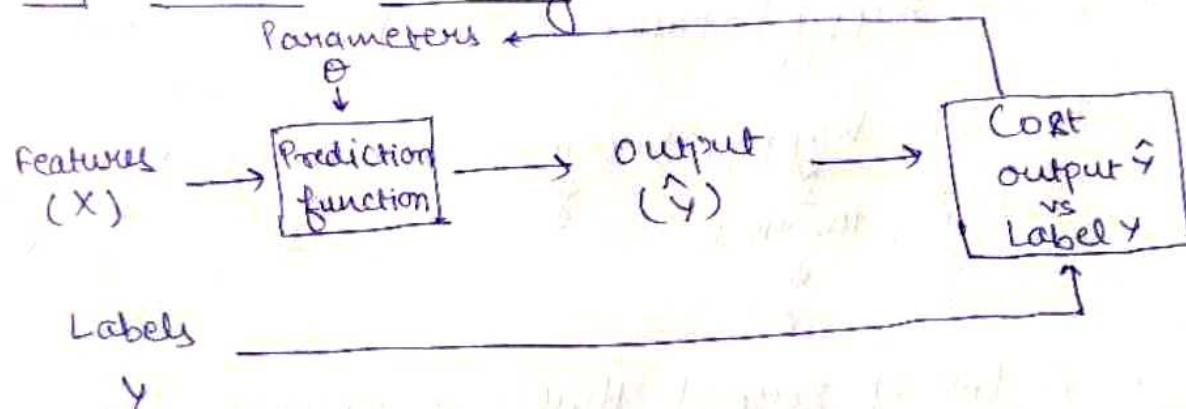01/07/2020

[ deeplearning.ai ]

## Natural Language Processing with classification and Vector Spaces :-

# WEEK-1 : Vocabulary & feature Extraction

~o Sentiment Analysis with logistic Regression

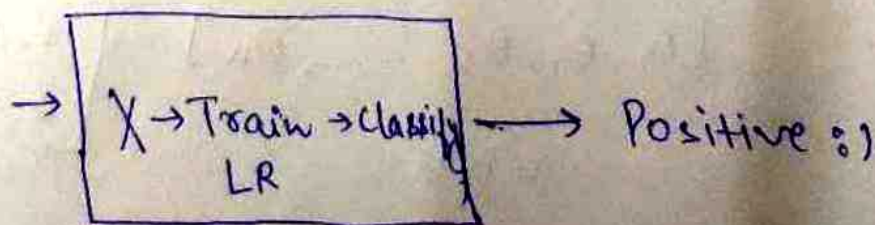→ Supervised ML & Sentiment Analysis

→ Supervised ML (training)

Parameters $\theta$

Features (X) → Prediction function → output $(\hat{y})$ → Cost output $\hat{y}$ vs Label Y

Labels Y

→ Sentiment Analysis

Tweet : I am happy because I am learning NLP.

Positive : 1

Negative : 0

↓

Logistic Regression

I am happy because I am learning NLP → [ X → Train → Classify LR ] → Positive :)

# Vocabulary & Feature Extraction

→ **Vocabulary:**

Tweets:

[tweet_1, tweet_2, ..., tweet_m]

> I am happy because I am learning NLP
> ...
> ...
> ...
> I hated the movie

$V = [$ I, am, happy, because, learning, NLP, hated, the, movie $]$

→ **Feature Extraction:**

- I am happy because I am learning NLP.

[I, am, happy, because, learning, NLP, ..., hated,
the, movie]

I → 1, am → 1, happy → 1, because → 1, learning → 1, NLP → 1, ... , hated → 0,
the → 0, movie → 0

- A lot of zeros! That's a sparse representation.

→ **Problems with sparse representations**

I am happy because I am learning NLP

$[1,1,1,1,1,1, ..., 0, ... , 0,0,0]$ → All zeros!

$1 \longleftrightarrow |V|$

$[\theta_0, \theta_1, \theta_2, ...., \theta_n]$ →
1. Large Training time
2. Large Prediction time

$n = |V|$

# Negative & Positive Frequencies

→ **Positive and Negative counts:**

### Corpus

| |
|---|
| I am happy because I am learning NLP |
| I am happy |
| I am sad, I am not learning NLP |
| I am sad |

**Vocabulary**
- I
- am
- happy
- because
- learning
- NLP
- sad
- not

### Positive Tweets

| |
|---|
| I am happy because I am learning NLP |
| I am happy |

→

| vocabulary | Posfreq $(c_i)$ |
|---|---|
| I | 3 |
| am | 3 |
| happy | 2 |
| because | 1 |
| learning | 1 |
| NLP | 1 |
| sad | 0 |
| not | 0 |

### Negative Tweets

| |
|---|
| I am sad, I am not learning NLP |
| I am sad |

→

| vocabulary | Negfreq $(0)$ |
|---|---|
| I | 3 |
| am | 3 |
| happy | 0 |
| because | 0 |
| learning | 1 |
| NLP | 1 |
| sad | 2 |
| not | 1 |

→ **Word frequency in classes:**

| vocabulary | Posfreq (1) | Negfreq (0) |
|---|---|---|
| I | 3 | 3 |
| am | 3 | 3 |
| happy | 2 | 0 |
| because | 1 | 0 |
| learning | 1 | 1 |
| NLP | 1 | 1 |
| sad | 0 | 2 |
| not | 0 | 1 |

freqs : dictionary mapping from (word, class) to frequency