

Text Classification

COSC 6336: Natural Language Processing
Spring 2020

Previous Lecture

★ Preprocessing text

Today's lecture

- ★ Text Classification
 - Example applications
 - Task definition
- ★ Classical approaches to Text Classification
 - Naive Bayes
 - Logistic Regression
 - Multinomial Logistic Regression or MaxEnt

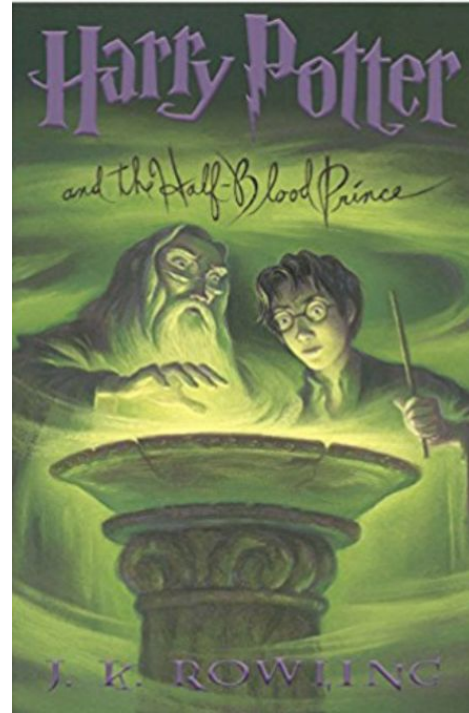
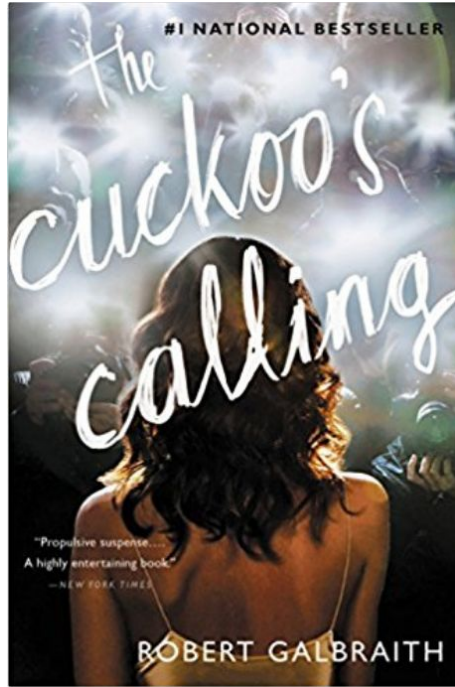
Hello Dear

My name is: Engineer Mrs.Eva Rose, I am a citizen of Austria, a business woman specialized in mining of raw Gold in Africa; but now I am critically sick with esophageal cancer which has damaged almost all the cells in my body system and I will soon die according to my doctors.

My late husband died in an accident with our two daughters few years ago leaving me with our only son whom is just 10 years old and he is my most concern now as he is still a child and does not know anything about live and has nobody to take care of him after I am dead; because I and my late husband does not have any relatives, we both grew up in the orphanage home and got married under orphanage as orphans. So if I die now my innocent child would be left alone in this wicked world and I do not wish to send him to any orphanage home, I want him to grow up in the hands of an individual, not orphanage.

Please, i am begging you in the name of God to sincerely accept my proposal; let me instruct my bank to wire transfer my fund worth the sum of US\$ 15,000,000.00 (FIFTEEN MILLION DOLLARS) to your account in your country immediately, then you take my son to your home and raise him as your own son. As you receive the fund into your account, you are entitled to take 30 percent and invest 70 percent for my son so that he will not suffer in his entire life.

What do these books have in common?



Other TC tasks

A Cancer Conundrum: Too Many Drug Trials, Too Few Patients

Breakthroughs in immunotherapy and a rush to develop profitable new treatments have brought a crush of clinical trials scrambling for patients.

By GINA KOLATA

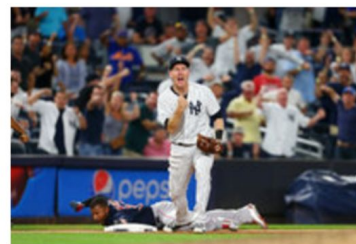
Yankees and Mets Are on Opposite Tracks This Subway Series

As they meet for a four-game series, the Yankees are playing for a postseason spot, and the most the Mets can hope for is to play spoiler.

By FILIP BONDY



→ Health



→ Sports

~20 classes

Other TC tasks: Entity Linking

Although he originally won the event, the United States Anti-Doping Agency announced in August 2012 that they had disqualified **Armstrong** from his seven consecutive Tour de France wins from 1999–2005.



Lance Edward Armstrong is an American former professional road cyclist



Armstrong County is a county in Pennsylvania...

?

?

- ▶ 4,500,000 classes (all articles in Wikipedia)

Other tasks that can be solved as TC

- ★ Sentiment classification
- ★ Native language identification
- ★ Author Profiling
- ★ Depression detection
- ★ Cyberbullying detection
- ★

Formal definition of the TC task

- ★ Input:
 - a document d
 - a fixed set of classes $C = \{c_1, c_2, \dots, c_J\}$
- ★ Output: a predicted class $c \in C$

Methods for TC tasks

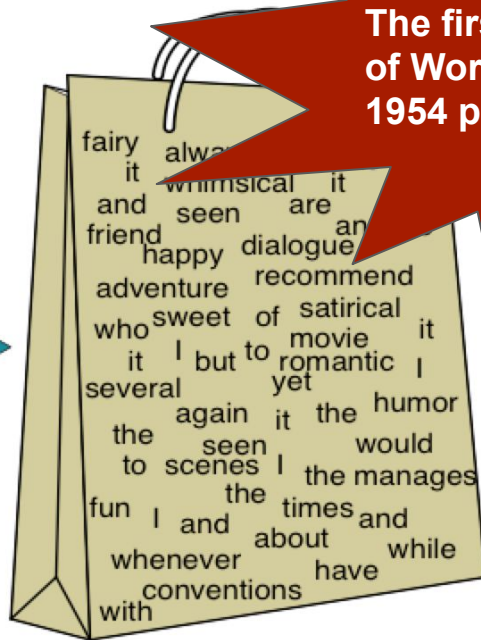
- ★ Rule based approaches
- ★ Machine Learning algorithms
 - Naive Bayes
 - Support Vector Machines
 - Logistic Regression
 - And now deep learning approaches

Naive Bayes for Text Classification

- ★ Simple approach
- ★ Based on the bag-of-words representation

Bag of words

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



The first reference to Bag of Words is attributed to a 1954 paper by Zellig Harris

3
3
seen
et
1
would
1
whimsical
1
times
1
sweet
1
satirical
1
adventure
1
genre
1
fairy
1
humor
1
have
1
great
1
...
...

Naive Bayes

Probabilistic classifier $\hat{c} = \operatorname{argmax}_{c \in C} P(c|d)$ (eq. 1)

According to Bayes rule: $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$ (eq. 2)

Replacing eq. 2 into eq. 1: $\hat{c} = \operatorname{argmax}_{c \in C} P(c|d) = \operatorname{argmax}_{c \in C} \frac{P(d|c)P(c)}{P(d)}$

Dropping the denominator:

$$\operatorname{argmax}_{c \in C} P(d|c)P(c)$$

Naive Bayes

A document d is represented as a set of features f_1, f_2, \dots, f_n

$$\hat{c} = \operatorname{argmax}_{c \in C} \overbrace{P(f_1, f_2, \dots, f_n | c)}^{\text{likelihood}} \overbrace{P(c)}^{\text{prior}}$$

How many parameters do we need to learn in this model?

Naive Bayes Assumptions

1. Position doesn't matter
2. Naive Bayes assumption: probabilities $P(f_i|c)$ are independent given the class c and thus we can multiply them:

$$P(f_1, f_2, \dots, f_n | c) = P(f_1 | c) \cdot P(f_2 | c) \cdot \dots \cdot P(f_n | c)$$

This leads us to:

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c) \prod_{f \in F} P(f | c)$$

Naive Bayes in Practice

We consider word positions:

positions \leftarrow all word positions in test document

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c) \prod_{i \in \text{positions}} P(w_i | c)$$

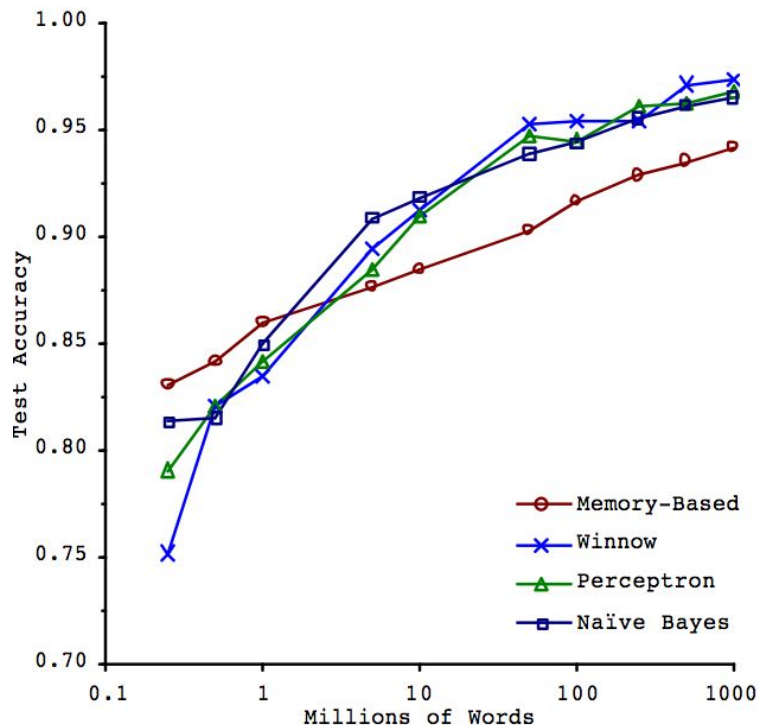
We also do everything in log space:

$$c_{NB} = \operatorname{argmax}_{c \in C} \log P(c) + \sum_{i \in \text{positions}} \log P(w_i | c)$$

Naive Bayes: Training

How do we compute $P(c)$ and $P(f_i|c)$?

Is Naive Bayes a good option for TC?



Scaling to Very Very Large Corpora for Natural Language Disambiguation

Michele Banko and Eric Brill

Microsoft Research

1 Microsoft Way

Redmond, WA 98052 USA

{mbanko,brill}@microsoft.com

Abstract

The amount of readily available on-line text has reached hundreds of billions of words and continues to grow. Yet for most core natural language tasks, algorithms continue to be optimized, tested and compared after training on corpora consisting of only one million words or less. In this paper, we

potentially large cost of annotating data for those learning methods that rely on labeled text.

The empirical NLP community has put substantial effort into evaluating performance of a large number of machine learning methods over fixed, and relatively small, data sets. Yet since we now have access to significantly more data, one has to wonder what conclusions that have been drawn on small data sets may carry over when these learning methods are trained using much larger corpora.

Evaluation in TC

Confusion table

	Gold Standard	
	True	False
True	TP = true positives	FP = False positives
False	FN = false negatives	TN = True negatives

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

Evaluation in TC: Issues with Accuracy?

Suppose we want to learn to classify each message in a web forum as “extremely negative”. The class “extremely negative” is the class of interest. We have a collected gold standard data:

- ★ 990 instances are labeled as negative
- ★ 10 instances are labeled as positive
- ★ Test data has 100 instances (99- and 1+)
- ★ A dumb classifier can get 99% accuracy by always predicting “negative” !

More Sensible Metrics: Precision, Recall and F-measure

$$P = TP / (TP + FP)$$

$$R = TP / (TP + FN)$$

	Gold Standard	
	True	False
True	TP = true positives	FP = False positives
False	FN = false negatives	TN = True negatives

$$\text{F-measure} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

What about Multi-class problems?

- Multi-class: $c > 2$
- P, R, and F-measure are defined for a single class
- We assume classes are mutually exclusive
- We use per class evaluation metrics

$P = \frac{c_{ii}}{\sum_j c_{ji}}$	$R = \frac{c_{ii}}{\sum_j c_{ij}}$
------------------------------------	------------------------------------

Micro vs Macro Average

- ★ **Macro** average: measure performance **per class** and **then average**
- ★ **Micro** average: collect predictions for **all classes** then compute TP, FP, FN, and TN
- ★ **Weighted** average: compute performance per label and then average where each label score is weighted by its support

Example

Class 1: Urgent

	true urgent	true not
system urgent	8	11
system not	8	340

$$\text{precision} = \frac{8}{8+11} = .42$$

Class 2: Normal

	true normal	true not
system normal	60	55
system not	40	212

$$\text{precision} = \frac{60}{60+55} = .52$$

Class 3: Spam

	true spam	true not
system spam	200	33
system not	51	83

$$\text{precision} = \frac{200}{200+33} = .86$$

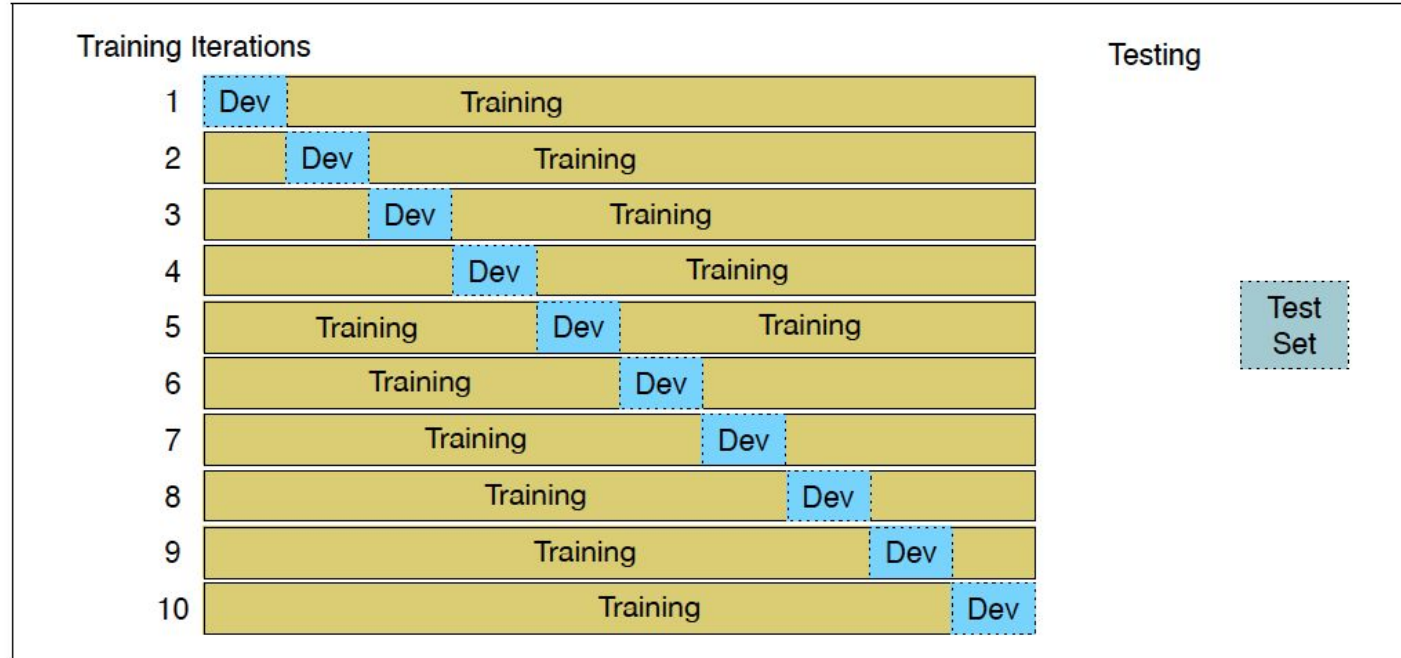
Pooled

	true yes	true no
system yes	268	99
system no	99	635

$$\text{microaverage precision} = \frac{268}{268+99} = .73$$

$$\text{macroaverage precision} = \frac{.42+.52+.86}{3} = .60$$

Train/Test Data Separation



Logistic Regression (MaxEnt)

Disadvantages of NB

- ★ Correlated features
 - Their evidence is overestimated
- ★ This will hurt classifier performance

Logistic Regression

- ★ No independence assumption
- ★ Feature weights are learned simultaneously

Logistic Regression

Computes $p(y|x)$ directly from training data:

$$\sum_{i=0}^{|X|} w_i x_i$$

Logistic Regression

Computes $p(y|x)$ directly from training data:

$$z = \sum_{i=0}^{|X|} w_i x_i$$

But since the result of this will not be a proper probability function we need to normalize this:

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

Logistic Regression

- ★ We then compute the probability $(y|x)$ for class 1 and 2
- ★ We choose the class with the highest probability

Features in Text Classification

★ We can have **indicator** functions of the form:

$$f_1(c, x) = \begin{cases} 1 & \text{if "great"} \in x \text{ \& } c = + \\ 0 & \text{otherwise} \end{cases}$$

$$f_2(c, x) = \begin{cases} 1 & \text{if "second-rate"} \in x \text{ \& } c = - \\ 0 & \text{otherwise} \end{cases}$$

$$f_3(c, x) = \begin{cases} 1 & \text{if "no"} \in x \text{ \& } c = - \\ 0 & \text{otherwise} \end{cases}$$

$$f_4(c, x) = \begin{cases} 1 & \text{if "enjoy"} \in x \text{ \& } c = - \\ 0 & \text{otherwise} \end{cases}$$

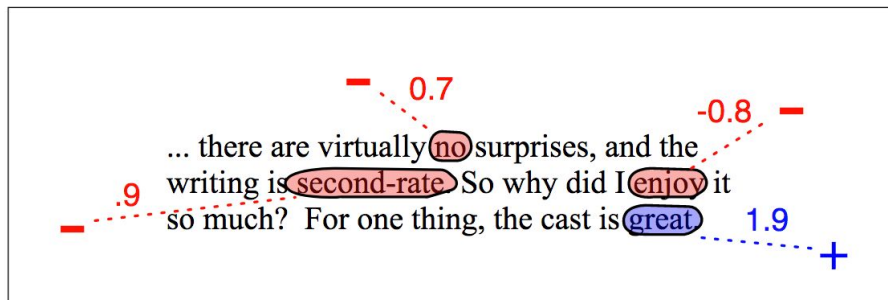
Features in Text Classification

★ Or other functions:

$$f(d) = \begin{pmatrix} \text{count}(\text{"boring"}) \\ \text{count}(\text{"not boring"}) \\ \text{length of document} \\ \text{author of document} \\ \vdots \end{pmatrix}$$

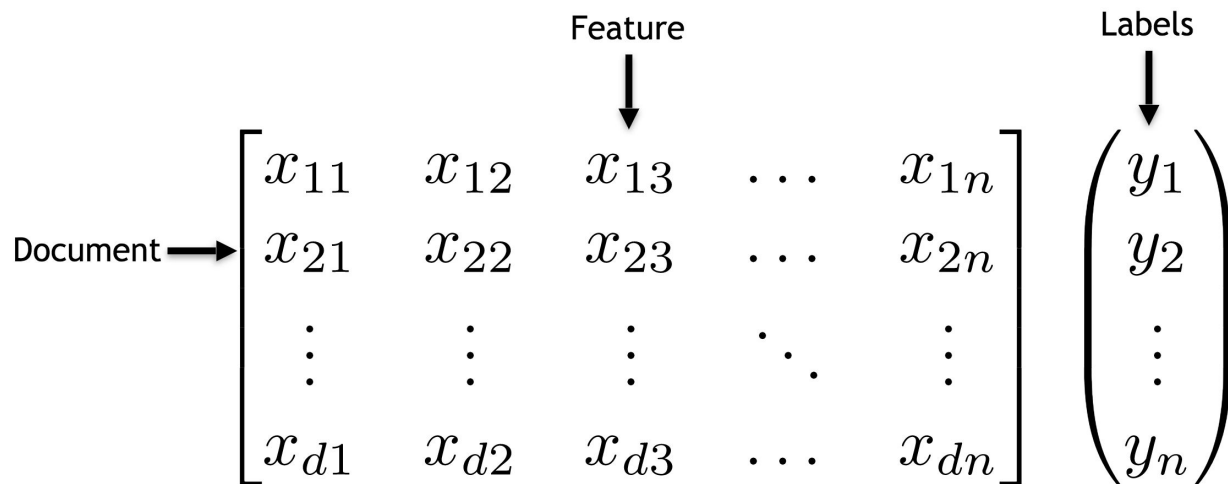
An Example

Assume weights are as follows: $[-5.0, 2.0]$, while $b = 0.1$



Red color is x_1 , blue color x_2 . What would this model predict?

Learning in Logistic Regression



Learning in Logistic Regression

How does training look like?

We need to find the right weights, **what do we mean by this?**

How?

By using the **Conditional Maximum Likelihood Principle** that measures how close $(\hat{y} = \sigma(w \cdot x + b))$ is to the actual true labels y . We focus on distance (loss, cost function):

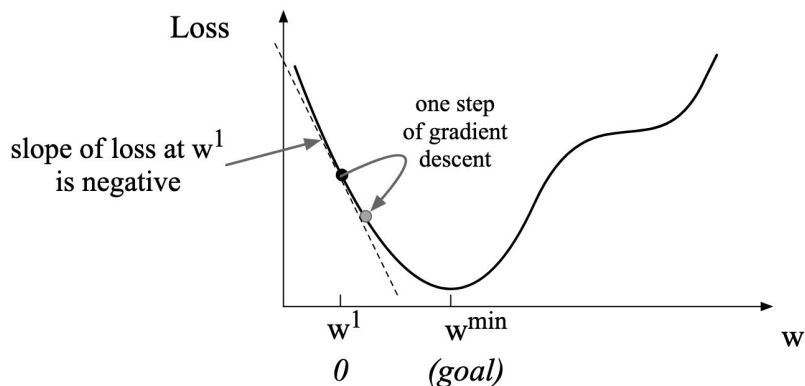
$$L_{CE}(w, b) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log (1 - \sigma(w \cdot x + b))]$$

Learning in Logistic Regression

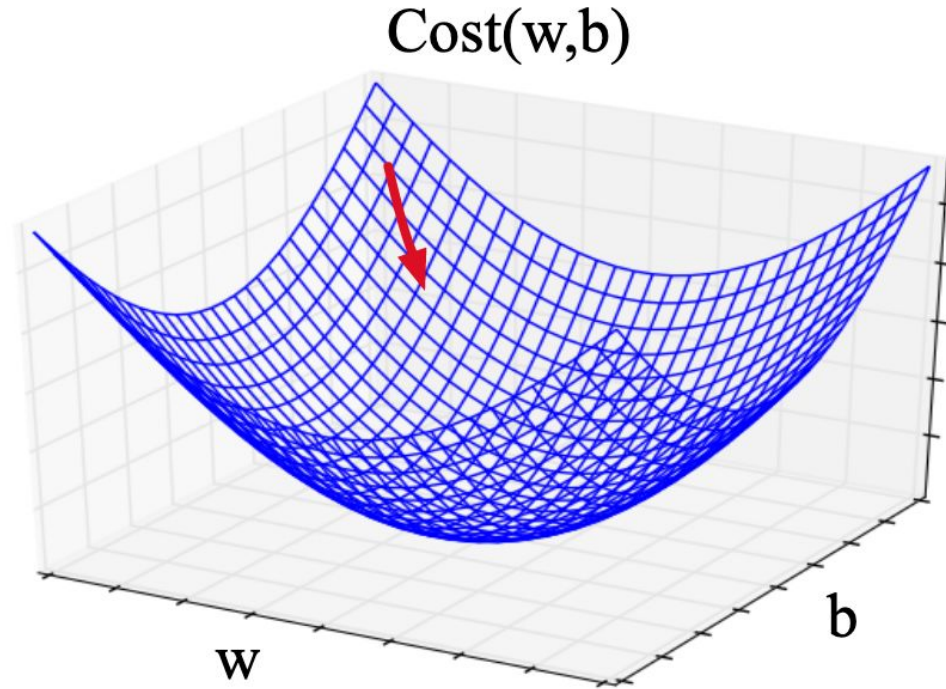
- ★ We use optimization methods like Stochastic Gradient Descent to find the weights that minimize $L(w)$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m L_{CE}(y^{(i)}, x^{(i)}; \theta)$$

- ★ We start with $w = w^1$ and move in the direction of the gradient (the partial derivative $L'(w)$)



Learning in Logistic Regression



Learning in Logistic Regression

More than 1 feature:

$$\nabla_{\theta} L(f(x; \theta), y) = \begin{bmatrix} \frac{\partial}{\partial w_1} L(f(x; \theta), y) \\ \frac{\partial}{\partial w_2} L(f(x; \theta), y) \\ \vdots \\ \frac{\partial}{\partial w_n} L(f(x; \theta), y) \end{bmatrix}$$

The update to the weights is computed as follows:

$$\theta_{t+1} = \theta_t - \eta \nabla L(f(x; \theta), y)$$

Learning in Logistic Regression

Definition of gradient in Logistic Regression:

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

Regularization

- ★ **Overfitting**: when a model “memorizes” the training data
- ★ To prevent overfitting we penalize large weights:

$$\hat{\theta} = \operatorname{argmax}_{\theta} \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}) - \alpha R(\theta)$$

- ★ Regularization forms:

- L2 norm (Euclidean)

$$R(\theta) = \|\theta\|_2^2 = \sum_{j=1}^n \theta_j^2$$

- L1 norm (Manhattan)

$$R(\theta) = \|\theta\|_1 = \sum_{i=1}^n |\theta_i|$$

What if we have more than two classes?

- ★ We use multinomial logistic regression
- ★ We replace the sigmoid function with the softmax

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \quad 1 \leq i \leq k$$

Input to softmax is still the dot product of features and weights. Now we keep weight vector for each class:

$$p(y = c|x) = \frac{e^{w_c \cdot x + b_c}}{\sum_{j=1}^k e^{w_j \cdot x + b_j}}$$

What if we have more than two classes?

Learning

$$\begin{aligned} L_{CE}(\hat{y}, y) &= - \sum_{k=1}^K 1\{y = k\} \log p(y = k|x) \\ &= - \sum_{k=1}^K 1\{y = k\} \log \frac{e^{w_k \cdot x + b_k}}{\sum_{j=1}^K e^{w_j \cdot x + b_j}} \end{aligned}$$

In Summary

- ★ NB is a generative model
 - Highly correlated features can result in bad results
- ★ Logistic Regression
 - Weights for feature are calculated simultaneously
- ★ NB seems to fare better in small data sets
- ★ Both NB and Logistic regression are linear classifiers
- ★ Regularization helps alleviate overfitting