

The background of the slide is a scenic landscape. It features a calm body of water in the foreground, with a small boat visible in the distance. The water reflects the light from the sky. In the background, there are misty mountains and a forested shoreline on the left. Several birds are flying in the sky. The overall atmosphere is peaceful and serene.

# Natural Language Processing with Classification and Vector Spaces

*Notes*

Junru Lin

# Contents

<b>1</b>	<b>Logistic Regression</b>	<b>3</b>
1.1	Learning Objectives . . . . .	3
1.2	Supervised ML (Training) . . . . .	3
1.3	Sentiment Analysis . . . . .	4
1.4	Vocabulary . . . . .	4
1.4.1	Vocabulary and Frequency . . . . .	4
1.4.2	Feature Extraction with Frequencies . . . . .	5
1.5	Preprocessing . . . . .	6
1.6	Putting it All Together . . . . .	6
1.7	Logistic Regression (LR) Overview . . . . .	7
1.8	LR: Training . . . . .	7
1.9	LR: Testing . . . . .	7
1.10	LR: Costing Function . . . . .	8
<b>2</b>	<b>Probability and Bay's Rule</b>	<b>9</b>
2.1	Probability . . . . .	9
2.2	Bayes' Rule . . . . .	9
2.3	Naive Bayes Introduction . . . . .	10
2.4	Laplacian Smoothing . . . . .	10

<i>CONTENTS</i>	3
2.5 Log Likelihood . . . . .	11
2.6 Training Naive Bayes . . . . .	12
2.7 Testing Naive Bayes . . . . .	12
2.8 Application of Naive Bayes . . . . .	12
2.9 Naive Bayes Assumptions . . . . .	12
2.10 Error Analysis . . . . .	13
<b>3 Vector Space Models</b>	<b>14</b>
3.1 Vector Space Models . . . . .	14
3.2 Word by Word and Word by Doc . . . . .	14
3.3 Euclidean Distance and Cosine Similarity . . . . .	15
3.4 Cosine Similarity . . . . .	15
3.5 Manipulating Words in Vector Space . . . . .	16
3.6 Visualization and PCA . . . . .	16
3.7 PCA Algorithm . . . . .	16
<b>4 Machine Translation</b>	<b>17</b>
4.1 Transforming Word Vectors . . . . .	17
4.2 K-nearest Neighbors . . . . .	18
4.3 Hash Tables and Hash Functions . . . . .	18
4.4 Locality Sensitive Hashing . . . . .	18
4.5 Multiple Planes . . . . .	19
4.6 Approximate Nearest neighbors . . . . .	19

# Chapter 1

## Logistic Regression

### 1.1 Learning Objectives

- Sentiment analysis
- Logistic regression
- Data pre-processing
- Calculating word frequencies
- Feature extraction
- Vocabulary creation
- Supervised learning

### 1.2 Supervised ML (Training)

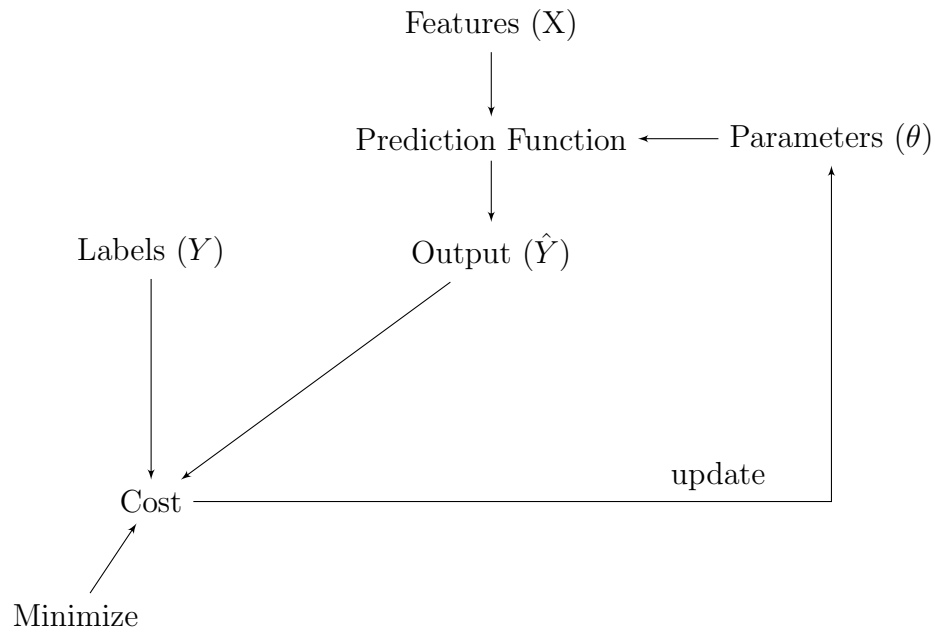


Figure 1.1: Flowchart of Supervised ML (Training)

## 1.3 Sentiment Analysis

### *Example*

I am happy because I am learning NLP.

### *Classification Label*

- Positive - 1
- Negative - 0

### *Training*

Use Logistic Regression

## 1.4 Vocabulary

### 1.4.1 Vocabulary and Frequency

#### *Vocabulary*

All unique words in a dataset.

#### *For Each data*

An array of 0's and 1's.

$|v|$  = size of vocabulary (number of features)

Logistic Regression:  $\theta = [\theta_0, \theta_1, \theta_2, \dots, \theta_n]$  ( $\theta_0 \rightsquigarrow b$  ( $y = ax + b$ ),  $n = |v|$ )

### ***Negative and Positive Frequencies***

*freqs*: dictionary mapping from [word, class] to frequency

example:  $\text{freqs}[(\text{'I'}, 1)] = 2$

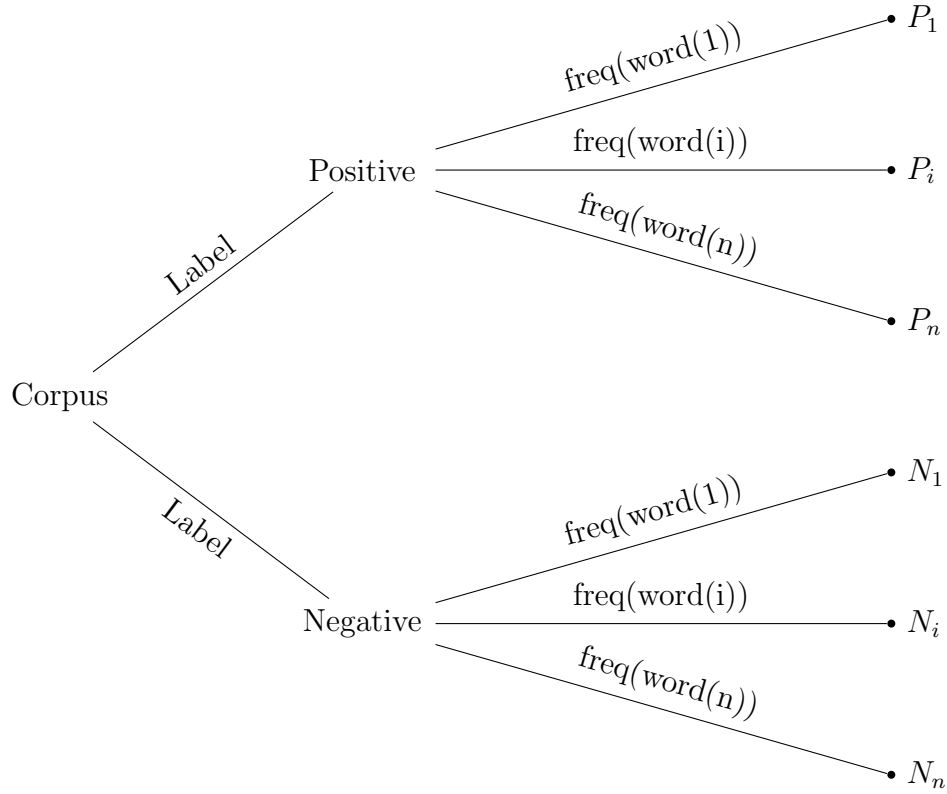


Figure 1.2: Word Frequency Tree

Vocabulary	PosFreq(1)	NegFreq(0)
word(1)	$P_1$	$N_1$
$\vdots$	$\vdots$	$\vdots$
word(n)	$P_n$	$N_n$

### **1.4.2 Feature Extraction with Frequencies**

$$X_m = \left[ 1, \sum_m \text{freqs}(w, 1), \sum_m \text{freqs}(w, 0) \right]$$

- $X_m$ : Feature of data  $m$  in the dataset

- 1: Bias
- $\sum_m freqs(w, 1)$ : Sum of Pos. Freq
- $\sum_m freqs(w, 0)$ : Sum of Neg. Freq

**Example**

I am sad.

I am not learning NLP.

$\hookrightarrow$  words (w): I, am, sad, not, learning, NLP

## 1.5 Preprocessing

**Stop Words**

and, is, a, at, has, for, of, ...

**Punctuation**

, . " ! " ' ' ...

**Handles and URLs**

@ ...

**Stemming**

tune, tuned, tuning  $\rightarrow$  tun

**lowercasing**

Great, GREAT  $\rightarrow$  great

## 1.6 Putting it All Together

$$X = \begin{bmatrix} 1 & X_1^{(1)} & X_2^{(1)} \\ 1 & X_1^{(2)} & X_2^{(2)} \\ \vdots & \vdots & \vdots \\ 1 & X_1^{(m)} & X_2^{(m)} \end{bmatrix}$$

where  $m$  is the sample size.

*Example*

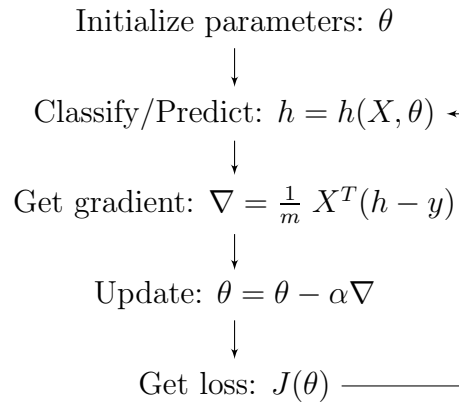
$$\begin{aligned}
 & [1, 40, 20], \\
 & [1, 20, 50], \\
 & \dots, \\
 & [1, 5, 35] ]
 \end{aligned}$$

**1.7 Logistic Regression (LR) Overview***Sigmoid Function*

$$h(x^{(i)}, \theta) = \frac{1}{1 + e^{-\theta^T x^{(i)}}} \in (0, 1)$$

$$\theta^T x^{(i)} \geq 0 : h(x^{(i)}, \theta) \geq 0.5 \rightarrow 1$$

$$\theta^T x^{(i)} < 0 : h(x^{(i)}, \theta) < 0.5 \rightarrow 0$$

**1.8 LR: Training****1.9 LR: Testing**

- Validation set:  $X_{val}, Y_{val}$
- After training:  $\theta$



- $\text{pred} = h(X_{\text{val}}, \theta) \geq 0.5$  (sigmoid function)
- $\text{Accuracy} = \sum_{i=1}^m \frac{\text{pred}^{(i)} == y_{\text{val}}^{(i)}}{m} = \frac{\sum_{i=1}^m (\text{pred}^{(i)} == y_{\text{val}}^{(i)})}{m}$
- $X_{\text{train}} : X_{\text{val}} : X_{\text{test}} = 8 : 1 : 1$

**Example**

$$\begin{bmatrix} 0.3 \\ 0.8 \\ \vdots \\ h_m \end{bmatrix} \geq 0.5 = \begin{bmatrix} 0.3 \geq 0.5 \\ 0.7 \geq 0.5 \\ \vdots \\ hm \geq 0.5 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ \text{pred}_m \end{bmatrix}$$

## 1.10 LR: Costing Function

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h(x^{(i)}, \theta) + (1 - y^{(i)}) \log(1 - h(x^{(i)}, \theta))]$$

**Want**  $\min_{\theta} J(\theta)$

# Chapter 2

## Probability and Bay's Rule

### 2.1 Probability

$A \rightarrow$  *Positive tweet*

$$P(A) = P(Positive) = N_{pos}/N$$

$$P(Negative) = 1 - P(Positive)$$

$B \rightarrow$  *tweet contains "happy"*

$$P(B) = P(happy) = N_{happy}/N$$

*Probability of the intersection*

$$P(A \cap B) = P(A, B)$$

### 2.2 Bayes' Rule

*Conditional Probability*

Probability of  $A$ , given  $B$  happened.  $\rightarrow P(A|B)$

*Bayes' Rule*

$$\begin{cases} P(A|B) = \frac{P(A \cap B)}{P(B)} \\ P(B|A) = \frac{P(A \cap B)}{P(A)} \end{cases}$$
$$\Rightarrow P(A|B) = P(B|A) \cdot \frac{P(A)}{P(B)}$$

## 2.3 Naive Bayes Introduction

### **Assumption**

$X_1, X_2, \dots, X_m$  are independent

### **By Bayes' Rule**

$$\begin{aligned}
 P(Positive|X) &= P(X|Positive) \cdot P(Positive)/P(X) \\
 &= P(X_1|Positive) \cdots P(X_m|Positive) \cdot P(Positive)/P(X) \\
 &= \frac{P(Positive)}{P(X)} \cdot \prod_{i=1}^m P(X_i|Positive) \\
 P(Negative|X) &= \frac{P(Negative)}{P(X)} \cdot \prod_{i=1}^m P(X_i|Negative)
 \end{aligned}$$

### **Assumption**

$P(Positive) = P(Negative)$

$$\begin{aligned}
 \Rightarrow \frac{P(Positive|X)}{P(Negative|X)} &= \frac{\prod_{i=1}^m P(X_i|Positive)}{\prod_{i=1}^m P(X_i|Negative)} \\
 &= \prod_{i=1}^m \frac{P(X_i|Positive)}{P(X_i|Negative)}
 \end{aligned}$$

### **Note**

- In a positive words table, we will have the probability for each word, which is  $P(X_i|Positive)$ .
- Same for  $P(X_i|Negative)$ .

## 2.4 Laplacian Smoothing

**To avoid**  $P(W_i|class) = 0$

- $N_{class}$  = frequency of all words in class
- $V$  = number of unique words in vocabulary
- $class \in \{ Positive, Negative \}$
- $P(W_i|class) = \frac{freq(w_i|class)}{N_{class}}$   
 $\downarrow$

- $P(W_i|class) = \frac{freq(w_i|class) + 1}{N_{class} + V}$

$$freq(w_i|class) + 1 \rightsquigarrow \text{make it } \neq 0$$

$$N_{class} + V \rightsquigarrow \text{make it add up to 1}$$

## 2.5 Log Likelihood

### *Ratio of probability*

$$ratio(w_i) = \frac{P(w_i|Pos)}{P(w_i|Neg)} \approx \frac{freq(w_i, 1) + 1}{freq(w_i, 0) + 1}$$

$$0(Negative) \longleftarrow 1(Neutral) \longrightarrow \infty(positive)$$

**Naive Bayes, without assumption**  $P(Pos) = P(Neg)$

$$\frac{P(Pos|w)}{P(Neg|w)} = \frac{P(Pos)}{P(Neg)} \cdot \prod_{i=1}^m \frac{P(w_i|Pos)}{P(w_i|Neg)}$$

$w$ : set of  $m$  words.

### **Log Likelihood**

$$\log \left( \frac{P(Pos|w)}{P(Neg|w)} \right) = \log \frac{P(Pos)}{P(Neg)} + \sum_{i=1}^m \log \left( \frac{P(w_i|Pos)}{P(w_i|Neg)} \right)$$

$\downarrow$   
log prior

$\downarrow$   
log likelihood

**Definition of**  $\lambda(w)$

$$\lambda(w) = \log \frac{P(pos)}{P(neg)} \quad (\text{stored in a table for reference})$$

$$\log \prod_{i=1}^m ratio(w_i) = \sum_{i=1}^m \lambda(w_i)$$

$$-\infty(Negative) \longleftarrow 0(Neutral) \longrightarrow \infty(positive)$$

## 2.6 Training Naive Bayes

### *Steps (Example of tweets)*

1. Get or annotate a dataset with positive and negative tweets
2. Preprocess the tweets
3. Compute  $freq(w, class)$
4. Get  $P(w|pos)$ ,  $P(w|neg)$
5. Get  $\lambda(w)$
6. Compute  $logprior = \log \frac{P(pos)}{P(neg)}$

## 2.7 Testing Naive Bayes

- $X_{val}, Y_{val}, \lambda, log\ prior$
- $score = predict(X_{val}, \lambda, log\ prior)$
- $pred = (score > 0)$
- $Accuracy = \frac{1}{m} \sum_{i=1}^m (pred_i == Y_{val\ i})$

## 2.8 Application of Naive Bayes

- Sentiment analysis
- Author identification
- Spam filtering
- Information retrieval
- Word disambiguation

## 2.9 Naive Bayes Assumptions

- Independence
- Relative frequency in corpus

## 2.10 Error Analysis

- Removing punctuation
- Removing words
- Words order
- Adversarial attacks

# Chapter 3

## Vector Space Models

### 3.1 Vector Space Models

- Know a word by the company it keeps

#### *Application*

- Information Extraction
- Machine Translation
- Chatbots

### 3.2 Word by Word and Word by Doc

#### *Word by Word Design*

Number of times they occur together within a certain distance

#### *Word by Document Design*

Number of times a word occurs within a certain category

#### *Vector Spaces*

Similarity between words/documents

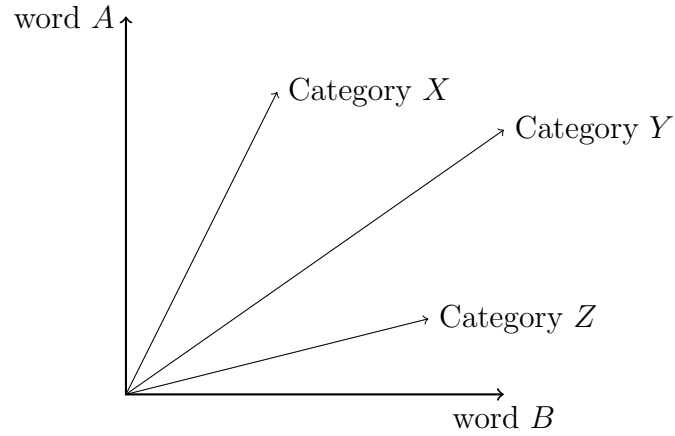


Figure 3.1: Vector Space

### 3.3 Euclidean Distance and Cosine Similarity

$$d(\vec{v}, \vec{w}) = \sqrt{\sum_{i=1}^n (v_i - w_i)^2}$$

### 3.4 Cosine Similarity

- Use cosine similarity when corpora are different sizes

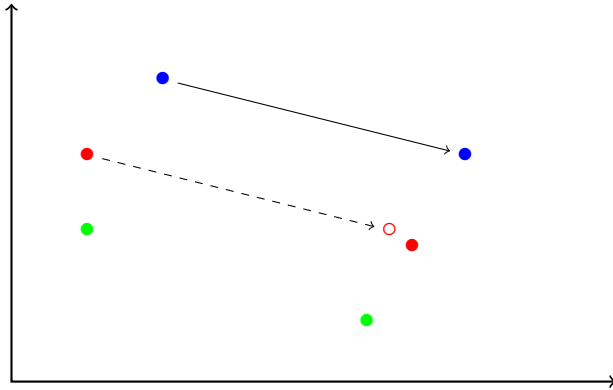
$$\cos \beta = \frac{\vec{v} \cdot \vec{w}}{|\vec{v}| |\vec{w}|}$$

- $\beta = 90^\circ \rightarrow \cos \beta = 0 \rightarrow$  dissimilar
- $\beta = 0^\circ \rightarrow \cos \beta = 1 \rightarrow$  similar



## 3.5 Manipulating Words in Vector Space

*Use Known Relationships to Make Prediction*



## 3.6 Visualization and PCA

- Original Space  $\rightarrow$  Uncorrelated features  $\rightarrow$  Dimension reduction

## 3.7 PCA Algorithm

### *Eigenvector*

Uncorrelated features for data gives the direction of uncorrelated features

### *Eigenvalue*

The amount of information retained by each feature

# Chapter 4

## Machine Translation

### 4.1 Transforming Word Vectors

*Transformation (using a matrix)*

$$\begin{array}{ccc} & XR \approx Y & \\ & \searrow \quad \swarrow & \\ & \text{subsets of the full vocabulary} & \end{array}$$

*Solving for  $R$*

- Initialize  $R$
- In a loop,

$$\begin{aligned} Loss &= ||XR - Y||_F \\ g &= \frac{d}{dR} Loss \quad (\text{gradient}) \\ R &= R - \alpha g \quad (\text{update}) \\ &\quad \downarrow \\ &\quad \text{learning rate} \end{aligned}$$

For a matrix  $A$ ,

$$||A||_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$

## 4.2 K-nearest Neighbors

### *Translation*

word  $\xrightarrow{R}$  [transformed]  $\xrightarrow[\text{words}]{\text{find similar}}$  word'

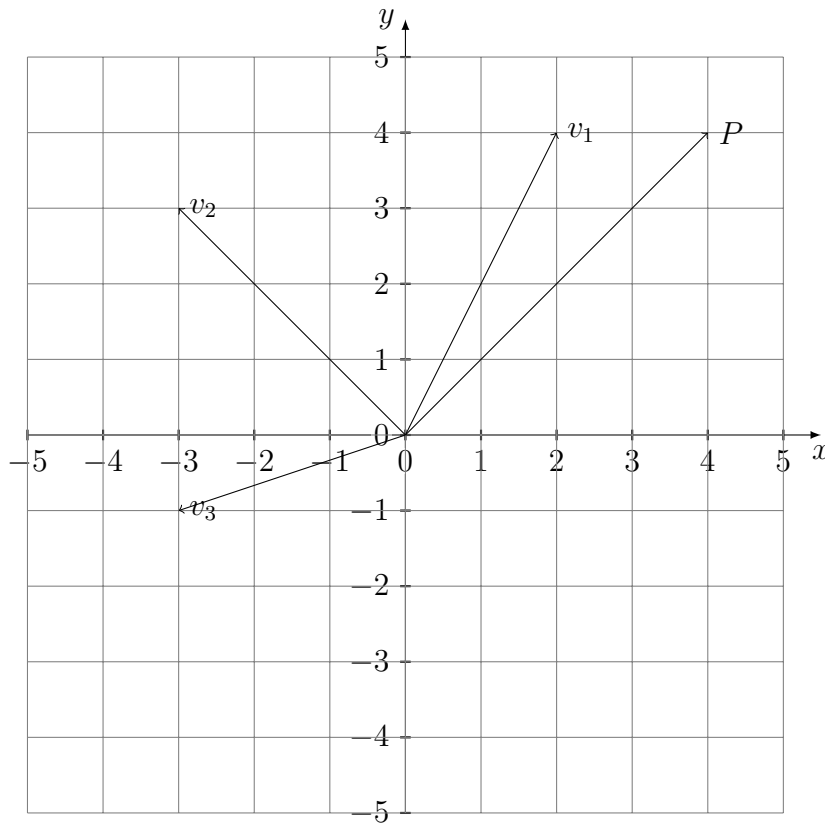
## 4.3 Hash Tables and Hash Functions

### *Hash Function*

A function that takes data of arbitrary sizes and maps it to a fixed value. The values returned are known as hash values or even hashes.

## 4.4 Locality Sensitive Hashing

*To hash similar inputs into the same buckets with high probability*



- $Pv_1^T > 0$
- $Pv_2^T = 0$

- $Pv_3^T < 0$

## 4.5 Multiple Planes

*Use multiple planes to get a single hash value*

*Example*

- $\vec{v}$ , planes  $\vec{P}_1, \vec{P}_2, \vec{P}_3$
- $hash = 2^0 \cdot h_1 + 2^1 \cdot h_2 + 2^2 \cdot h_3$
- $h_i = \begin{cases} 1 & \text{if } \vec{P}_i \cdot \vec{v} > 0 \\ 0 & \text{else} \end{cases}$

## 4.6 Approximate Nearest neighbors

- Does not give the full nearest neighbors
- Trade off accuracy for efficiency