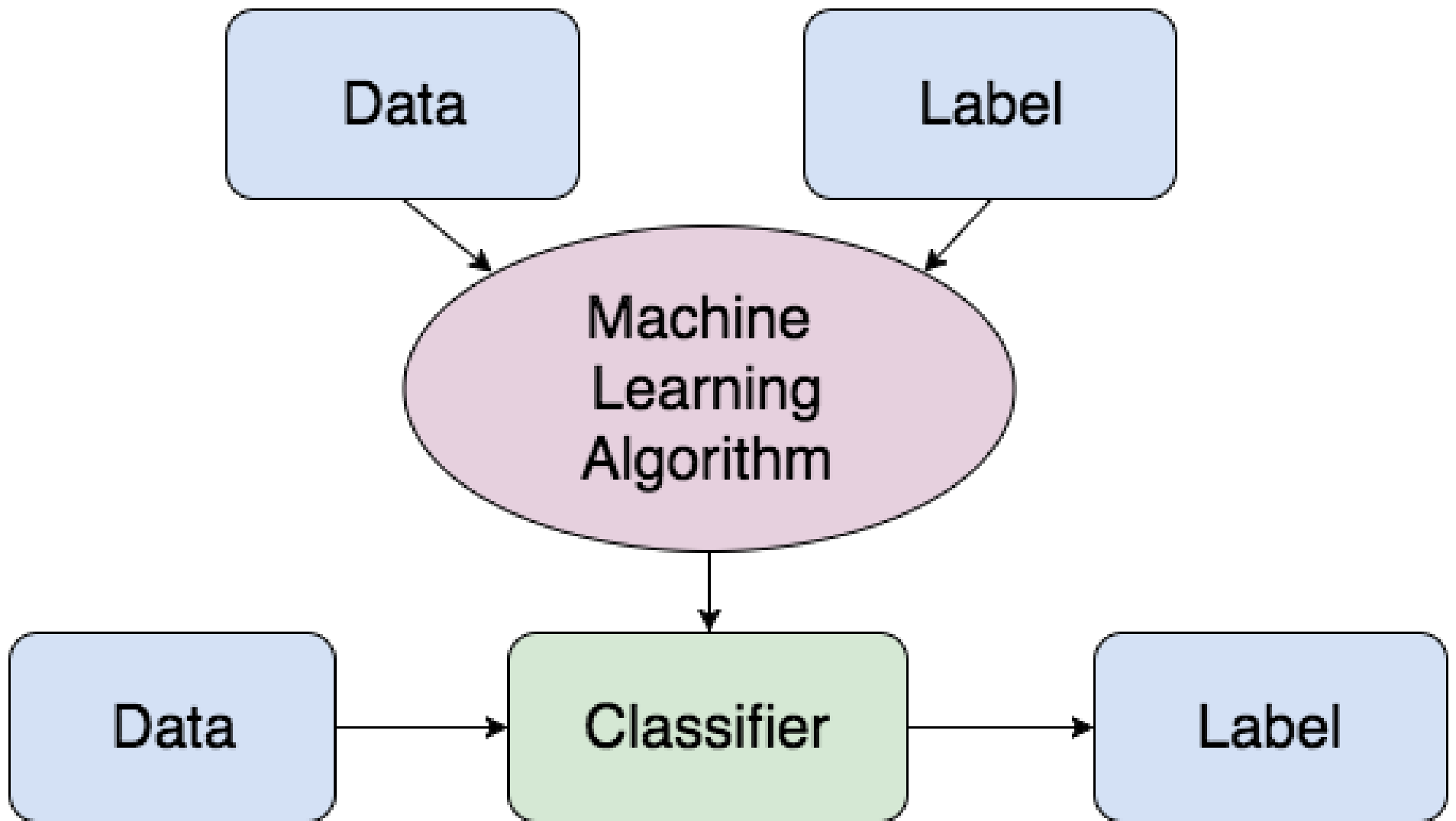
Here, we classified the reviews into 2 categories

- 1) Positive: which indicates the customer is satisfied with the product
- 2) Negative: which indicates the customer is not satisfied with the product

The main propaganda of this project is, it can be used in case of e-commerce websites where the review plays a vital role based on which the interest of the new customers depends.

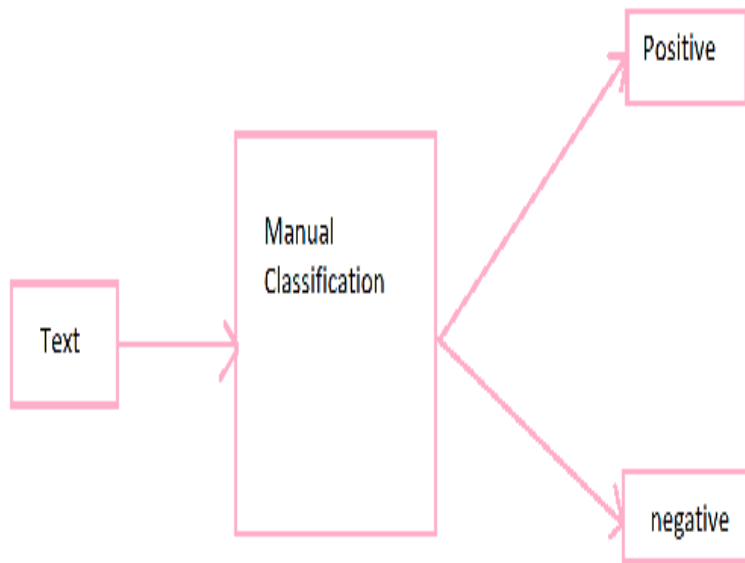






# EXISTING SYSTEM

---



In early days  
review are  
classified  
manually but

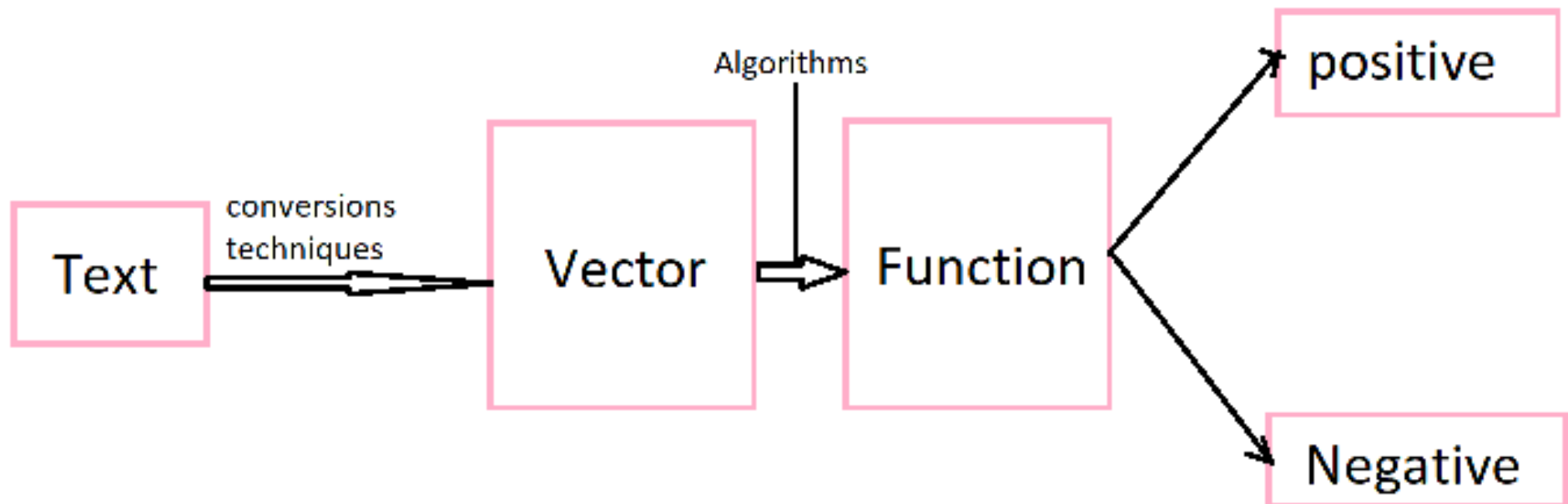
at present the as  
review has been  
increased



exponentially, it is  
impossible to do  
manually

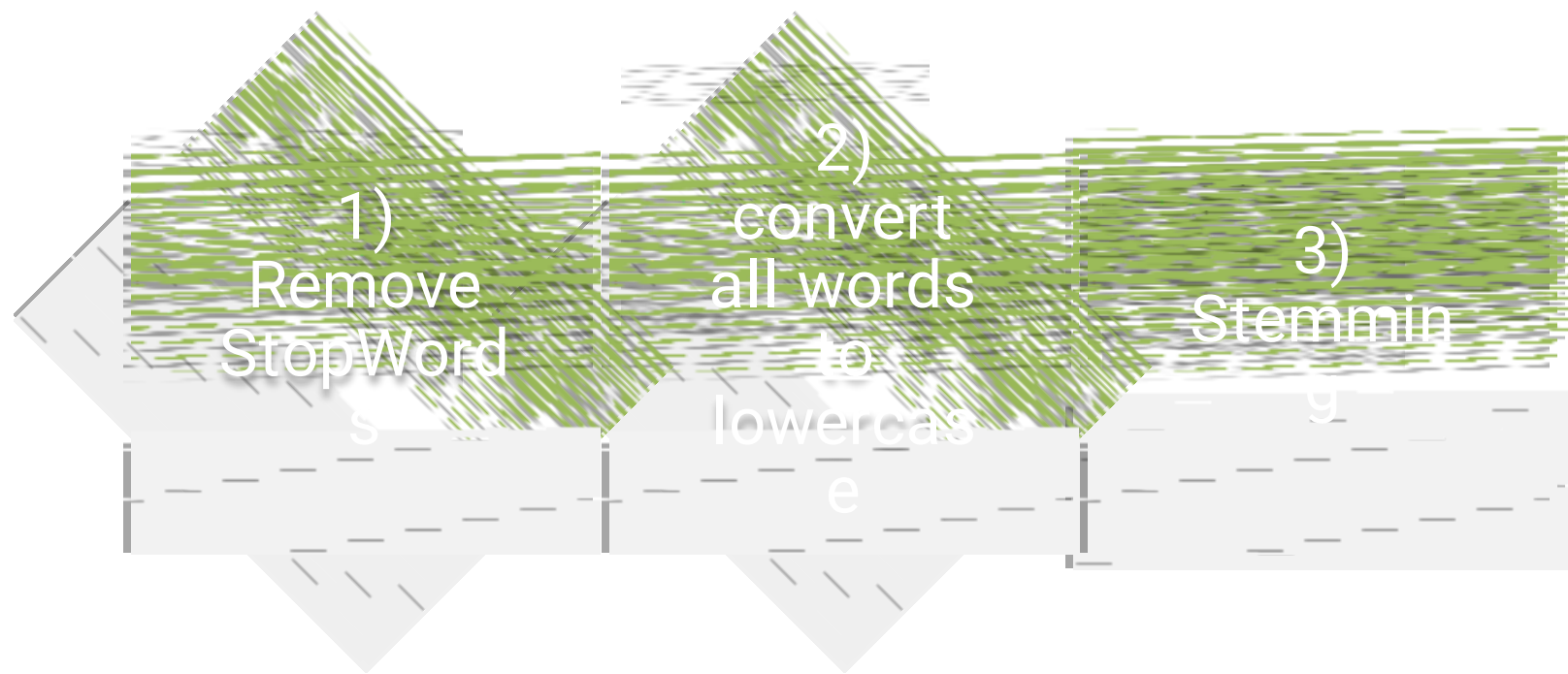
so automation is  
required.

# PROPOSED SYSTEM:



# Text Pre\_Processing

---



# Text to Vector techniques



1. BOW(Bag Of Words)



2. TF-IDF(Term frequency  
Inverse Document  
Frequency)



3. W2V (word to vector)

# Bag Of Words

- T1: the dog is on the table
- T2: now the cats are on the table

the dog is on the table

Limitation example

- T3: This is a jntu college
- T4: This is not a jntu college

0	0	1	1	0	1	1	1
are	cat	dog	is	now	on	table	the



## TF-IDF

$$w_{i,j} = tf_{i,j} \times \log \left( \frac{N}{df_i} \right)$$

$tf_{ij}$  = number of occurrences of  $i$  in  $j$

$df_i$  = number of documents containing  $i$

$N$  = total number of documents

# Formulas

1)  $TF(W,R) = \frac{\text{No of times } w \text{ occur in } R}{\text{Total no of word in } R}$

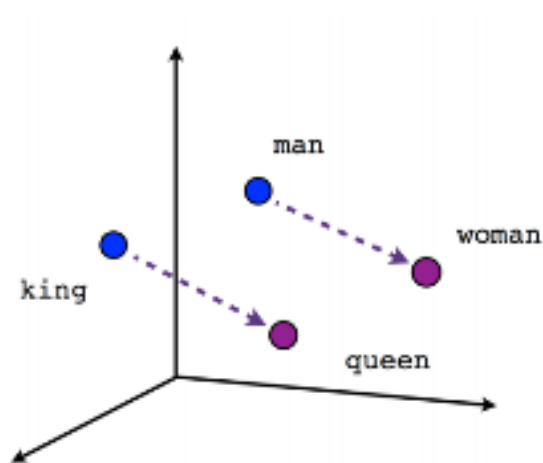
$0 \leq TF(W,R) \leq 1$

2)  $IDF(W,D) = \log(N/n_i)$

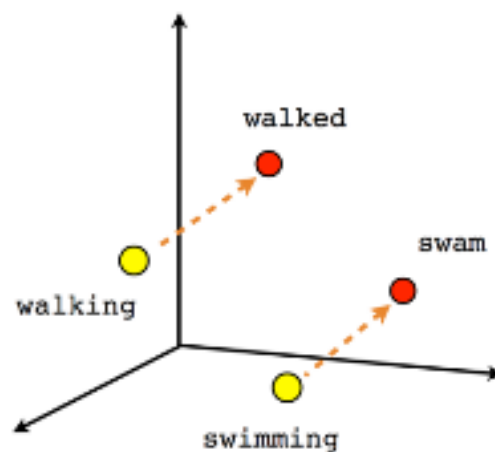
$N \rightarrow$  Total no of Documents

$n_i \rightarrow$  No of documents contain  $w$

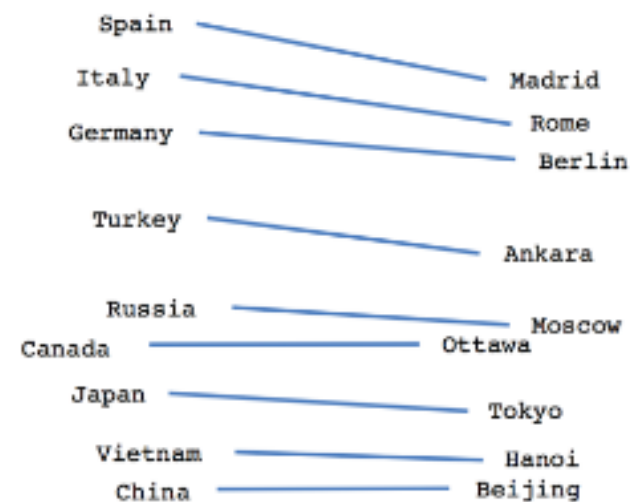
# W2V(Word to Vector)



Male-Female



Verb tense



Country-Capital

# Algorithms



KNN(k\_nearest-neighbors)

Logistic Regression

Decision Tree

Random Forest

Naive Bayes

# Parameter Tuning

- GridSearchCV
- RandomizedSearch

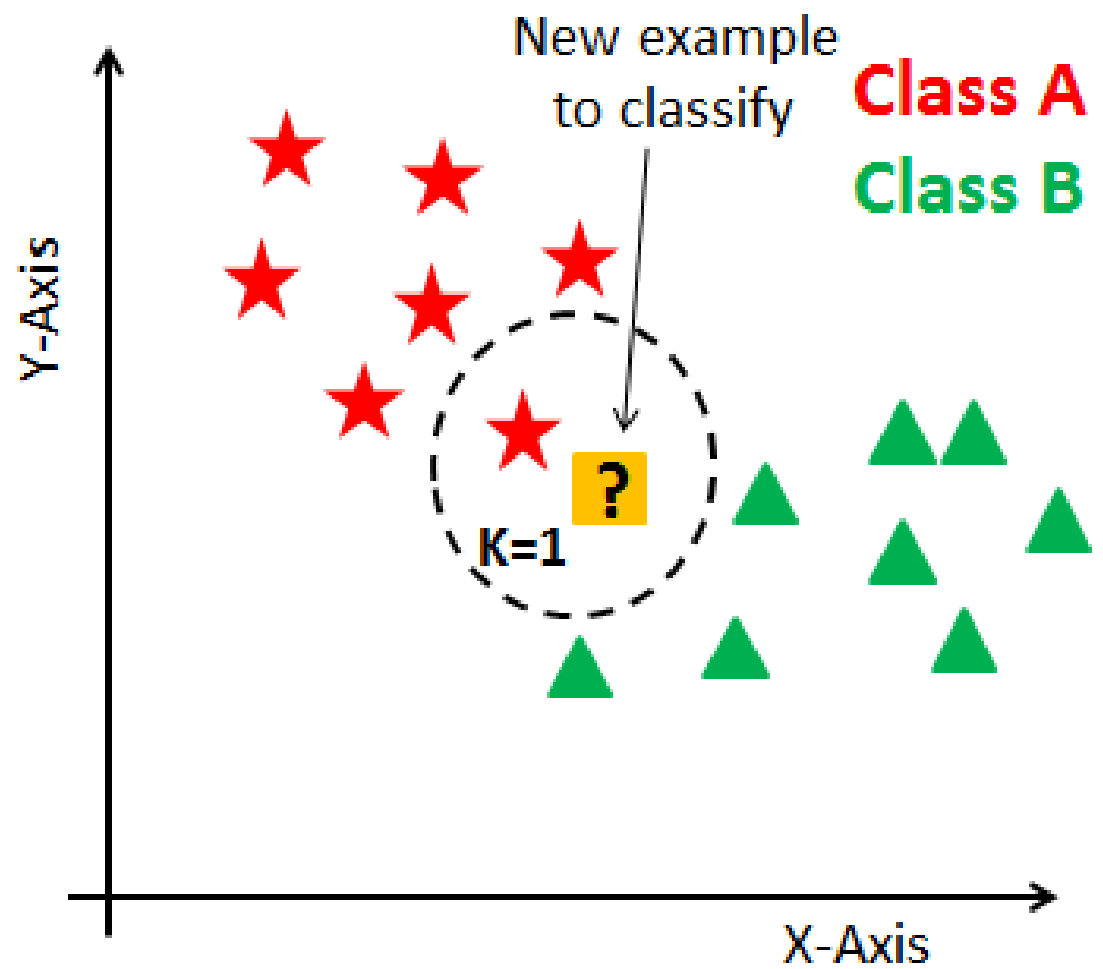
**fit KNN model to data after w2v**

```
In [69]: from sklearn.model_selection import GridSearchCV
          from sklearn.neighbors import KNeighborsClassifier
          parameters = {'n_neighbors':[1,2,3,4,5,6,7,8]}
          neigh = KNeighborsClassifier()
          clf = GridSearchCV(neigh, parameters, cv=5, scoring="accuracy")
          clf.fit(x2_train,y2_train)
          clf.best_params_
```

```
Out[69]: {'n_neighbors': 7}
```

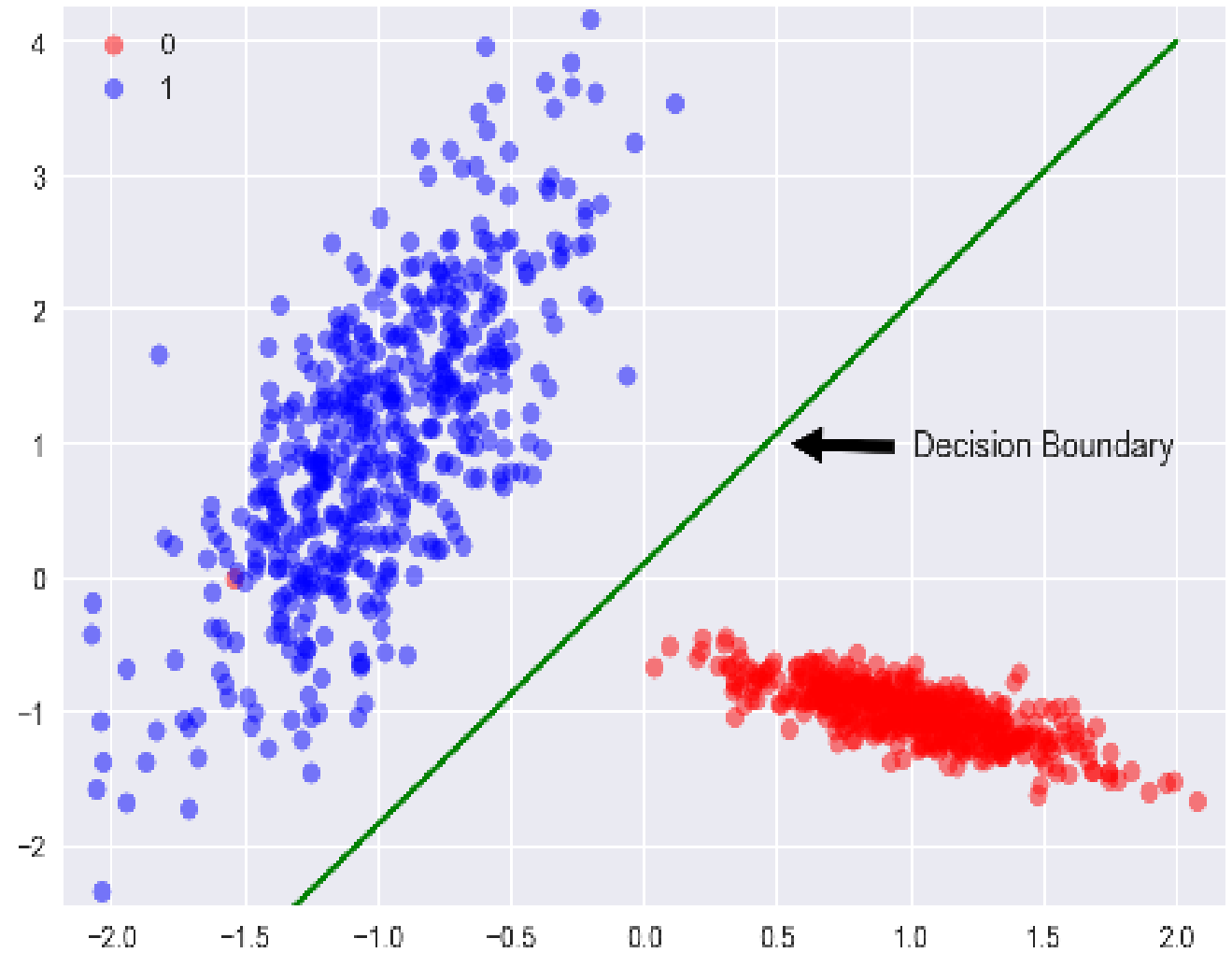
---

KNN

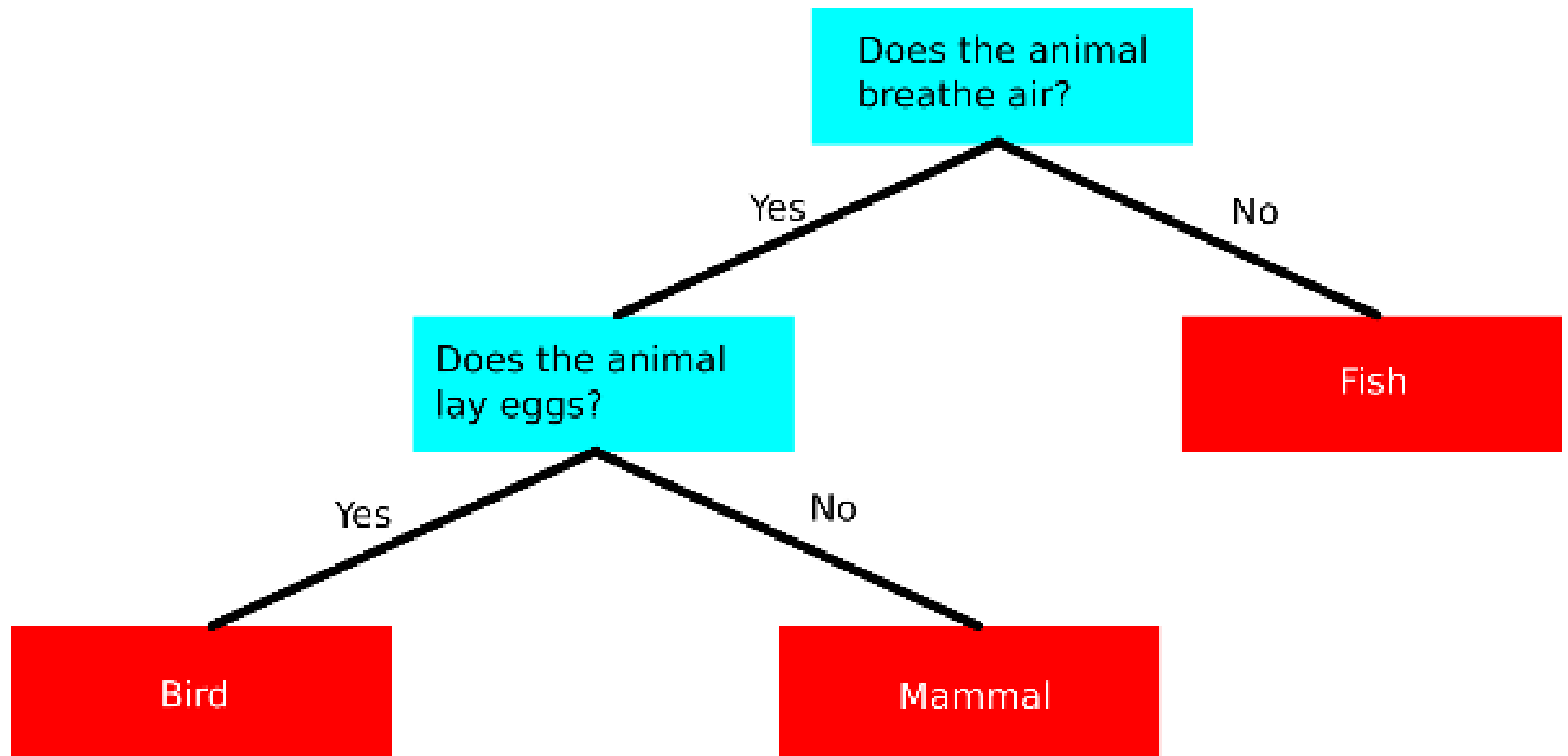




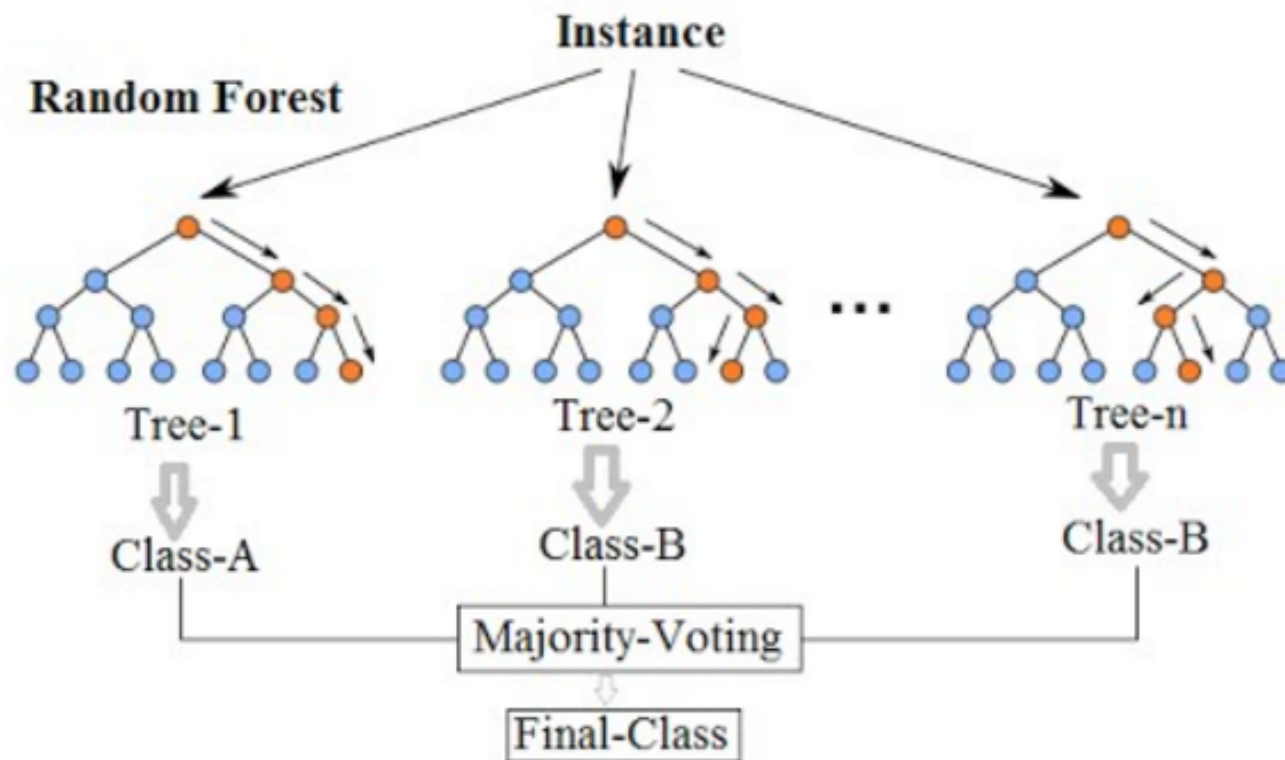
# LogisticRegression



# Decision Tree



## Random Forest Simplified



# Naïve Bayes

The diagram shows the Naïve Bayes formula with four labels and arrows pointing to the corresponding parts of the equation:

- Likelihood** points to  $P(x|c)$
- Class Prior Probability** points to  $P(c)$
- Posterior Probability** points to  $P(c|x)$
- Predictor Prior Probability** points to  $P(x)$

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \cdots \times P(x_n|c) \times P(c)$$

# Performance Metrics

- ACCURACY
- F1\_SCORE

		Actual	
		Positives(1)	Negatives(0)
Predicted	Positives(1)	TP	FP
	Negatives(0)	FN	TN

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

$$F_1 = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

# Results

## Model Evaluation and Validation

Algorithms	BOW	TFIDF
W2V		
naive Bayes	80(97)	88(100)
57(50)		
KNN	82(90)	90(92)
69(74)		
Logistic regression	90(100)	88(100)
57(85)		
Decision Tree	83(85)	83(92)
71(73)		
Random Forest	90(94)	88(95)
70(90)		



# Scope of Development



We may use some advanced machine learning techniques to improve the performance , like



Deep Learning

## REQUIRMENTS

OS: Windows,Mac OS,unix,linux,  
etc.

RAM-1)Minimum(8gb)

2)Recommended(16gb)

PROCESSOR: 1)Minimum(i5)

2)Recommended(i7)

ROM:1)Minimum(512gb)

2)Recommended(1 TB)