



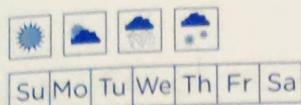
Su Mo Tu We Th Fr Sa

Date: / /

## Appendix

### Table of Contents

<u>Topic</u>	<u>Page</u>	<u>Date</u>
i) Title for research	— 01	— 02/07/2021
ii) Research question	— 02	— 03/07/2021
iii) Regression	— 03	— 05/07/2021
iv) Feature Extraction	— 04	— 06/07/2021
v) Naive Bayes	— 05	— 08/07/2021
vi) Neural Net	— 06/07	— 09/07/2021
vii) AND & XOR&OR	— 08	— 11/07/2021
viii) Decision Boundaries	— 09	— 12/07/2021
ix) Multi layered Perceptron	— 10	— 13/07/2021
X) Bag of Words	— 11	— 14/07/2021
xii) Max Likelihood Estimate	— 12	— 16/07/2021
xii) Random Forest	— 12	— 16/07/2021
xiii) SGD Classifier	— 13	— 17/07/2021
xiv) Model	— 14	— 19/07/2021
xv) Preprocessing	— 14	— 19/07/2021
xvi) Dataset Analyzation	— 15	— 20/07/2021
xvii) Strategy to Label	— 16	— 21/07/2021
xviii) Final Results	— 17	— 22/07/2021



Date: 02/07/2021

Finding a suitable title for research  
Title: proposal.

- i) Assessing the performances of various algorithms for sentiment analysis.
- ii) Comparing different sentiment analysis algorithms for optimal performance.
- iii) Implementing different sentiment analysis algorithms to find the best one.
- iv) Assessing product reviews.
- v) Assessing customer reviews to better
- vi) Assessing sentiment analysis algorithms to harness online product reviews.



Su Mo Tu We Th Fr Sa

Date: 03/07/2021

## Research question:

With respect to title (iv),

- i) Which methodology is the most effective to analyze the sentiments from product reviews?
- ii) Which algorithms is the most efficient in terms of analyzing sentiments from online product reviews? when
- iii) Which algorithms provide the best accuracies when it comes to analyzing sentiment from online reviews?
- iv) Which algorithm is the most accurate to better analyze sentiment from online product



Su Mo Tu We Th Fr Sa

Date: 05 / 07 / 2021

Linear regression → Linear regression is for handling regression, continuous output

Logistic regression → Logistic regression is for handling classification, discrete output

$$z = wx + b$$

$$\hat{y} = \sigma(wx + b)$$

$$= \sigma(z)$$

$$= \frac{1}{1 + e^{-z}}$$

$$\frac{\partial L_{CE}(\hat{y}, y)}{\partial w_j} = [\sigma(wx + b) - \hat{y}]x_j$$

$$\theta_{t+1} = \theta_t - \eta x^j$$



Su Mo Tu We Th Fr Sa

Date: 06/07/2021

## Bag of Words:

Step i) Other preprocessing

Step ii) Create vocabulary from unique words

Step iii) Create a matrix, also known as text vectorization

\* Drawback is that the order of occurrences is lost

\* This is why create bigrams.

TF-IDF (Term Frequency - Inverse Document Frequency)

TF → Specifies how frequently a term appears.  
— Calculates how many times it appears, with respect to total number of words.

IDF →  
— whether a term is rare or frequent.  
— Rare words have high IDF scores.  
— It's a log-normalized value



Su Mo Tu We Th Fr Sa

Date: 08 / 07 / 2021

Naive Bayes:

Let document  $d \in$

$$\text{classes } C = \{c_1, c_2, c_3, \dots\}$$

in labeled documents. ( $d_{C_1}, d_{C_2}, d_{C_3}, \dots$ )

Output  $\hat{h}: d \rightarrow C$

$$\text{Bow: } \hat{h}(?) = \hat{c}$$

$$\text{Bow: } h(d) = \hat{c}$$

For doc  $d$ , class  $C$ ,

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

$$c_{MAP} = \arg \max P(c|d)$$
$$= \arg \max \frac{P(d|c) P(c)}{P(d)}$$

$$= \arg \max P(d|c) P(c)$$

$$= \arg \max P(x_1, x_2, x_3, \dots, x_n | c) P(c)$$

$$d = \{$$
  
$$\text{neutral : 3}$$
  
$$\text{happy : 6}$$
  
$$\text{sad : 3}$$
  
$$\text{surprised : 2}\}$$

05



SANOFI

Empowering Life



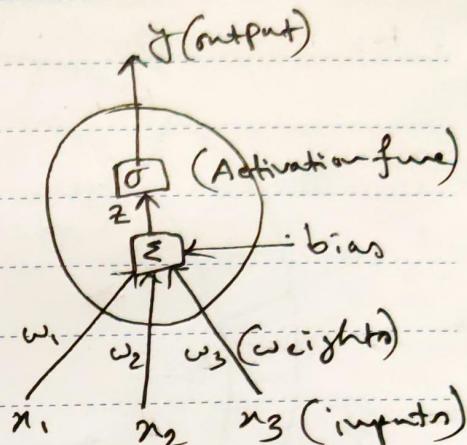
Su Mo Tu We Th Fr Sa

Date: 09/07/2021

## Neural net:

$$z = b + \sum w_i x_i \\ = w \cdot n + b$$

$$y = a = f(z)$$



$$y = \sigma(z) = \frac{1}{1+e^{-z}}$$

Other activation functions:

- tanh
- softmax
- ReLU

\* Sigmoid makes the value range from 0 to 1.

$$* \tanh(y) = \frac{e^y - e^{-y}}{e^y + e^{-y}}$$

tanh makes the value range from -1 to 1

06



Su Mo Tu We Th Fr Sa

Date: 09/07/2021

\* ReLU  $\rightarrow z$  when  $z$  is positive, 0 otherwise.  
Almost linear..

$$\text{ReLU}(y) = \max(z, 0)$$

### Perceptron:

Binary output (0 or 1)

- non linear activation function

$$- f = \begin{cases} 0, & \text{if } w \cdot x + b \leq 0 \\ 1, & \text{if } w \cdot x + b > 0 \end{cases}$$

AND, OR with Perceptron

AND

$x_1$	$x_2$	$f$
0	0	0
0	1	0
1	0	0
1	1	1

OR

$x_1$	$x_2$	$f$
0	0	0
0	1	1
1	0	1
1	1	1

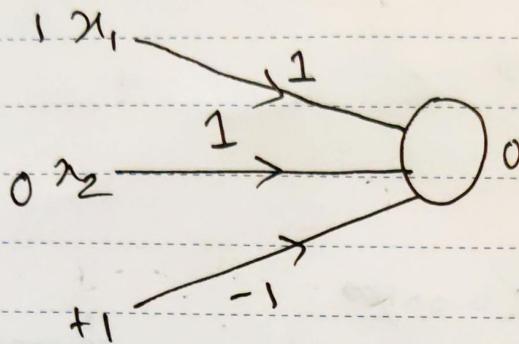
OR



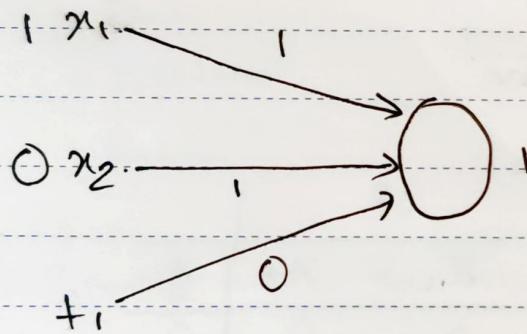
Su Mo Tu We Th Fr Sa

Date: 11 / 07 / 2021

AND



OR



\* XOR cannot be built with perception.

\* Given input  $x_1, x_2$ ,  $w_1x_1 + w_2x_2 + b = 0$

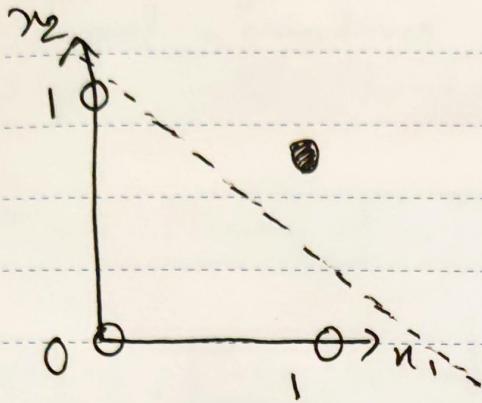
\* In standard linear format  $x_2 = (-\frac{w_1}{w_2})x_1 + (\frac{-b}{w_2})$



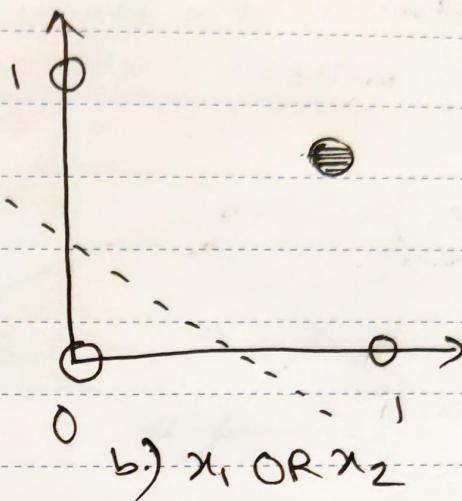
Su Mo Tu We Th Fr Sa

Date: 12 / 07 / 2021

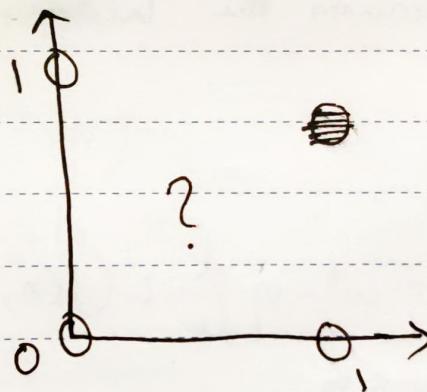
## Decision boundaries:



a)  $x_1 \text{ AND } x_2$



b)  $x_1 \text{ OR } x_2$



c)  $x_1 \text{ XOR } x_2$

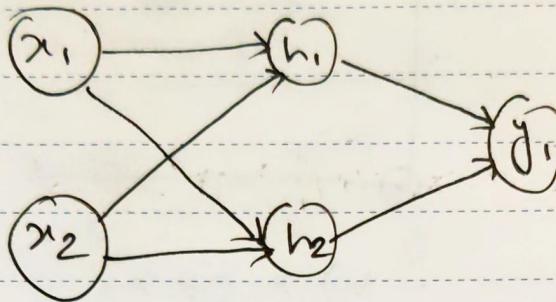
\* we can see that XOR is not linearly separable.



Su Mo Tu We Th Fr Sa

Date: 13 / 07 / 2021

\* XOR cannot be calculated with single perceptron.  
However, it is possible by combining a layer of units.



\*  $h_1$  and  $h_2$  represents the hidden layers.

Gradient descent:

$$\omega^{t+1} = \omega^t - \eta \frac{d}{d\omega} L(f(x; \omega), y)$$

For logistic regression

$$\frac{\partial L_{CE}(\hat{y}, y)}{\partial w_j} = [\sigma(wx + b) - y] x_j$$

10



Su Mo Tu We Th Fr Sa

Date: 19 / 02 / 2021

Bag of words (B.o.w) with N-grams:

$N=1$  : This is a sentence  $\rightarrow$  unigrams

$N=2$  : This is a sentence  $\rightarrow$  bigrams

$N=3$  : This is a sentence  $\rightarrow$  trigrams

\* last word in a N gram sentence given previous words.

$$P(w_n | w_{n-1}, w_{n-2} \dots, w_{n-N+1})$$

e.g.  $N=2$ ,  $P(w_n | w_{n-1})$

N-gram sentence

$$P(w_n, w_{n-1}, w_{n-2} \dots, w_{n-N+1})$$

e.g  $N=2$   $P(w_n, w_{n-1})$



Su Mo Tu We Th Fr Sa

Date: 16/02/2021

Maximum likelihood estimate,

$$P(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

E.g.:

$\langle s \rangle$  I am Anfin  $\langle /s \rangle$

$\langle s \rangle$  Anfin I am  $\langle /s \rangle$

$\langle s \rangle$  I do not like eggs and ham  $\langle /s \rangle$

$$P(I | \langle s \rangle) = 0.67, P(\text{Anfin} | \langle s \rangle) = 0.33$$

$$P(\text{am} | I) = 0.67, P(\text{Anfin} | \text{am}) = 0.5$$

$$P(\langle /s \rangle | \text{Anfin}) = 0.5, P(\text{dot} | I) = 0.33$$

Random forest:

Let,

training set  $\mathcal{X} = x_1, \dots, x_n$

" labels  $\mathcal{Y} = y_1, \dots, y_n$

For,  $b = 1, \dots, B$  (bagging repeatedly  $B$  times)

$$\hat{f} = \frac{1}{B} \sum_{b=1}^B f_b(x')$$

$$\sigma = \sqrt{\frac{\sum_{b=1}^B (f_b(x') - \hat{f})^2}{B-1}}$$





Su Mo Tu We Th Fr Sa

Date: 12/02/2021

Relationship to nearby neighbors:

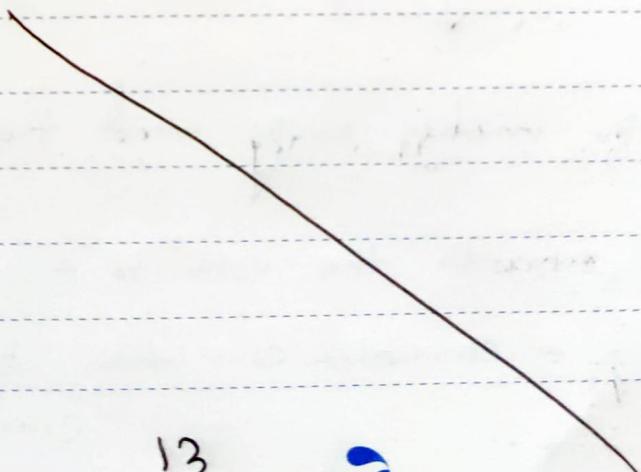
$$\hat{y} = \sum_{i=1}^n w_j(x_i, x) y_i$$

Since, a forest averages the predictions of a set of  $m$  trees, with individual weight function  $w_j$ , its predictions are:

$$\hat{y} = \frac{1}{m} \sum_{j=1}^m \sum_{i=1}^n w_j(x_i, x) y_i = \sum_{i=1}^n \left( \frac{1}{m} \sum_{j=1}^m w_j(x_i, x) \right) y_i$$

SGD: (Stochastic Gradient classifier):

- \* It's just a linear classifier that uses gradient descent on a loss function unlike logistic regression.



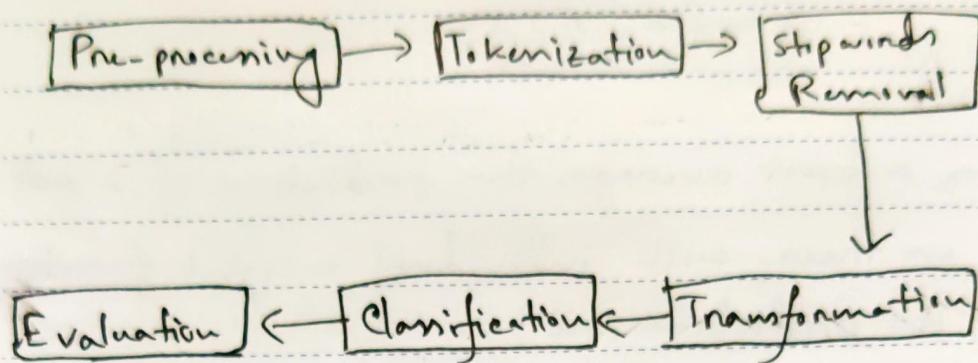
13



Su Mo Tu We Th Fr Sa

Date: 19/07/2021

Model:



① Pre-processing:

- Removing special characters, numbers, punctuations
- Tokenization
- Stemming
- Removing stop words

② Tokenization makes every word token

③ Stemming converts the verb to its original form. E.g. - connected, connection, becomes connect.



Su Mo Tu We Th Fr Sa

Date: 20/02/2021

\* Stopwords { are common English words, like, "a", "the", "are", "is"

By removing these the training data becomes less ambiguous.

~~flat~~

Dataset: Amazon reviews (Labeled data)

Column names:

- ID
- Product ID
- User ID
- Profile Name
- Helpfulness Numerator
- Helpfulness Denominator
- Score
- Time
- Summary
- Text



Su Mo Tu We Th Fr Sa

Date: 21/07/2021

\*Marking the dataset "Text" column based on Scene.

- Score  $\geq 4 \Rightarrow$  Positive
- Score = 3  $\Rightarrow$  Neutral
- Score  $\leq 2 \Rightarrow$  Negative

Visualizing label count:

Positive  $\rightarrow 443772$

Negative  $\rightarrow 82032$

Neutral  $\rightarrow 42640$

\*Splitting the dataset into training and testing.

X\_train.shape : (454763, 72127)

X\_test.shape : (113691, 72127)

y\_train.shape : (454763,)

y\_test.shape : (113691,)



Su Mo Tu We Th Fr Sa

Date: 22/07/2021

i) Naive Bayes:-

Training accuracy : 0.834 [4th]

Testing accuracy : 0.828

ii) Stochastic Gradient Classifier (SGD) :-

Training accuracy : 0.867 [3rd]

Testing accuracy : 0.863

iii) Random Forest:-

Training accuracy : 0.999 [1st]

Testing accuracy : 0.886

iv) Logistic Regression:-

Training accuracy : 0.877 [2nd]

Testing accuracy : 0.867