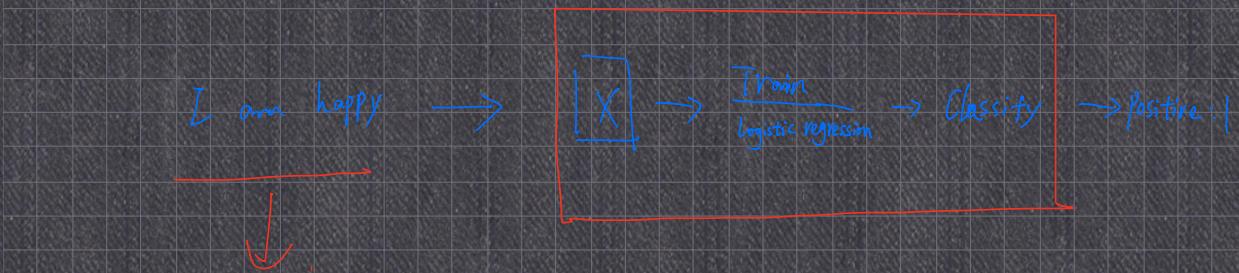


To perform sentiment analysis



Vocabulary & Feature Extraction

$$V = [1, 1, 1, 0 \dots 0] \rightarrow \begin{cases} 1. \text{ Large training time} \\ 2. \text{ Large prediction time.} \end{cases}$$

对应出现的就有1，没有的就是0。

(parameter: $n+1$.
单词数.)

Feature Extraction with Frequencies

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

单词在 Pos 中的频率 在 Neg 中的频率

Preprocessing
数据预处理 .

如. / 去掉 the, a, an 等这类不重要的 word
将 -ing, -ed, -s 等的转过去成原型 .

NLTK

```
from nltk.corpus import stopwords #停词
from nltk.stem import PorterStemmer #词干提取
from nltk.tokenize import TweetTokenizer #将字符串拆分为每个单词.
```

1. 删除字符串中不需要的字符、文本(如:链接、脚注)
可用 `re.sub()`

2. 将字符串拆分为每个单词.

```
tokenizer = TweetTokenizer(preserve_case=False, strip_handles=True, reduce_len=True)
tweet_tokens = tokenizer.tokenize(tweet)
```

3. 删除停用词和标点符号.

```
stop_words = stopwords.words('english')
if word not in stop_words:
    pass
```

4. 词干提取 .

```
stemmer = PorterStemming()
stem_word = stemmer.stem(word)
```

Putting it all together.

文本 → 预处理 → 提取特征 → 文本转化为数字 .

I am Happy Because i am learning NLP @deeplearning

↓ Preprocessing

[happy, learn, nlp]

↓ Feature Extraction

Bias ← [1, 4, 2] → Sum negative frequencies

↓ Sum positive frequencies

将每条 tweet 与矩阵 X = $\begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \\ \vdots & \vdots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} \end{bmatrix}$ 相乘

```
freqs = build_freqs(tweets, labels) #Build frequencies dictionary

X = np.zeros((m, 3)) #Initialize matrix X

for i in range(m): #For every tweet
    p_tweet = process_tweet(tweets[i]) #Process tweet
    X[i, :] = extract_features(p_tweet, freqs) #Extract Features
```

Building and Visualizing word frequencies

将数据叠加，变成频率矩阵。

创建一个矩阵表示 tweet 的标签 (0, 1)

$\text{labels} = \text{np.append}(\text{np.ones}(\text{len(all_positive)}), \text{np.zeros}(\text{len(all_negative)}))$

`def build_freqs(tweets, labels):`

`label_list = np.squeeze(labels).tolist()`

去掉多余的维度

转换为 list() 是为了迭代。

`freqs = {}`

`for y, words in zip(label_list, tweets):`

`for word in process_tweet(words):`

`pair = (word, y)`

`if pair in freqs:`

`freqs[pair] += 1`

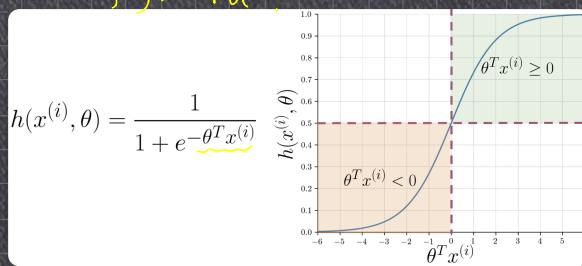
`else:`

`freqs[pair] = 1`

`return freqs.`

Logistic Regression Overview

Sigmoid:

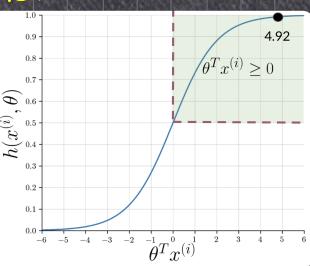


如何对 tweet 进行预测：

@Ymourri and
@AndrewYNg are tuning a
GREAT AI model

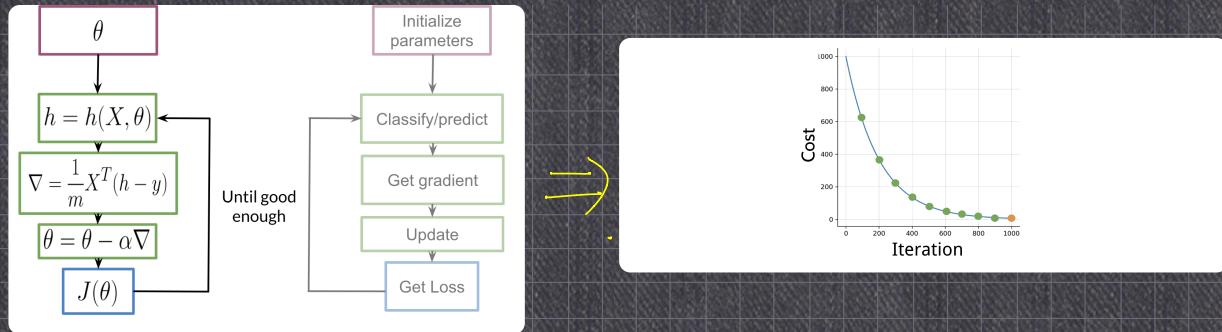
[tun, ai, great, model]

$$x^{(i)} = \begin{bmatrix} 1 \\ 3476 \\ 245 \end{bmatrix} \quad \theta = \begin{bmatrix} 0.00003 \\ 0.00150 \\ -0.00120 \end{bmatrix}$$



Training

初始化参数 θ → 更新 θ → 迭代降低 cost.



Testing

$> 0.5 \rightarrow \text{positive}$

根据 sigmoid 转出 pred \Rightarrow $< 0.5 \rightarrow \text{negative}$

- $X_{val} Y_{val} \theta$

$h(X_{val}, \theta)$

$pred = h(X_{val}, \theta) \geq 0.5$

$$\begin{bmatrix} 0.3 \\ 0.8 \\ 0.5 \\ \vdots \\ h_m \end{bmatrix} \geq 0.5 = \begin{bmatrix} \frac{0.3 \geq 0.5}{0.3 \geq 0.5} \\ \frac{0.8 \geq 0.5}{0.8 \geq 0.5} \\ \frac{0.5 > 0.5}{0.5 > 0.5} \\ \vdots \\ pred_m \geq 0.5 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ \vdots \\ pred_m \end{bmatrix}$$

将预测的值与真实的结果比较，得到准确率。

$$\text{Accuracy} \longrightarrow \sum_{i=1}^m \frac{(pred^{(i)} == y_{val}^{(i)})}{m}$$

Cost Function

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

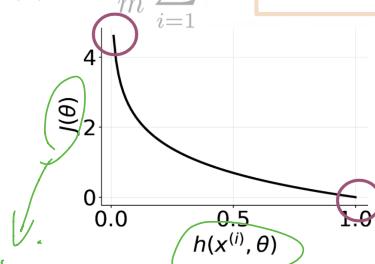
$y=1$, $h(x^{(i)}, \theta)=0$ 时, $J(\theta)=\infty$

$y=0$, $h(x^{(i)}, \theta)=1$ 时, $J(\theta)=\infty$

$y=1$, $h(x^{(i)}, \theta)=1$ 时, $J(\theta)=0$

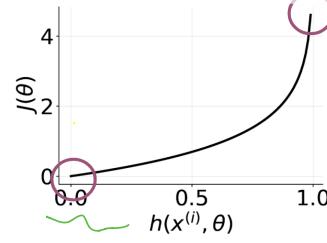
$y=0$, $h(x^{(i)}, \theta)=0$ 时, $J(\theta)=0$

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



Cost 减小

当 $y=1$, 预测为 1 时



当 $y=0$ 时, 预测为 0.
Cost 最小.