

Introduction to NLP

511.
Classification

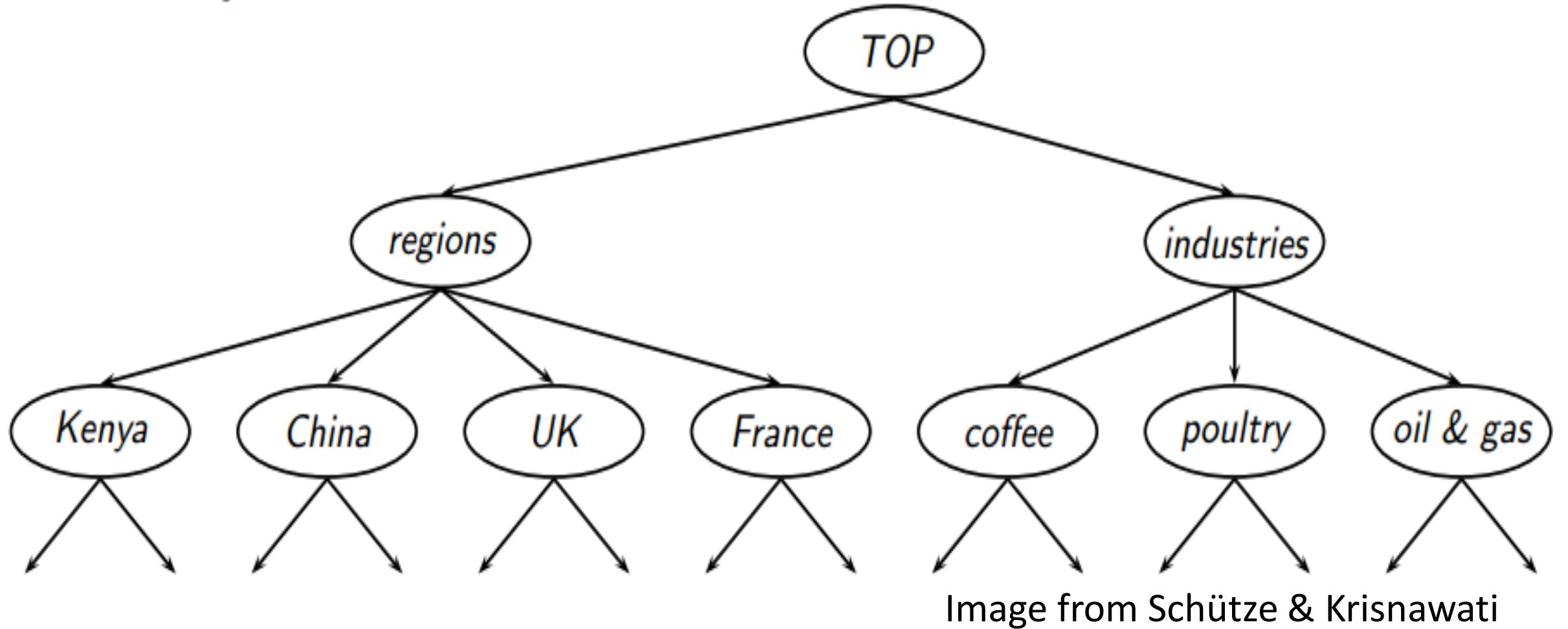
Classification

- Assigning documents or sentences to predefined categories
 - topics, languages, users ...
- *Input:*
 - a document (or sentence) d
 - a fixed set of classes $C = \{c_1, c_2, \dots, c_J\}$
- *Output:* a predicted class $c \in C$

Variants of the Problem Formulation

- Binary categorization: only two categories
 - Retrieval: {relevant-doc, irrelevant-doc}
 - Spam filtering: {spam, non-spam}
 - Opinion: {positive, negative}
- K-category categorization: more than two categories
 - Topic categorization: {sports, science, travel, business,...}
 - Word sense disambiguation:{bar1, bar2, bar3, ...}
- Hierarchical vs. flat
- Overlapping (soft) vs non-overlapping (hard)

Hierarchical Classification



Components of a Classification System

1. A **feature representation** of the input. For each input observation $x^{(i)}$, this will be a vector of features $[x_1, x_2, \dots, x_n]$. We will generally refer to feature i for input $x^{(j)}$ as $x_i^{(j)}$, sometimes simplified as x_i , but we will also see the notation f_i , $f_i(x)$, or, for multiclass classification, $f_i(c, x)$.
2. A classification function that computes \hat{y} , the estimated class, via $p(y|x)$. In the next section we will introduce the **sigmoid** and **softmax** tools for classification.
3. An objective function for learning, usually involving minimizing error on training examples. We will introduce the **cross-entropy loss function**
4. An algorithm for optimizing the objective function. We introduce the **stochastic gradient descent** algorithm.

Mathematical Formulation

- ▶ Input: a numeric representation x of text d
- ▶ Output: scores $\Psi(x, y; \theta) \in \mathbb{R}, \forall y \in \mathcal{Y}$

Classification

$$\hat{y} = \arg \max_{y' \in \mathcal{Y}} \Psi(x, y'; \theta)$$

Mathematical Formulation

- ▶ Training set $\mathcal{T} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
- ▶ Development set $\mathcal{D} = \{(\mathbf{x}^{(j)}, y^{(j)})\}_{j=1}^M$
- ▶ Test set $\mathcal{U} = \{(\mathbf{x}^{(l)}, y^{(l)})\}_{l=1}^L$

Mathematical Formulation

Score function

$$\Psi(x, y; \theta) = \theta^\top f(x, y)$$

- ▶ How to design a feature function $f(x, y)$?
 - ▶ Bag-of-words representation
- ▶ How to learn θ ?
 - ▶ Perceptron algorithm

Bag of Words Representation

- Doc 1= “information retrieval”
 - Doc 2 = “computer information retrieval information”
 - Doc 3 = “computer computer retrieval”
-
- Vocabulary: information, retrieval, computer
 - Doc 1 = $\langle 1, 1, 0 \rangle$
 - Doc 2 = $\langle 2, 1, 1 \rangle$
 - Doc 3 = $\langle 0, 1, 2 \rangle$
- information, retrieval, computer
- $$D = \begin{bmatrix} 1 & 1 & 0 \\ 2 & 1 & 1 \\ 0 & 1 & 2 \end{bmatrix}$$

Question: Doc 4 = “retrieval information retrieval” ?

Mathematical Formulation

Decision function

$$\Psi(x, y) = w_y^\top f(x, \theta)$$

- ▶ x : data point
- ▶ y : label
- ▶ w_y : classification weights with respect to label y
- ▶ $f(x, \theta)$: feature function
- ▶ θ : parameter of feature function

Features for Classification

- Vector-based
 - Words: “cat”, “dog”, “great”, “horrible”, etc.
 - Meta features: document length, author name, etc.
 - Each document (or sentence) is represented as a vector in an n -dimensional space
 - Similar documents appear nearby in the vector space
 - (more later)

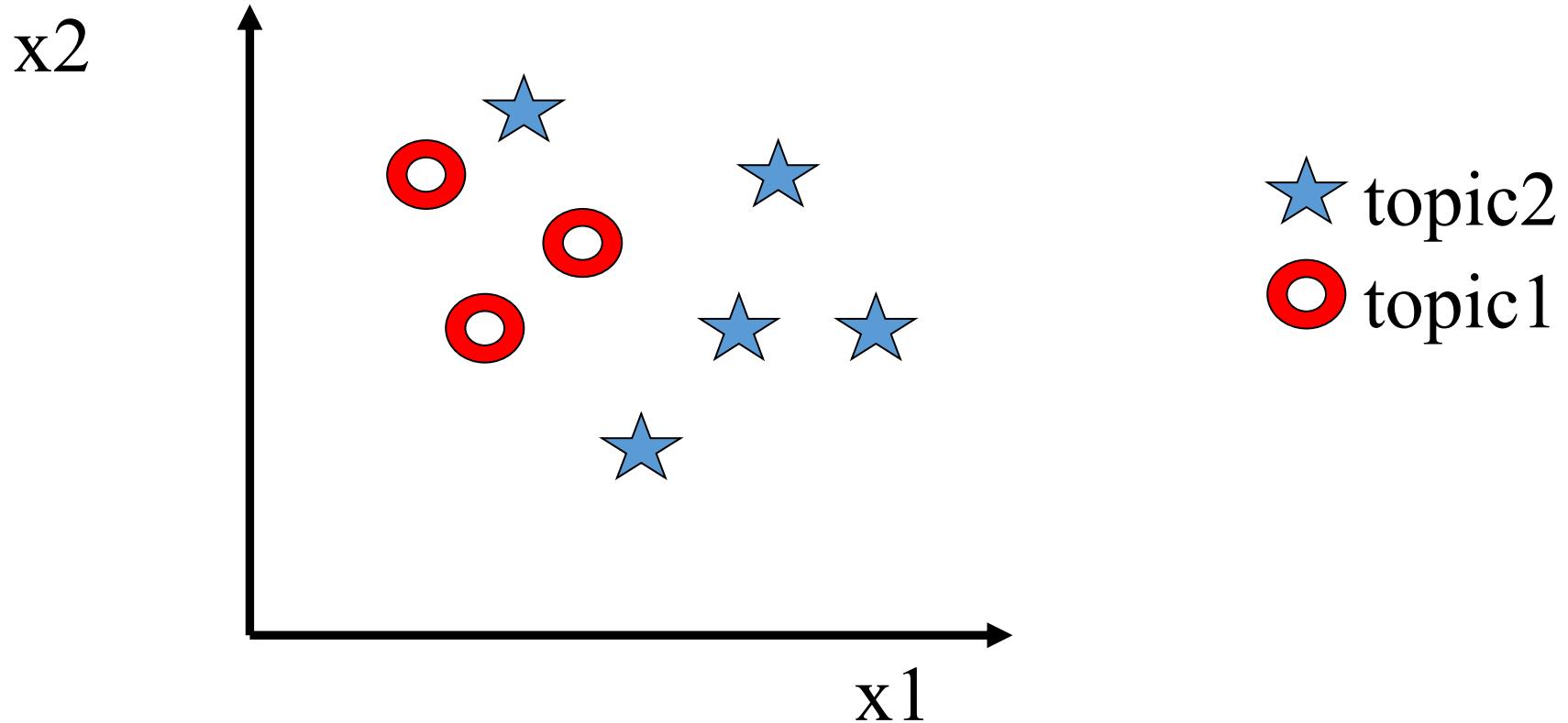
Classification in NLP

- Part of speech tagging
 - Sentiment analysis
 - Word sense disambiguation
 - Parsing
 - Optical character recognition
 - Spelling correction
 - Named entity classification
- assign the correct word or phrase to a given misspelled input by choosing the best candidate from a set of alternatives.

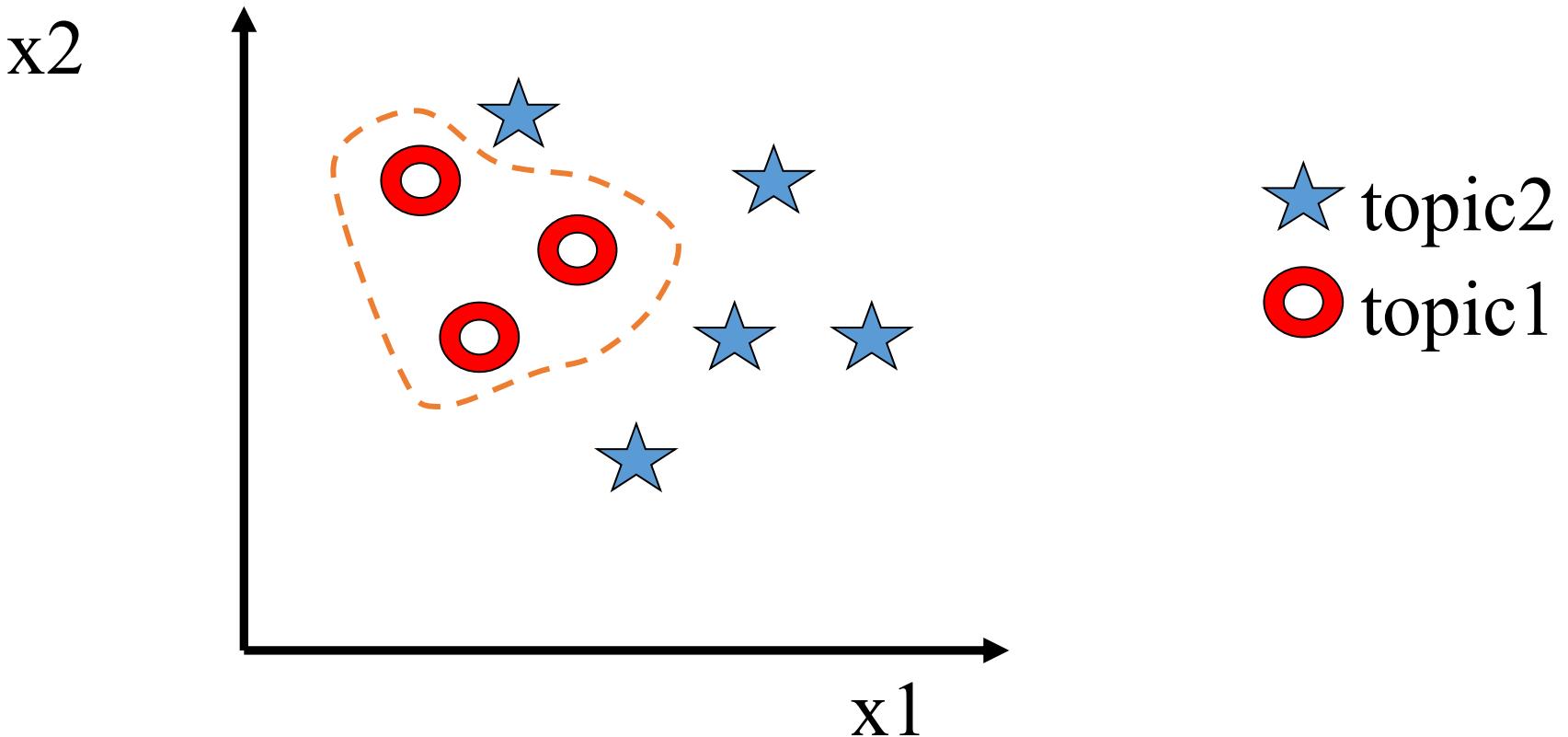
Introduction to NLP

Vector Space Classification

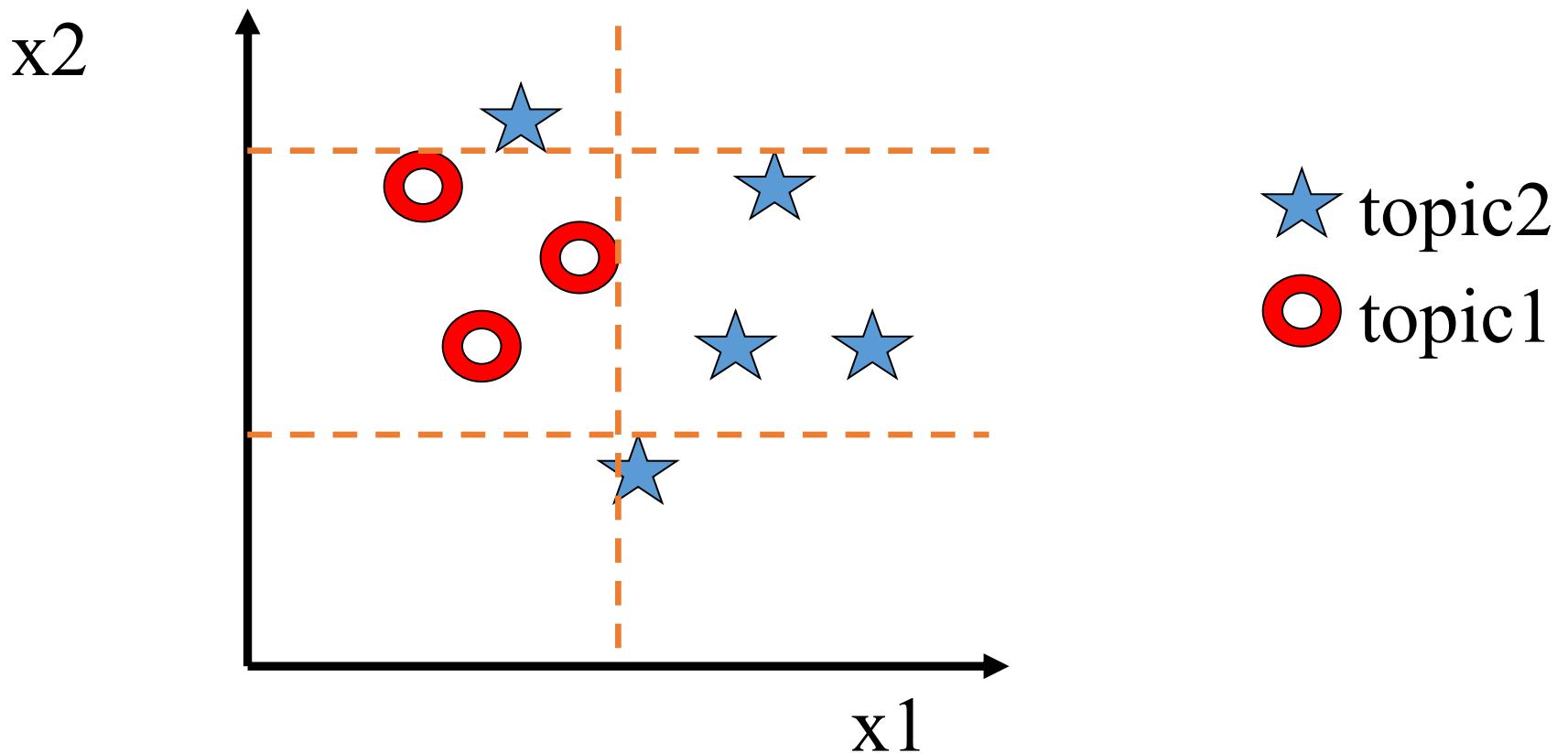
Vector Space Classification



Decision surfaces



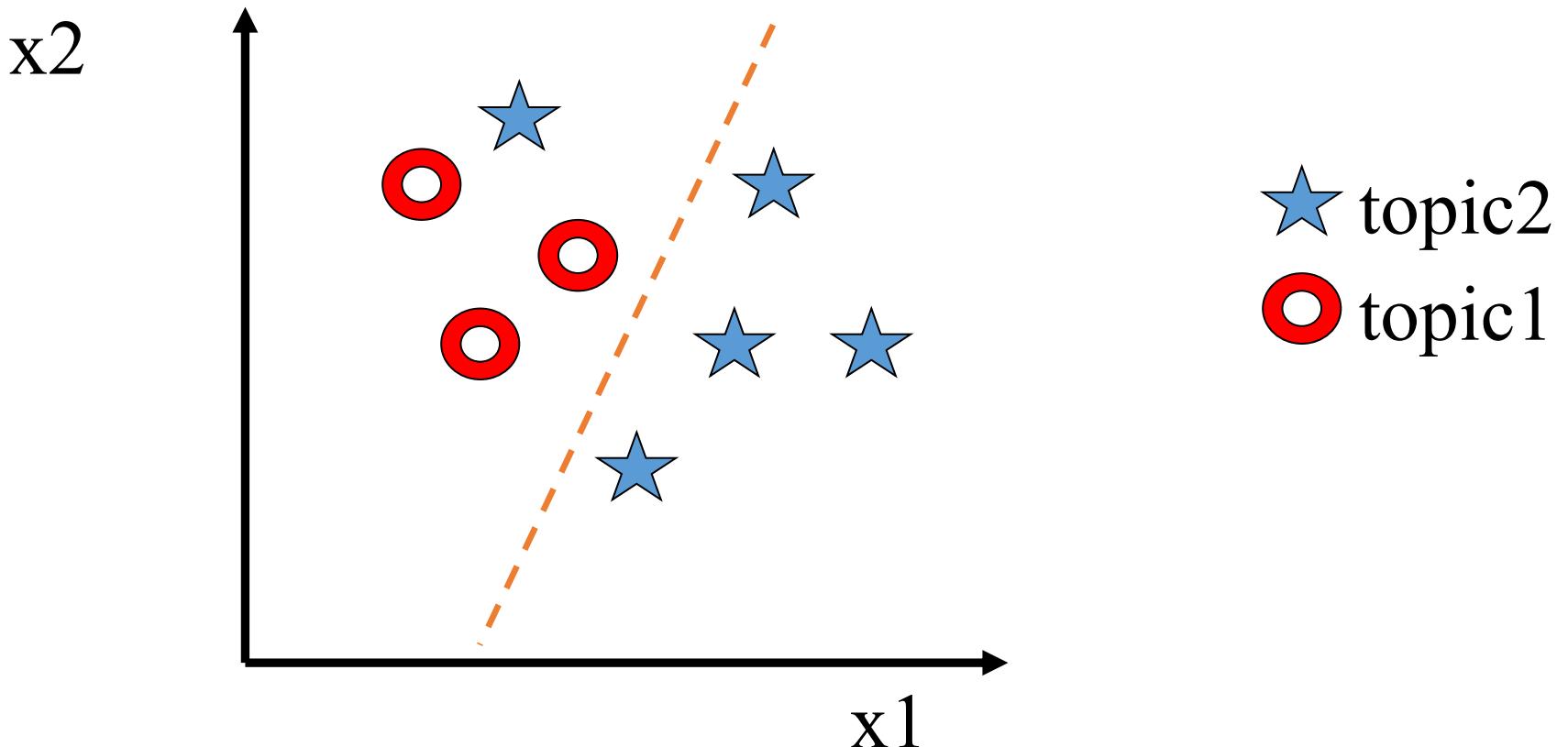
Decision trees



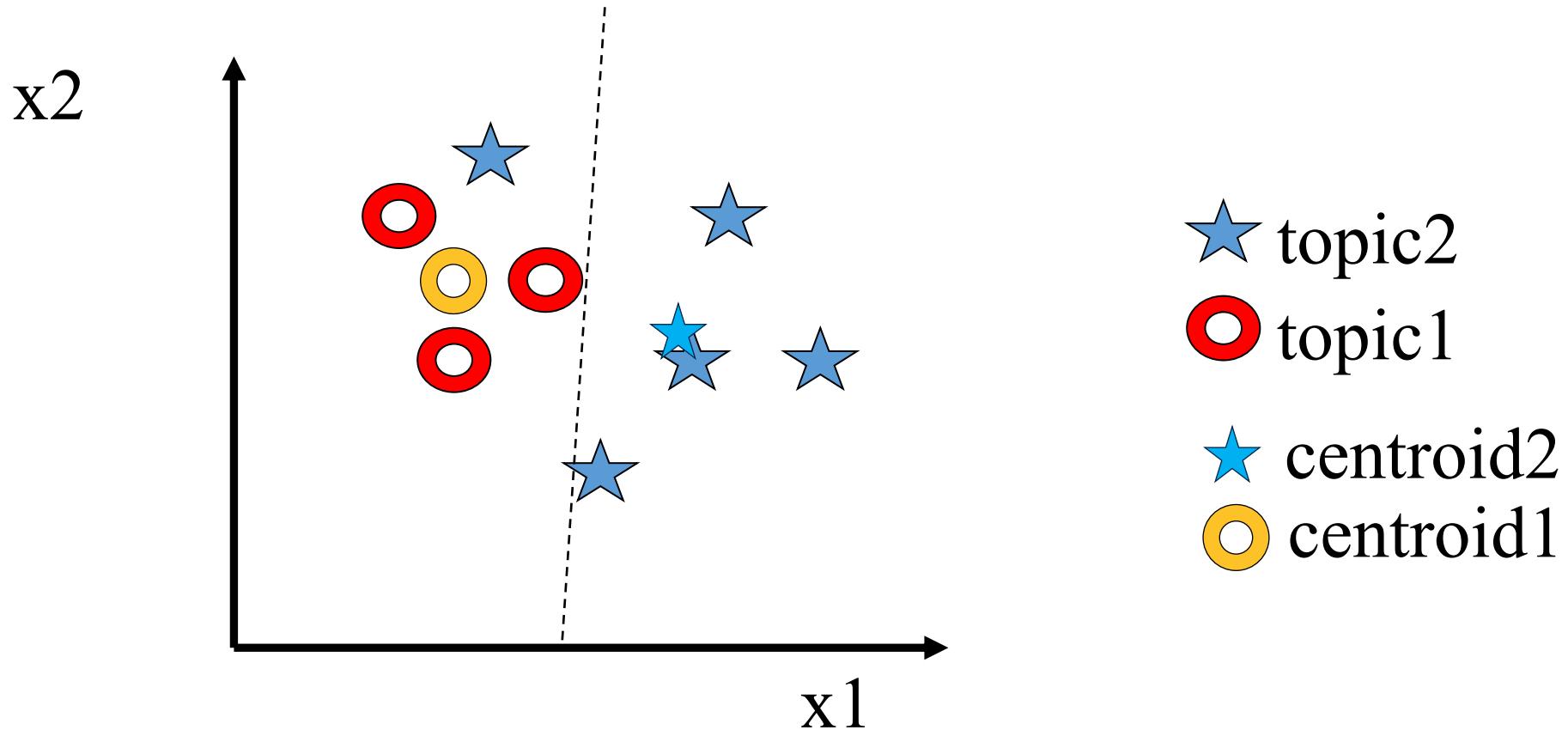
Classification Using Centroids

- Centroid
 - the point most representative of a class
- Compute centroid by finding vector average of known class members
- Decision boundary is a line that is equidistant from two centroids.
- New document on one side of the goes in one class; new document on the other side goes in the other.

Linear boundary



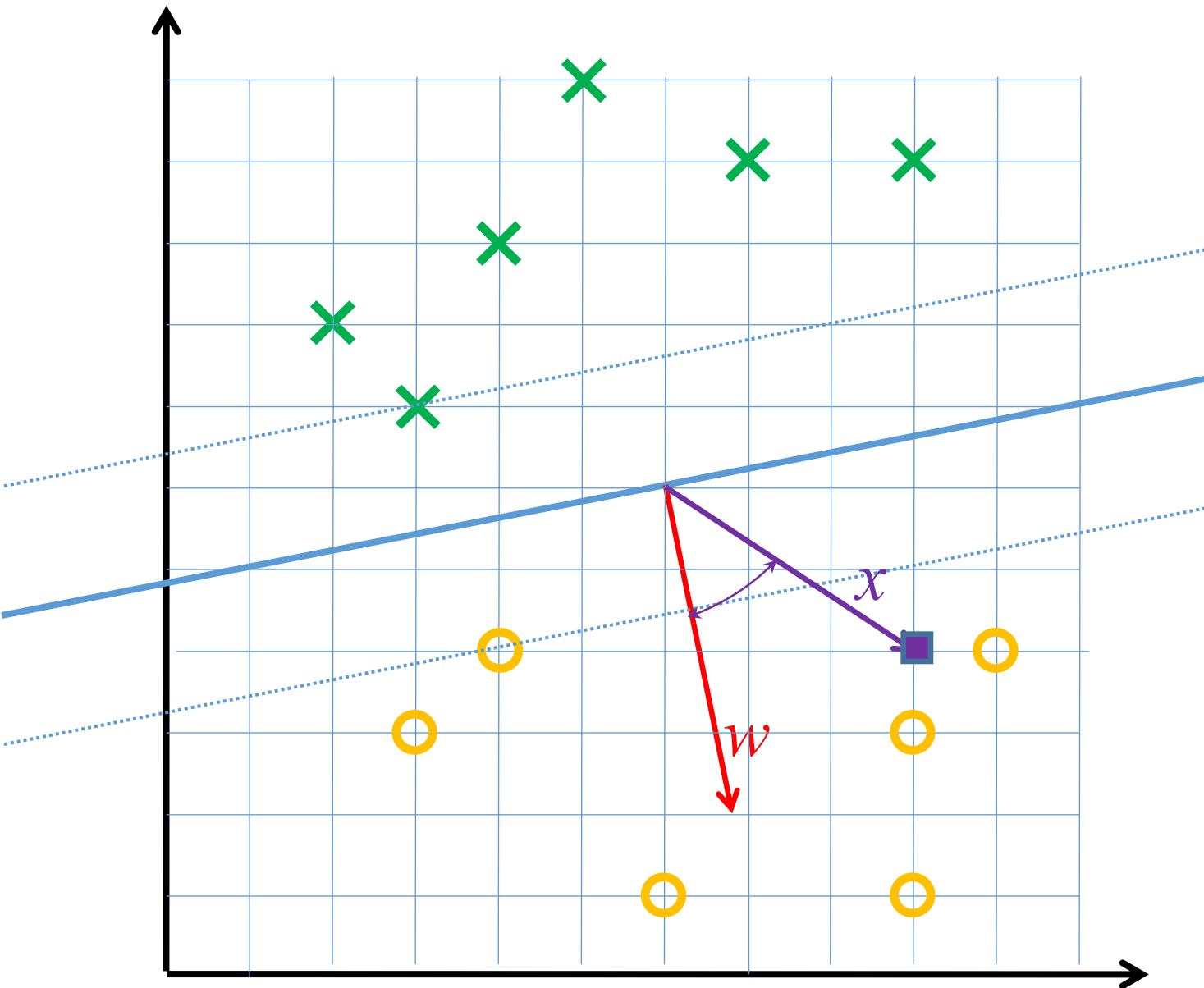
Classification Using Centroids



Introduction to NLP

Linear Classifiers

Decision Boundary



Linear Separators

- Two-dimensional line:

$w_1x_1 + w_2x_2 = b$ is the linear separator

$w_1x_1 + w_2x_2 > b$ for the positive class

- In n -dimensional spaces:

$$\vec{w}^T \vec{x} = \sum_{i=1}^n w_i x_i = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n$$

- One can also add $w_0 = 1$, $x_0 = b$ (constant)
- w is the weight vector
- x is the feature vector

Example

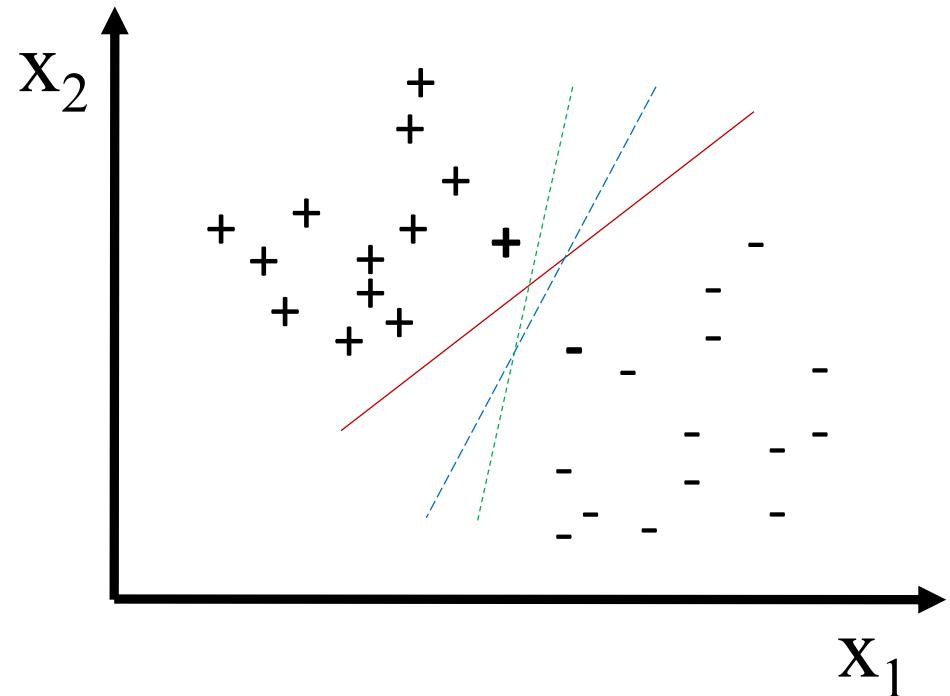
$$\vec{w}^T \vec{x} = b$$

- Bias $b=0$ (in this example)
- Sentence is “A D E H”
- Its score will be
 $0.6*1+0.4*1+0.4*1+(-0.5)*1 = 0.9 > 0$

w _i	x _i	w _i	x _i
0.6	A	-0.7	G
0.5	B	-0.5	H
0.5	C	-0.3	I
0.4	D	-0.2	J
0.4	E	-0.2	K
0.3	F	-0.2	L

How to Find the Linear Boundary?

- Find the linear boundary = find \vec{W}
- Many methods
 - E.g., Perceptron
- Problem:
 - There are infinite number of linear boundaries if the two classes are linearly separable!
 - Maximum margin: Support Vector Machines (SVM)



General Training Idea

- Go through the training data
- Predict the class y (1 or -1)

$$\hat{y} = \arg \max_{y'} \theta^\top f(x, y')$$

- If the prediction is wrong, update ϑ

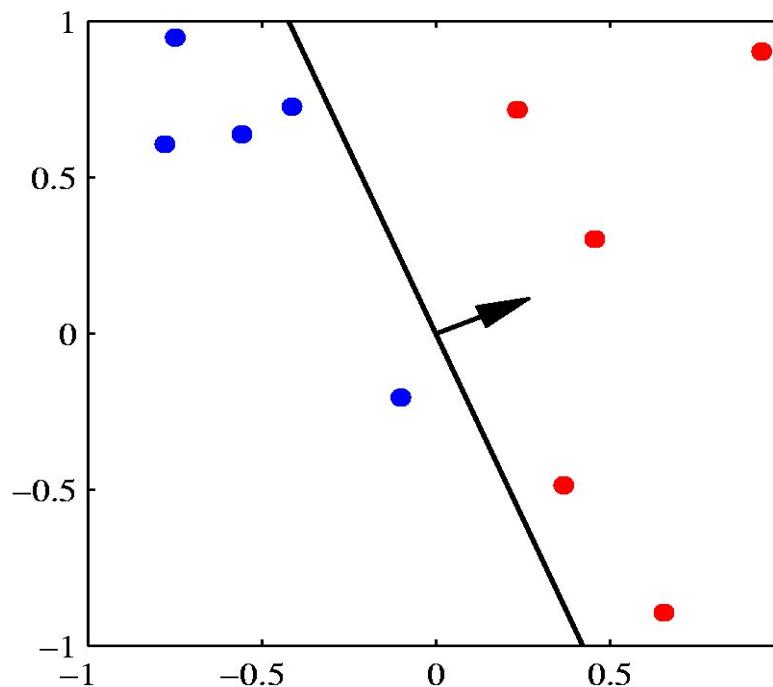
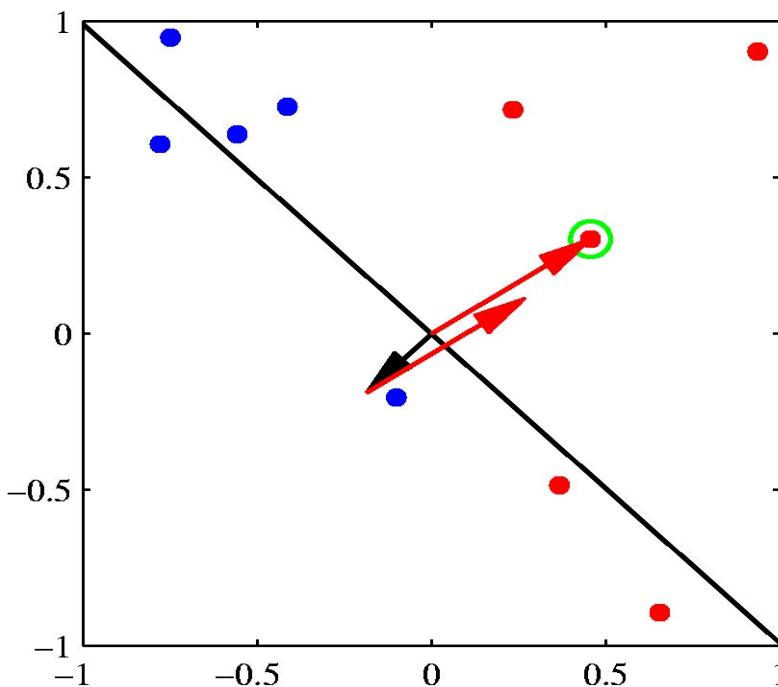
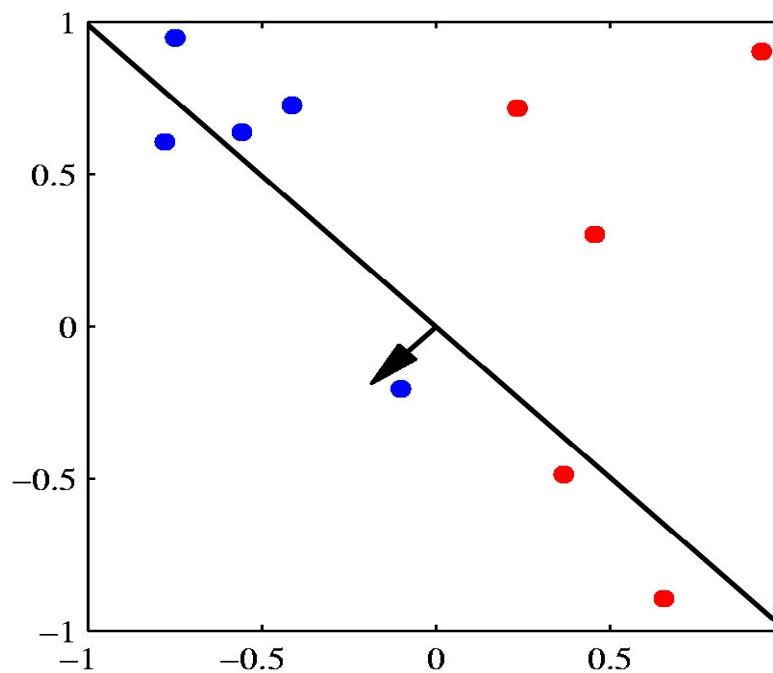
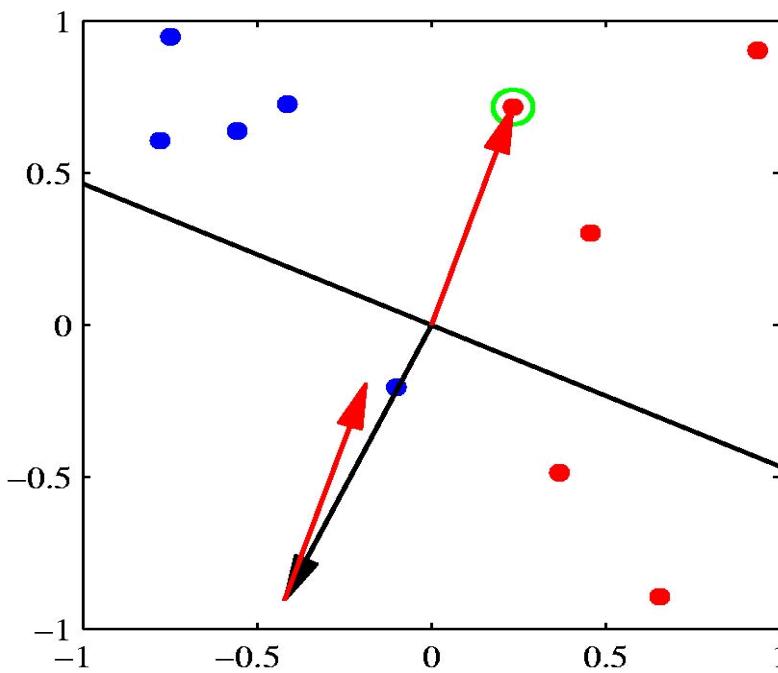
If $\hat{y} \neq y$:

$$\begin{aligned}\theta^{(\text{new})} &\leftarrow \theta^{(\text{old})} + f(x, y) &\rightsquigarrow &\text{ground truth} \\ &- f(x, \hat{y}) &\rightsquigarrow &\text{predicted label}\end{aligned}$$

Perceptron Algorithm (preview)

Algorithm 3 Perceptron learning algorithm

```
1: procedure PERCEPTRON( $\mathbf{x}^{(1:N)}, y^{(1:N)}$ )
2:    $t \leftarrow 0$ 
3:    $\boldsymbol{\theta}^{(0)} \leftarrow \mathbf{0}$ 
4:   repeat
5:      $t \leftarrow t + 1$ 
6:     Select an instance  $i$ 
7:      $\hat{y} \leftarrow \text{argmax}_y \boldsymbol{\theta}^{(t-1)} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y)$ 
8:     if  $\hat{y} \neq y^{(i)}$  then
9:        $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \mathbf{f}(\mathbf{x}^{(i)}, \hat{y})$ 
10:    else
11:       $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)}$ 
12:    until tired
13:    return  $\boldsymbol{\theta}^{(t)}$ 
```



[Example: Chris Bishop]

Interpolation using Features

- Linear interpolation for language modeling
 - Estimating the trigram probability $P(c|ab)$ using $P_{MLE}(c|a)$, $P_{MLE}(c)$, etc.
 - Weights λ_1, λ_2 , etc.
- We may want to consider other features
 - E.g., POS tags of previous words, heads, word endings, etc.
- General idea
 - Compute the conditional probability $P(y|x)$
 - $P(y|x) = \text{sum of weights} * \text{features}$
- Label Y for a given history x in X

Logistic Regression (preview)

- ▶ Model: $P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$
- ▶ Inference: $\operatorname{argmax}_y P_w(y|x)$
- ▶ Learning: gradient ascent on the discriminative log-likelihood
$$f(x, y^*) - \mathbb{E}_y[f(x, y)] = f(x, y^*) - \sum_y [P_w(y|x) f(x, y)]$$
"towards gold feature value, away from expectation of feature value"

[discriminative method]

[Greg Durrett]

Naïve Bayes (preview)

- Multinomial Naïve Bayes is a linear model

$$x = [1, \quad \quad \quad x_1, \quad \quad \quad x_2 \quad \quad \quad \dots]$$

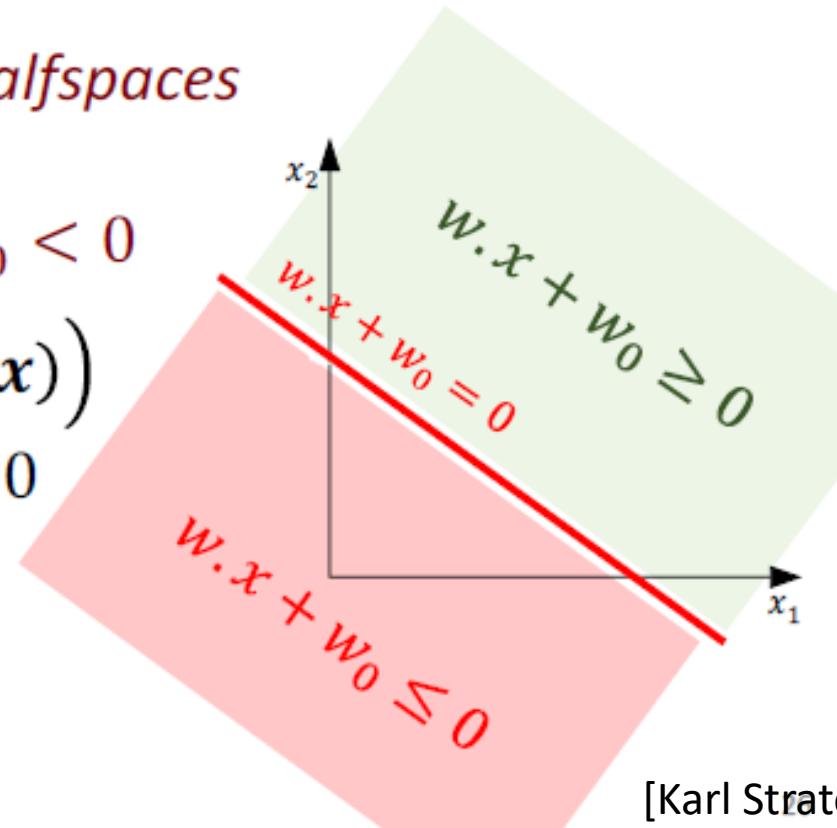
$$w = [\log P(y), \quad \log P(w_1|y), \log P(w_2|y) \dots]$$

Parametric classifiers

- What is the equivalent of linear regression?
 - something easy to train
 - something easy to use at test time
- $f(\mathbf{x}) = f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + w_0$
- $\mathcal{H} = \{f_{\mathbf{w}} = \mathbf{x} \rightarrow \mathbf{w} \cdot \mathbf{x} + w_0 : \mathbf{w} \in \mathbb{R}^d, w_0 \in \mathbb{R}\}$
- but $f(\mathbf{x}) \notin \{-1, 1\}$! how do we get labels?
 - reasonable choice
 - $\hat{y}(\mathbf{x}) = 1 \text{ if } f_{\widehat{\mathbf{w}}}(\mathbf{x}) \geq 0 \text{ and } \hat{y}(\mathbf{x}) = -1 \text{ otherwise}$
 - linear classifier: $\hat{y}(\mathbf{x}) = \text{sign}(\widehat{\mathbf{w}} \cdot \mathbf{x} + \widehat{w}_0)$

Parametric classifiers

- $\mathcal{H} = \{f_w = \mathbf{x} \rightarrow \mathbf{w} \cdot \mathbf{x} + w_0 : \mathbf{w} \in \mathbb{R}^d, w_0 \in \mathbb{R}\}$
- $\hat{y}(\mathbf{x}) = \text{sign}(\hat{\mathbf{w}} \cdot \mathbf{x} + \hat{w}_0)$
- $\hat{\mathbf{w}} \cdot \mathbf{x} + \hat{w}_0 = 0$ (linear) decision boundary or separating *hyperplane*
 - that separates \mathbb{R}^d into two *halfspaces* (*regions*)
 $\hat{\mathbf{w}} \cdot \mathbf{x} + \hat{w}_0 > 0$ and $\hat{\mathbf{w}} \cdot \mathbf{x} + \hat{w}_0 < 0$
- more generally, $\hat{y}(\mathbf{x}) = \text{sign}(\hat{f}(\mathbf{x}))$
→ decision boundary is $\hat{f}(\mathbf{x}) = 0$



Classification vs Regression

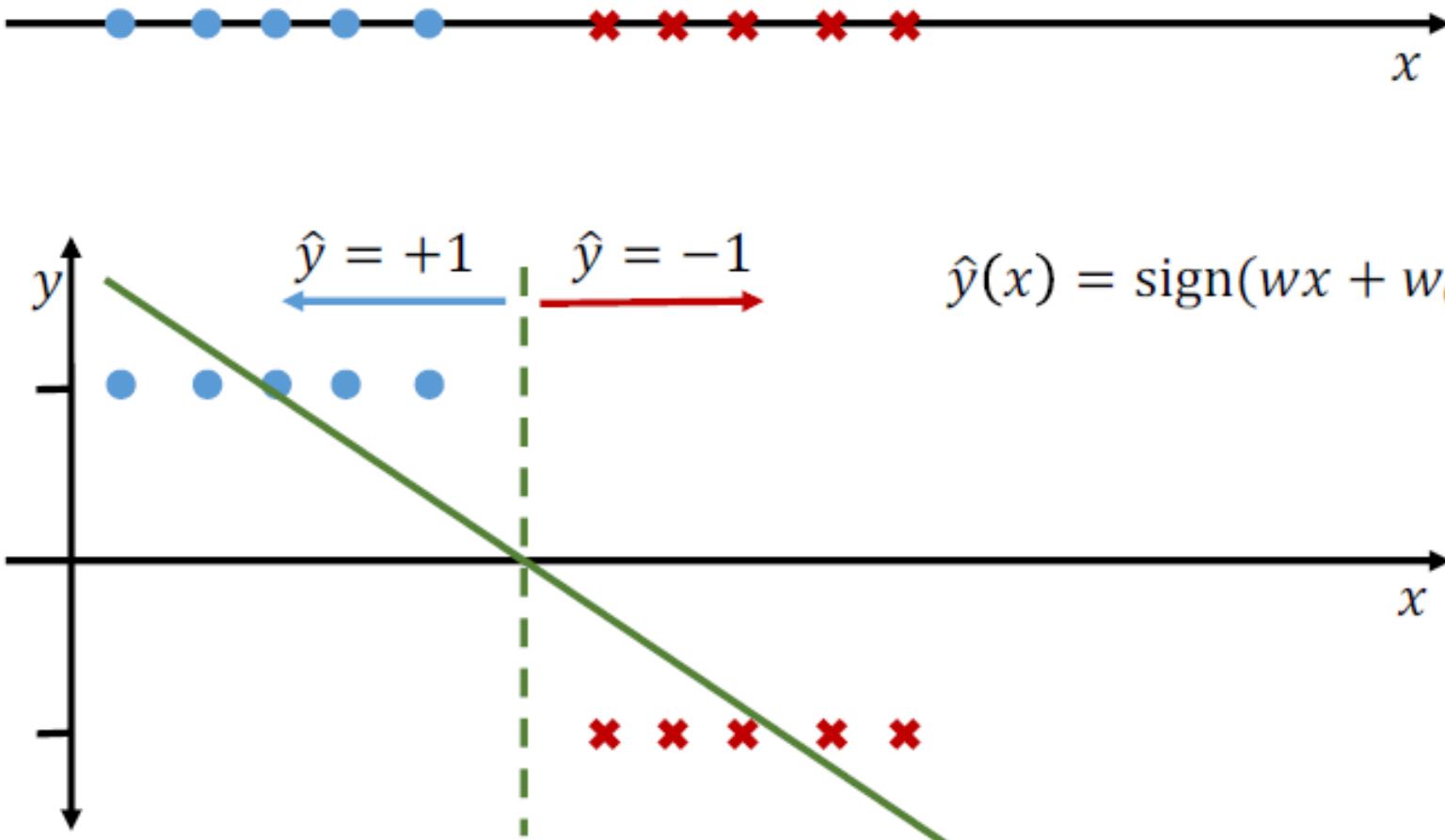
- Label-values do not have meaning
 - $\mathcal{Y} = \{\text{spam, nospam}\}$ or $\mathcal{Y} = \{0,1\}$ or $\mathcal{Y} = \{-1,1\}$
- Ordering of labels does not matter (for most parts)
 - $f(x) = "0"$ when $y = "1"$ is as bad as $f(x) = "9"$ when $y = "1"$
- Often $f(x)$ does not return labels y
 - e.g. in binary classification with $\mathcal{Y} = \{-1,1\}$ we often estimate $f: \mathcal{X} \rightarrow \mathbb{R}$ and then post process to get $\hat{y}(f(x)) = \mathbf{1}[f(x) \geq 0]$
 - mainly for computational reasons
 - remember, we need to solve $\min_{f \in \mathcal{H}} \sum_i \ell(f(x^{(i)}), y^{(i)})$
 - discrete values → combinatorial problems → hard to solve
 - more generally $\mathcal{H} \subset \{f: \mathcal{X} \rightarrow \mathbb{R}\}$ and loss $\ell: \mathbb{R} \times \mathcal{Y} \rightarrow \mathbb{R}$
 - compare to regression, where typically $\mathcal{H} \subset \{f: \mathcal{X} \rightarrow \mathcal{Y}\}$ and loss $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$

Classification as regression

- Binary classification $\mathcal{Y} = \{-1, 1\}$ and $\mathcal{X} \in \mathbb{R}^d$
- Treat it as regression with squared loss, say linear regression
 - Training data $S = \{(\mathbf{x}^{(i)}, y^{(i)}): i = 1, 2, \dots, N\}$
 - ERM

$$\hat{\mathbf{w}}, \hat{w}_0 = \operatorname{argmin}_{w, w_0} \sum_i (\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0 - y^{(i)})^2$$

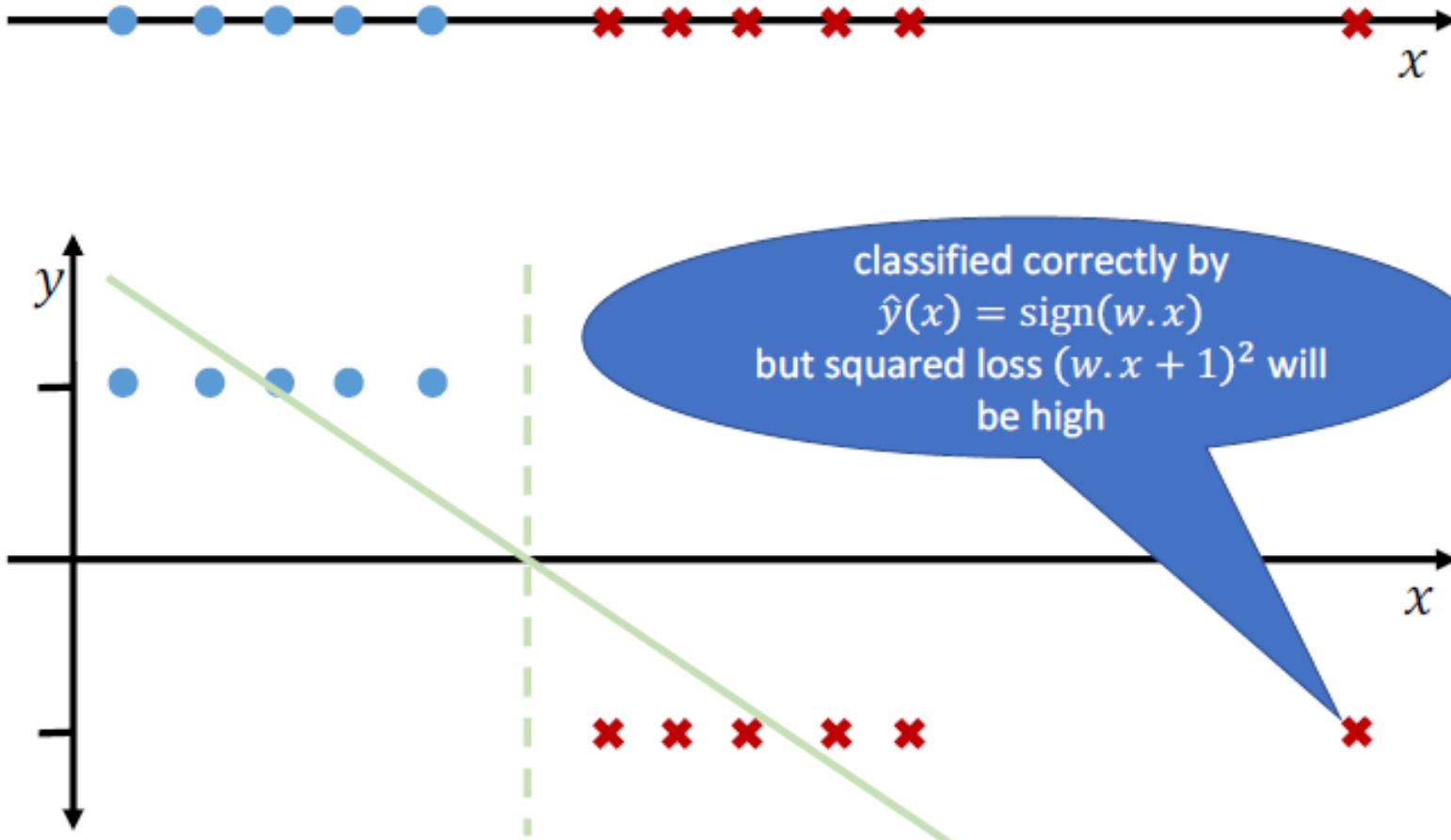
Classification as regression



Example credit: Greg Shaknarovich

[Karl Stratos]

Classification as regression

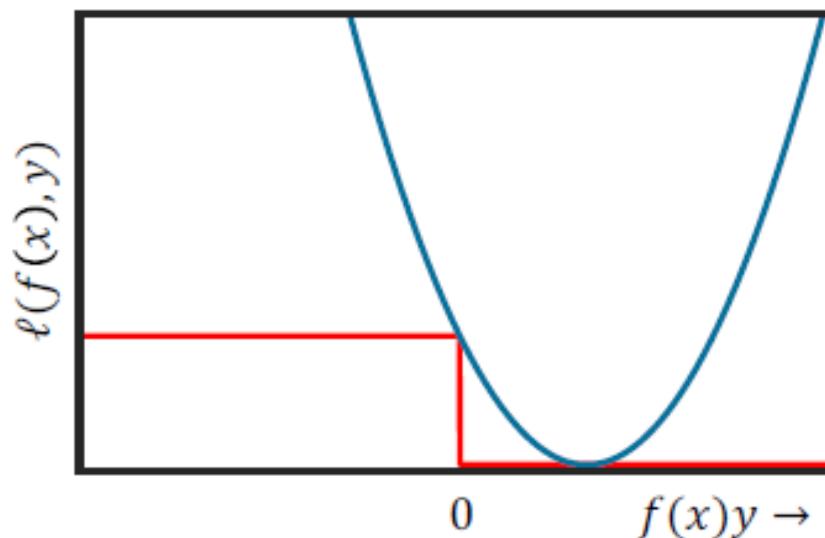


Surrogate Losses

- The correct loss to use is 0-1 loss *after* thresholding

$$\begin{aligned}\ell^{01}(f(x), y) &= \mathbf{1}[\text{sign}(f(x)) \neq y] \\ &= \mathbf{1}[\text{sign}(f(x)y) < 0]\end{aligned}$$

- Linear regression uses $\ell^{LS}(f(x), y) = (f(x) - y)^2$



- Why not do ERM over $\ell^{01}(f(x), y)$ directly?
 - non-continuous, non-convex

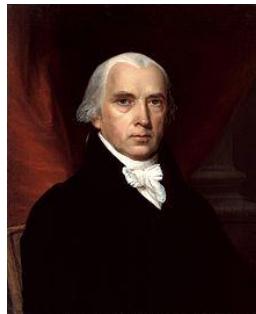
Introduction to NLP

541.

Text Classification

Who wrote which Federalist papers?

- 1787-8: anonymous essays try to convince New York to ratify U.S Constitution: Jay, Madison, Hamilton.
- Authorship of 12 of the letters in dispute
- 1963: solved by Mosteller and Wallace using Bayesian methods



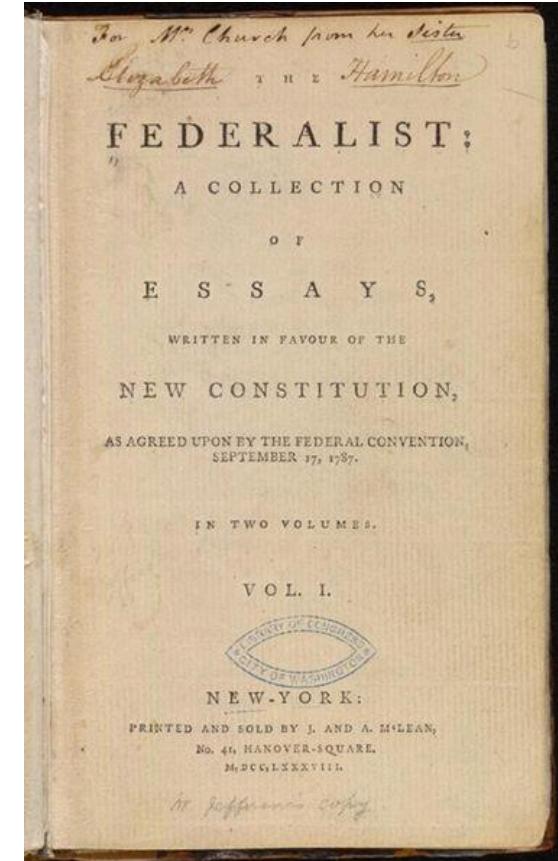
James Madison



Alexander Hamilton



John Jay



Positive or negative movie review?



- unbelievably disappointing



- Full of zany characters and richly applied satire, and some great plot twists



- this is the greatest screwball comedy ever filmed



- It was pathetic. The worst part about it was the boxing scenes.

The bag of words representation

y(

seen	2
sweet	1
whimsical	1
recommend	1
happy	1
...	...

) = c

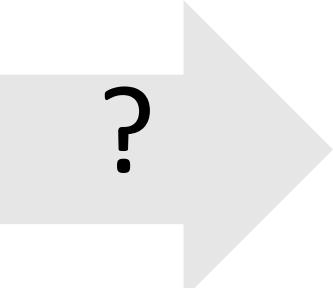


What is the subject of this article?

MEDLINE Article



MeSH Subject Category Hierarchy

- 
- Antagonists and Inhibitors
 - Blood Supply
 - Chemistry
 - Drug Therapy
 - Embryology
 - Epidemiology
 - ...

Introduction to NLP

Evaluation of Classification

Accuracy

- Percentage of items correctly classified
- Not a great metric for imbalanced data sets. Why?

Contingency Table

		<i>gold standard labels</i>		
		gold positive	gold negative	
<i>system output labels</i>	system positive	true positive	false positive	precision = $\frac{tp}{tp+fp}$
	system negative	false negative	true negative	
		recall = $\frac{tp}{tp+fn}$		accuracy = $\frac{tp+tn}{tp+fp+tn+fn}$

Figure 4.4 Contingency table

A combined measure: F

- A combined measure that assesses the P/R tradeoff is F measure (weighted harmonic mean):

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

- The harmonic mean is a very conservative average
- People usually use balanced F1 measure
 - i.e., with $\beta = 1$ (that is, $\alpha = \frac{1}{2}$):

$$F1 = 2PR/(P+R)$$

More than two classes

		gold labels		
		urgent	normal	spam
system output	urgent	8	10	1
	normal	5	60	50
	spam	3	30	200
recall _u =		recall _n =	recall _s =	
$\frac{8}{8+5+3}$		$\frac{60}{10+60+30}$	$\frac{200}{1+50+200}$	

precision_u = $\frac{8}{8+10+1}$
precision_n = $\frac{60}{5+60+50}$
precision_s = $\frac{200}{3+30+200}$

Figure 4.5 Confusion matrix for a three-class categorization task, showing for each pair of classes (c_1, c_2) , how many documents from c_1 were (in)correctly assigned to c_2

Micro-averaging and Macro-averaging

If we have more than one class, how do we combine multiple performance measures into one quantity?

- **Macroaveraging:** Compute performance for each class, then average.
- **Microaveraging:** Collect decisions for all classes, compute contingency table, evaluate.

Micro-averaging and Macro-averaging

Class 1: Urgent		Class 2: Normal		Class 3: Spam		Pooled	
true	true	true	true	true	true	true	true
urgent	not	normal	not	spam	not	yes	no
system	urgent	8	11	60	55	200	33
system	not	8	340	40	212	51	83

precision = $\frac{8}{8+11} = .42$ precision = $\frac{60}{60+55} = .52$ precision = $\frac{200}{200+33} = .86$ microaverage precision = $\frac{268}{268+99} = .73$

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

Figure 4.6 Separate contingency tables for the 3 classes from the previous figure, showing the pooled contingency table and the microaveraged and macroaveraged precision.

Classic Reuters-21578 Data Set

- Most (over)used data set, 21,578 docs (each 90 types, 200 tokens)
- 9603 training, 3299 test articles (ModApte/Lewis split)
- 118 categories
 - An article can be in more than one category
 - Learn 118 binary category distinctions
- Average document (with at least one category) has 1.24 classes
- Only about 10 out of 118 categories are large

Common categories
(#train, #test)

- Earn (2877, 1087)
- Acquisitions (1650, 179)
- Money-fx (538, 179)
- Grain (433, 149)
- Crude (389, 189)
- Trade (369, 119)
- Interest (347, 131)
- Ship (197, 89)
- Wheat (212, 71)
- Corn (182, 56)

Reuters-21578 document

<REUTERS TOPICS="YES" LEWISSPLIT="TRAIN" CGISPLIT="TRAINING-SET" OLDID="12981" NEWID="798">

<DATE> 2-MAR-1987 16:51:43.42</DATE>

<TOPICS><D>livestock</D><D>hog</D></TOPICS>

<TITLE>AMERICAN PORK CONGRESS KICKS OFF TOMORROW</TITLE>

<DATELINE> CHICAGO, March 2 - </DATELINE><BODY>The American Pork Congress kicks off tomorrow, March 3, in Indianapolis with 160 of the nations pork producers from 44 member states determining industry positions on a number of issues, according to the National Pork Producers Council, NPPC.

Delegates to the three day Congress will be considering 26 resolutions concerning various issues, including the future direction of farm policy and the tax law as it applies to the agriculture sector. The delegates will also debate whether to endorse concepts of a national PRV (pseudorabies virus) control and eradication program, the NPPC said.

A large trade show, in conjunction with the congress, will feature the latest in technology in all areas of the industry, the NPPC added. Reuter

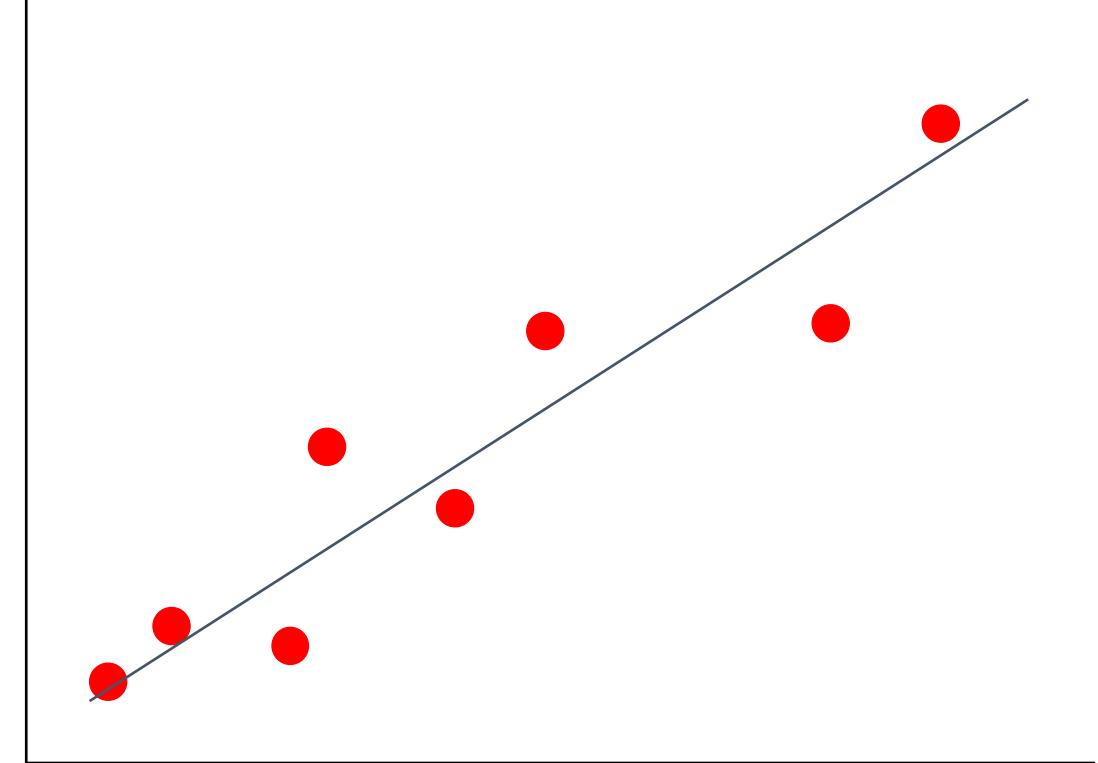
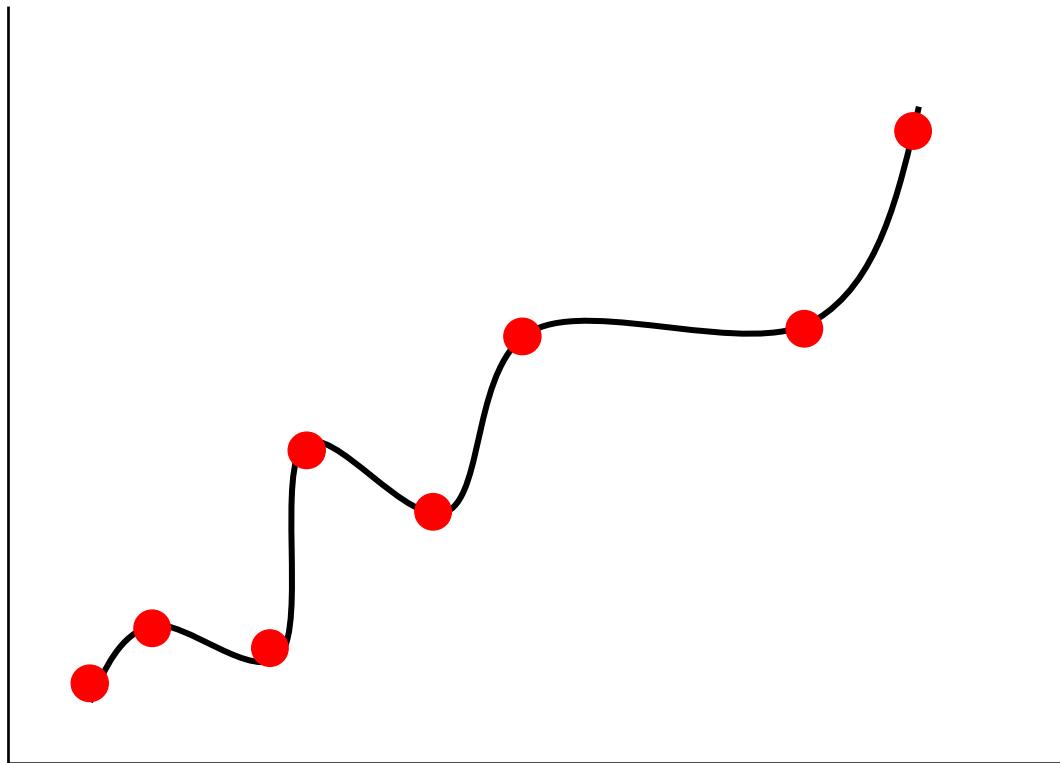
</BODY></TEXT></REUTERS>

Confusion Matrix

- For each pair of classes $\langle c_1, c_2 \rangle$ how many documents from c_1 were incorrectly assigned to c_2 ?
 - $c_{3,2}$: 90 wheat documents incorrectly assigned to poultry

Docs in test set	Assigned UK	Assigned poultry	Assigned wheat	Assigned coffee	Assigned interest	Assigned trade
True UK	95	1	13	0	1	0
True poultry	0	1	0	0	0	0
True wheat	10	90	0	1	0	0
True coffee	0	0	0	34	3	7
True interest	-	1	2	13	26	5
True trade	0	0	2	14	5	10

Pitfall: Overfitting



What happens when your model learns your training data a little too well?

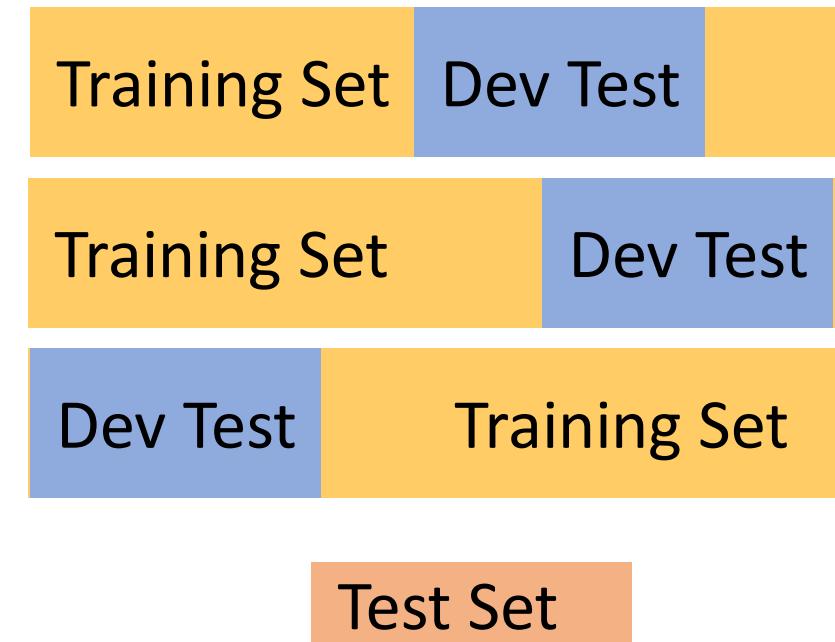
Development Test Sets and Cross-validation

Training set

Development Test Set

Test Set

- Metric: P/R/F1 or Accuracy
- Unseen test set
 - avoid overfitting ('tuning to the test set')
 - more conservative estimate of performance
- Cross-validation over multiple splits
 - Handle sampling errors from different datasets
 - Pool results over each split
 - Compute pooled dev set performance



Cross-validation

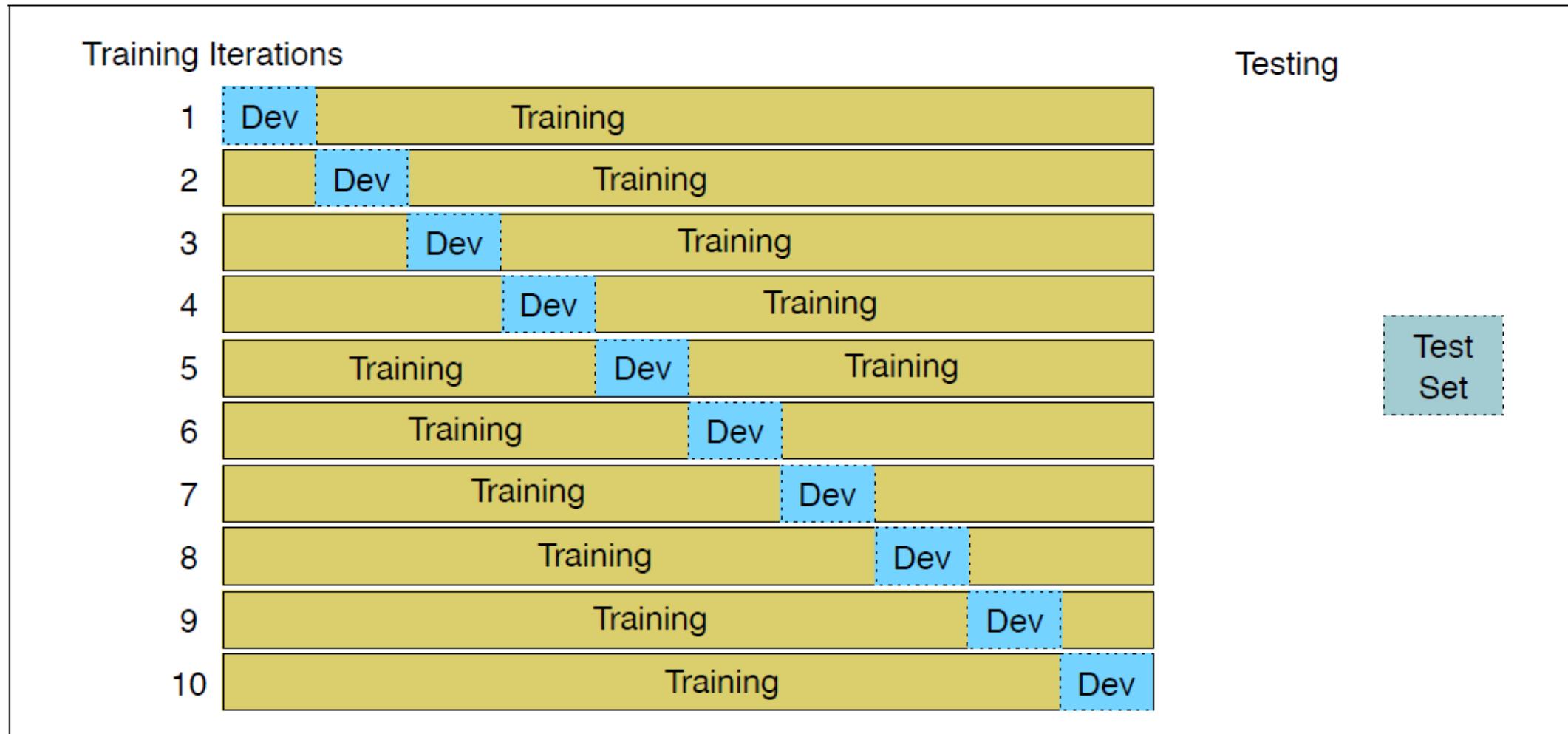


Figure 4.7 10-fold cross-validation

Introduction to NLP

516.

Logistic Regression

Combining Features

- Linear interpolation for language modeling
 - Estimating the trigram probability $P(c|ab)$?
 - features: $P_{MLE}(c|a)$, $P_{MLE}(c)$, etc.
 - Weights: λ_1, λ_2 , etc.
- We may want to consider other features
 - E.g., POS tags of previous words, heads, word endings, etc.
- General idea
 - Compute the conditional probability $P(y|x)$
 - $P(y|x) = \text{sum of weights} * \text{features}$
- Label Y for a given history x in X

Logistic Regression

- Similar to Naïve Bayes (but discriminative!)
 - Log-linear model
 - Features don't have to be independent
- Examples of features
 - Anything of use
 - Linguistic and non-linguistic
 - Count of “good”
 - Count of “not good”
 - Sentence length
- Convex optimization problem
 - It can be solved using gradient descent

Feature Templates

- For each word $\$w$
 - $\$w_count$
 - $\$w_ends_in_er$
 - $\$w_is_negated$
 - $\$w_is_all_caps$

Logistic Regression

- An example of a discriminative classifier
- Input:
 - Training example pairs of (\vec{x}, y) where \vec{x} is the feature vector and y is the label
- Goal:
 - Build a model that predicts the probability of the label
- Output:
 - Set of weights \vec{w} that maximizes likelihood of correct labels on training examples

Loglinear Models

- Naive Bayes, HMM, and Logistic Regression can all be expressed in a similar way:

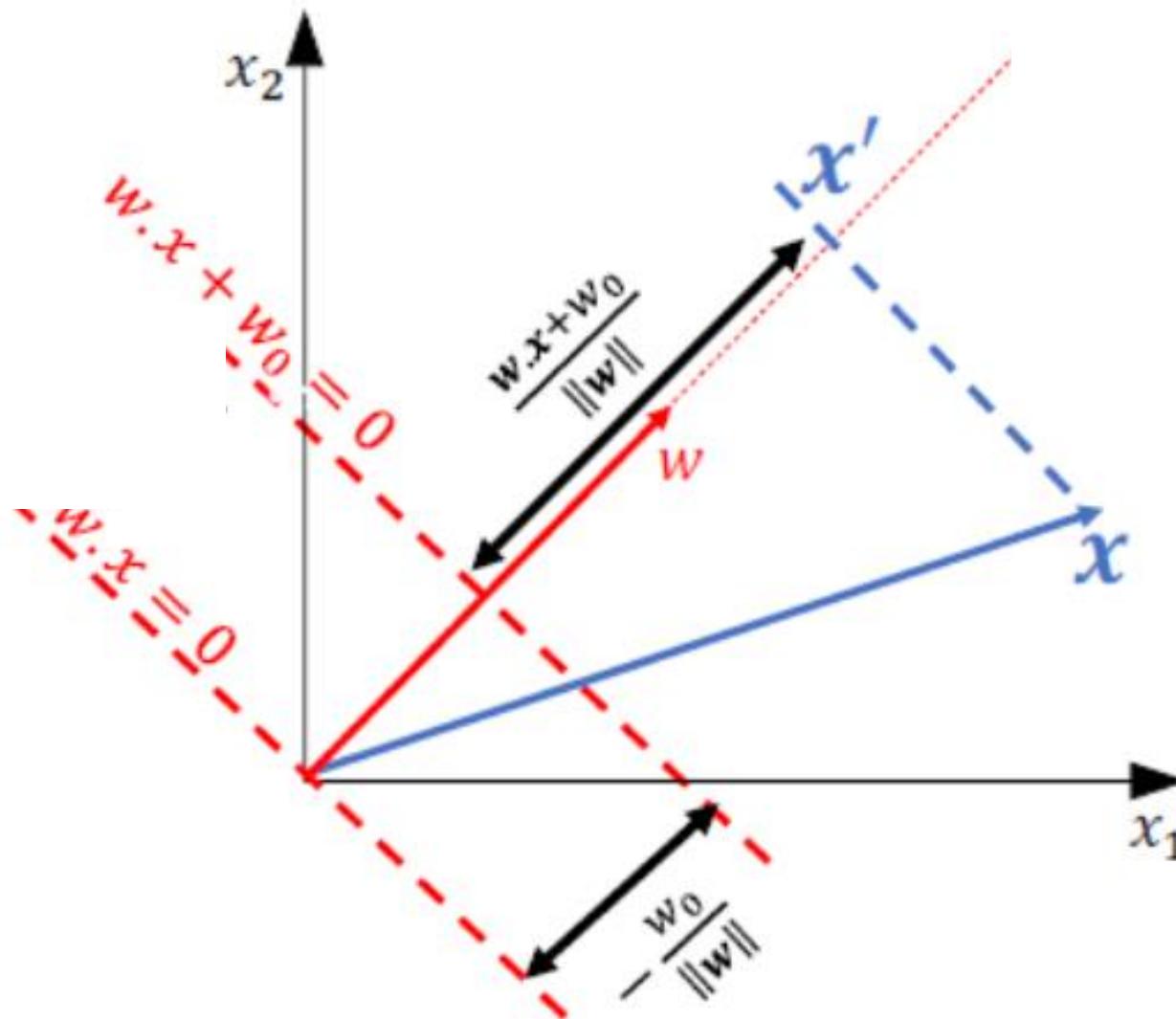
$$G(y, \theta) = \theta^T f(x, y)$$

$$p(y|x, \theta) \propto \exp G(y, \theta) \Leftrightarrow \log p(y|x) = C + G(y, \theta)$$

- Decision rule:

$$\operatorname{argmax}_{y^*} G(y^*, \theta)$$

Mapping x to a 1-D coordinate



[image from Greg Shakhnarovich]

Mapping x to a 1-D coordinate

$$\hat{w}, \hat{w}_0 = \operatorname{argmin}_{w, w_0} \sum_i \log \left(1 + \exp(-(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)y^{(i)}) \right)$$

- Learns a linear decision boundary
 - $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 = 0\}$ is a hyperplane in \mathbb{R}^d - decision boundary
 - $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 = 0\}$ divides \mathbb{R}^d into two halfspace (regions)
 - $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 \geq 0\}$ will get label +1 and
 $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 < 0\}$ will get label -1
- Maps \mathbf{x} to a 1D coordinate

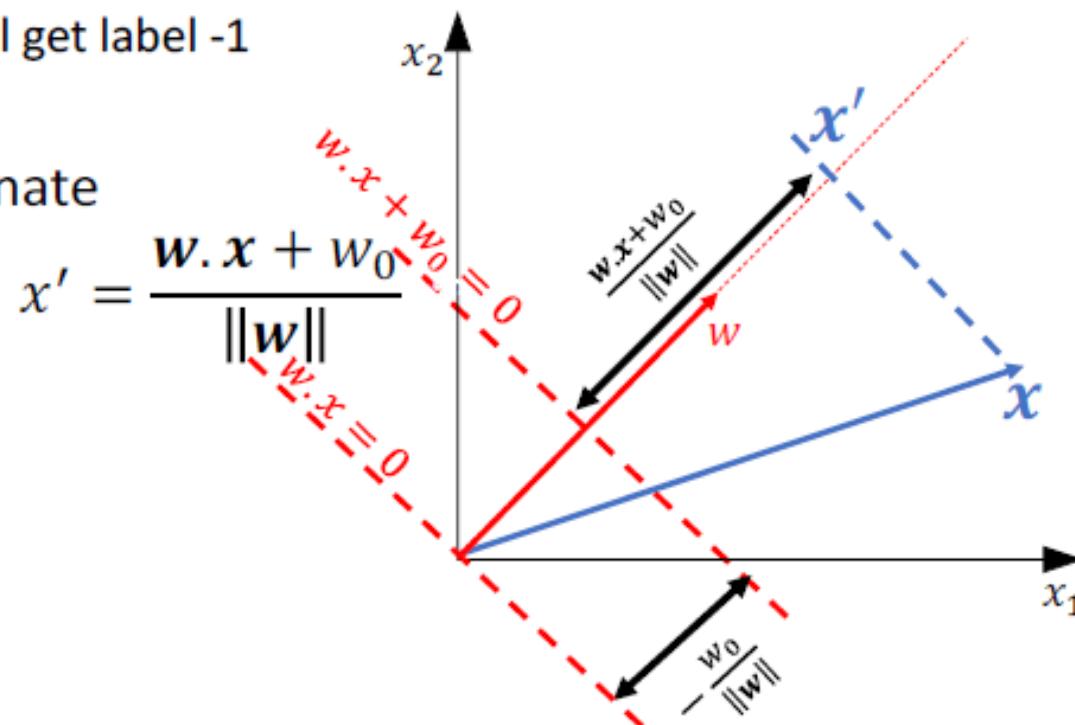


Figure credit: Greg Shakhnarovich

Logistic Regression

- Model: $P_w(y|x) = \frac{\exp(w^\top f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^\top f(x, y'))}$
- Inference: $\operatorname{argmax}_y P_w(y|x)$
- Learning: gradient ascent on the discriminative log-likelihood

$$f(x, y^*) - \mathbb{E}_y[f(x, y)] = f(x, y^*) - \sum_y [P_w(y|x)f(x, y)]$$

“towards gold feature value, away from expectation of feature value”

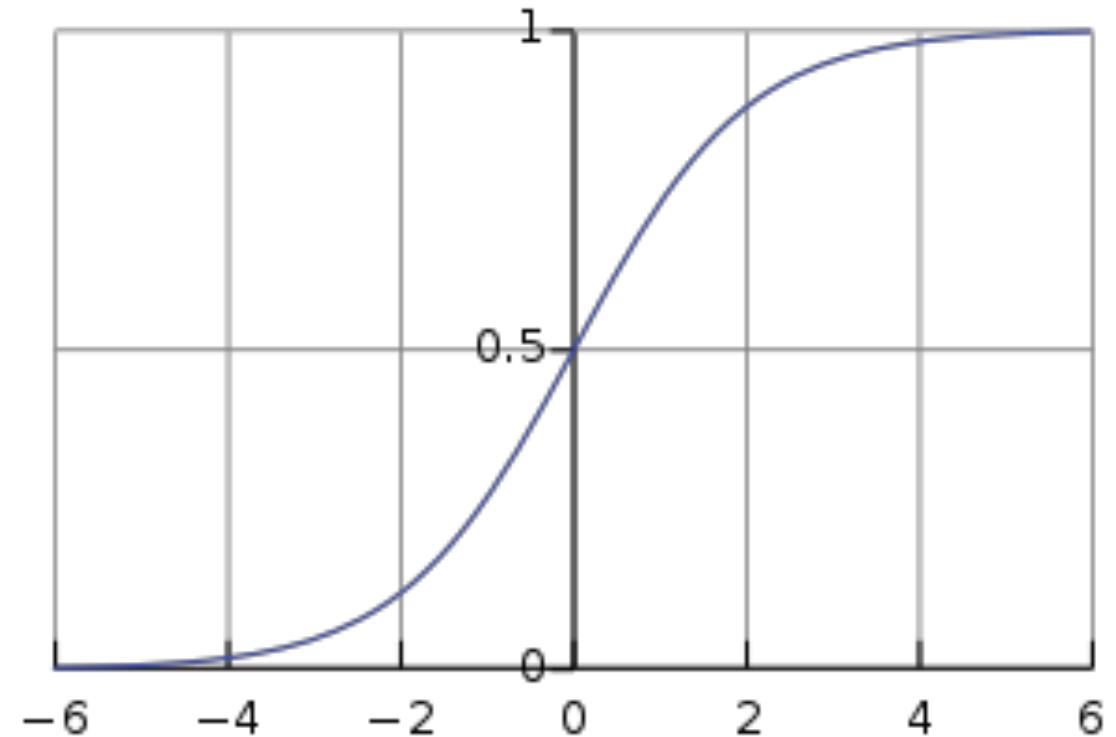
Including a Nonlinearity

- Compute the feature vector x
- Multiply with weight vector w

$$z = \sum w_i x_i$$

- Compute the logistic function (sigmoid)

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



Sigmoid Function

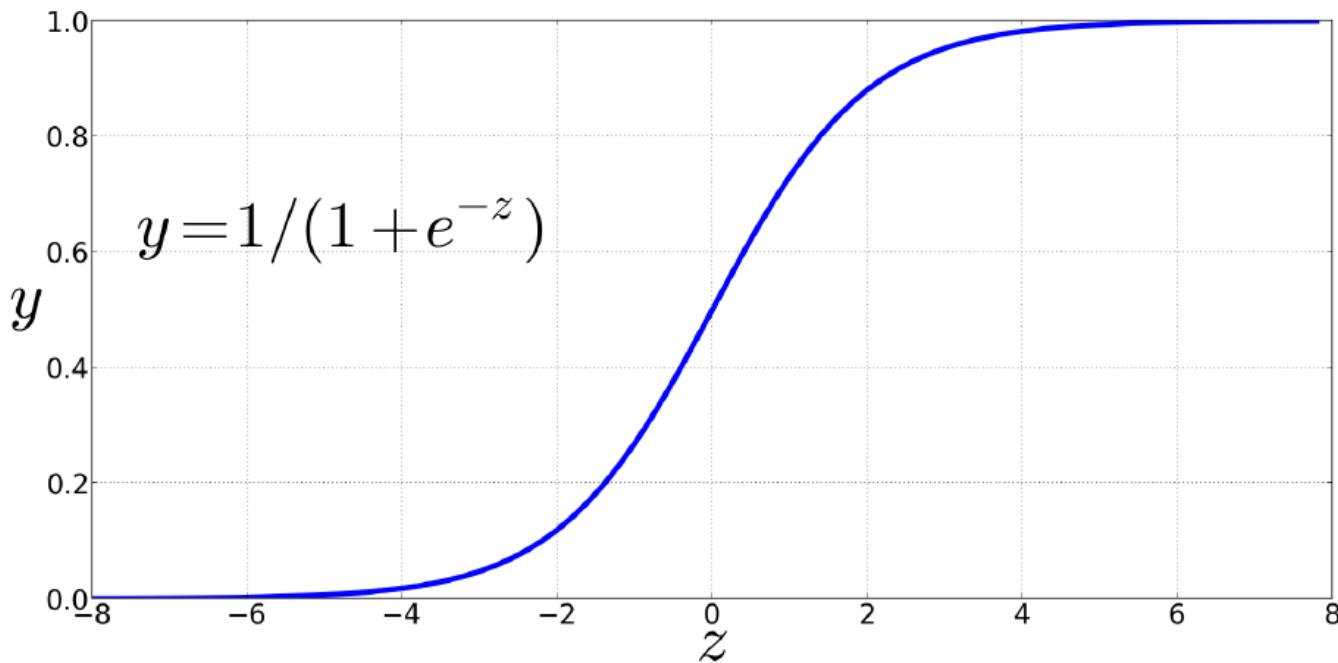


Figure 5.1 The sigmoid function $y = \frac{1}{1+e^{-z}}$ takes a real value and maps it to the range $[0, 1]$. It is nearly linear around 0 but outlier values get squashed toward 0 or 1.

Non-linear Activation Functions

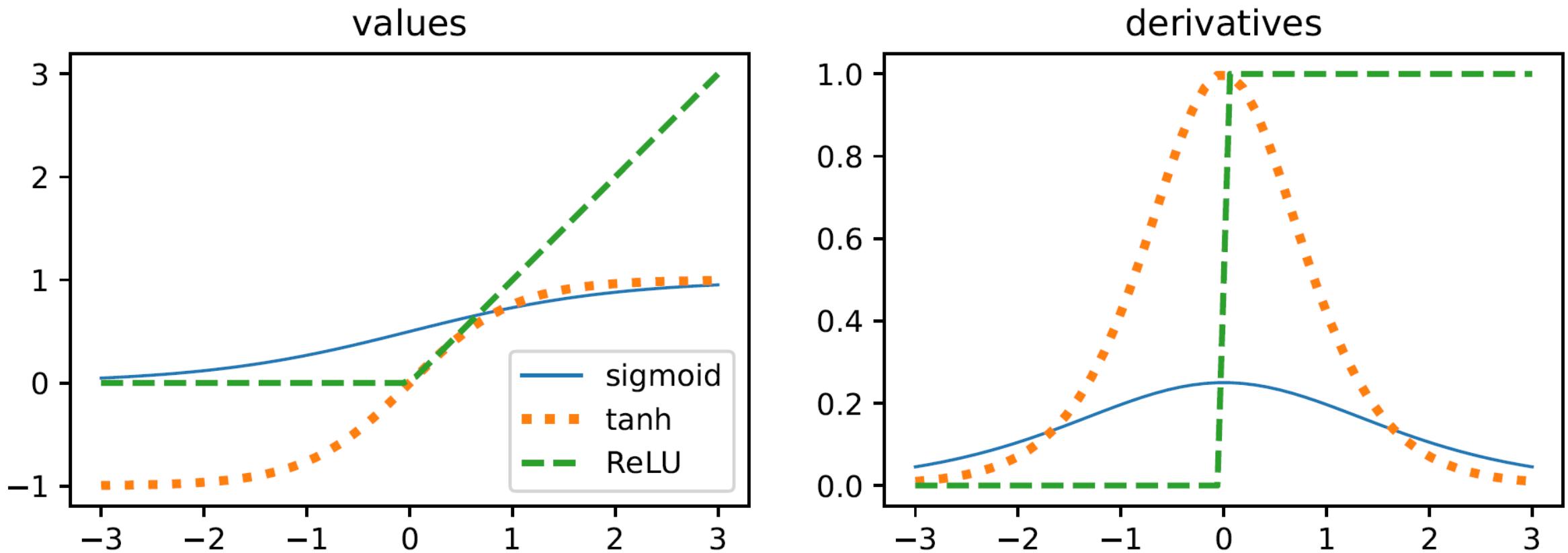


Figure 3.2: The sigmoid, tanh, and ReLU activation functions

Training and Testing LR

Logistic regression has two phases:

training: we train the system (specifically the weights w and b) using stochastic gradient descent and the cross-entropy loss.

test: Given a test example x we compute $p(y|x)$ and return the higher probability label $y = 1$ or $y = 0$.

$$z = \left(\sum_{i=1}^n w_i x_i \right) + b$$

Classification using LR

$$\begin{aligned}y &= \sigma(z) = \frac{1}{1 + e^{-z}} \\P(y = 1) &= \sigma(w \cdot x + b) \\&= \frac{1}{1 + e^{-(w \cdot x + b)}}\end{aligned}$$

$$\hat{y} = \begin{cases} 1 & \text{if } P(y = 1|x) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned}P(y = 0) &= 1 - \sigma(w \cdot x + b) \quad (5.5) \\&= 1 - \frac{1}{1 + e^{-(w \cdot x + b)}} \\&= \frac{e^{-(w \cdot x + b)}}{1 + e^{-(w \cdot x + b)}}\end{aligned}$$

Examples

- Example 1

$$x = (2, 1, 1, 1)$$

$$w = (1, -1, -2, 3)$$

$$z = 2 - 1 - 2 + 3 = 2$$

$$f(z) = 1/(1+e^{-2})$$

- Example 2

$$x = (2, 1, 0, 1)$$

$$w = (0, 0, -3, 0)$$

$$z = 0$$

$$f(z) = 1/(1+e^0) = 1/2$$

Example

Var	Definition	Value in Fig. 5.2
x_1	$\text{count}(\text{positive lexicon}) \in \text{doc}$)	3
x_2	$\text{count}(\text{negative lexicon}) \in \text{doc})$	2
x_3	$\begin{cases} 1 & \text{if “no”} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	1
x_4	$\text{count}(1\text{st and 2nd pronouns} \in \text{doc})$	3
x_5	$\begin{cases} 1 & \text{if “!”} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	0
x_6	$\log(\text{word count of doc})$	$\ln(66) = 4.19$

Example

It's hokey. There are virtually no surprises , and the writing is second-rate .
So why was it so enjoyable ? For one thing , the cast is
great . Another nice touch is the music . I was overcome with the urge to get off
the couch and start dancing . It sucked me in , and it'll do the same to you .

$x_1=3$ $x_2=2$ $x_3=1$ $x_4=3$ $x_5=0$ $x_6=4.19$

Figure 5.2 A sample mini test document showing the extracted features in the vector x .

Per-class Probabilities

$$\begin{aligned} p(+|x) = P(Y = 1|x) &= \sigma(w \cdot x + b) \\ &= \sigma([2.5, -5.0, -1.2, 0.5, 2.0, 0.7] \cdot [3, 2, 1, 3, 0, 4.19] + 0.1) \\ &= \sigma(.833) \\ &= 0.70 \end{aligned} \tag{5.6}$$

$$\begin{aligned} p(-|x) = P(Y = 0|x) &= 1 - \sigma(w \cdot x + b) \\ &= 0.30 \end{aligned}$$

Period Disambiguation

$$x_1 = \begin{cases} 1 & \text{if } \text{"Case}(w_i) = \text{Lower"} \\ 0 & \text{otherwise} \end{cases}$$

$$x_2 = \begin{cases} 1 & \text{if } \text{"}w_i \in \text{AcronymDict}\text{"} \\ 0 & \text{otherwise} \end{cases}$$

$$x_3 = \begin{cases} 1 & \text{if } \text{"}w_i = \text{St.} \& \text{Case}(w_{i-1}) = \text{Cap}\text{"} \\ 0 & \text{otherwise} \end{cases}$$

Cross-Entropy Loss (5.9)

$L(\hat{y}, y) =$ How much \hat{y} differs from the true y

$$p(y|x) = \hat{y}^y (1 - \hat{y})^{1-y}$$

$$\begin{aligned}\log p(y|x) &= \log [\hat{y}^y (1 - \hat{y})^{1-y}] \\ &= y \log \hat{y} + (1 - y) \log (1 - \hat{y})\end{aligned}\tag{5.9}$$

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

$$L_{CE}(w, b) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log (1 - \sigma(w \cdot x + b))] \tag{5.11}$$

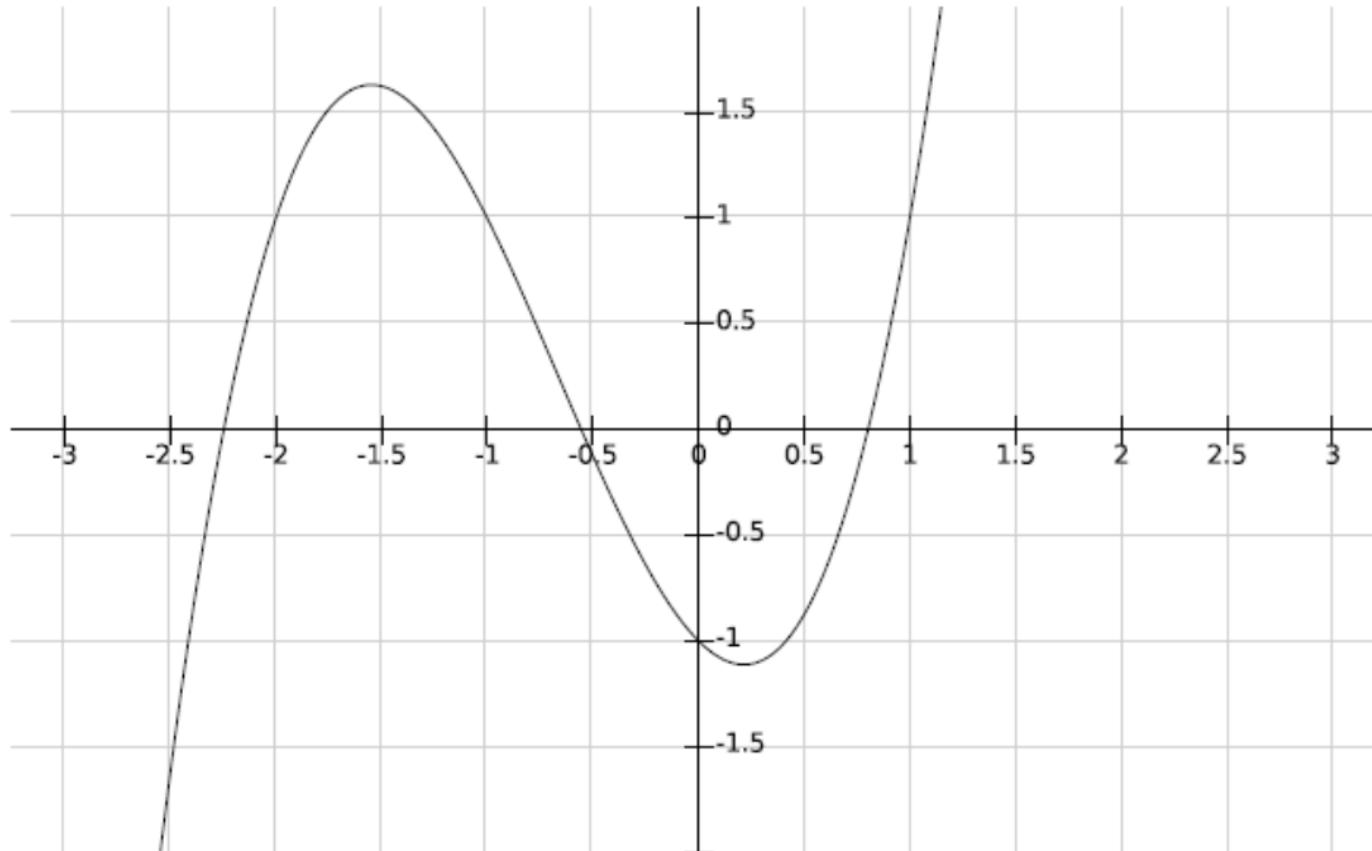
Two examples (.69 vs .31)

$$\begin{aligned} L_{CE}(w, b) &= -[y \log \sigma(w \cdot x + b) + (1 - y) \log (1 - \sigma(w \cdot x + b))] \\ &= -[\log \sigma(w \cdot x + b)] \\ &= -\log(.69) \\ &= .37 \end{aligned}$$

$$\begin{aligned} L_{CE}(w, b) &= -[y \log \sigma(w \cdot x + b) + (1 - y) \log (1 - \sigma(w \cdot x + b))] \\ &= -[\log (1 - \sigma(w \cdot x + b))] \\ &= -\log (.31) \\ &= 1.17 \end{aligned}$$

Function derivatives (refresher)

Minimize $f(x) = x^3 + 2x^2 - x - 1$ over x



$$ax^2 + bx + c = 0$$

$$a = 3$$

$$b = 4$$

$$c = -1$$

$$x = -\frac{2}{3} \pm \frac{\sqrt{7}}{3}$$

$$x = 0.21525$$

$$x = -1.54858$$

(Courtesy to FooPlot)

Gradient Descent

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

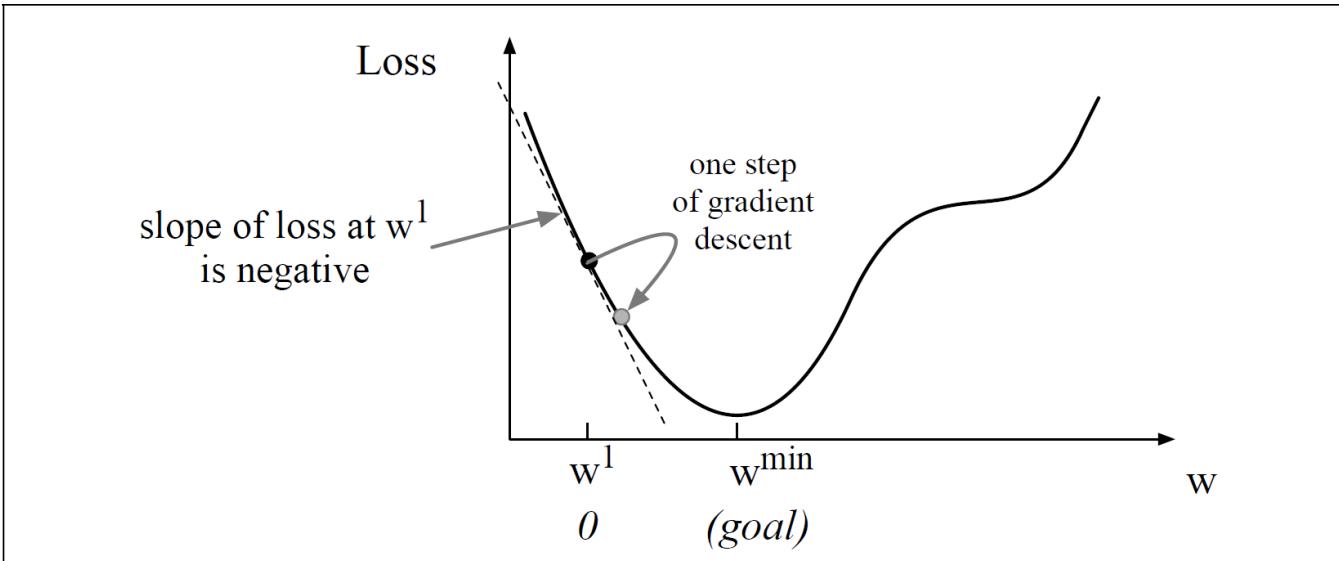


Figure 5.3 The first step in iteratively finding the minimum of this loss function, by moving w in the reverse direction from the slope of the function. Since the slope is negative, we need to move w in a positive direction, to the right. Here superscripts are used for learning steps, so w^1 means the initial value of w (which is 0), w^2 at the second step, and so on.

$$w^{t+1} = w^t - \eta \frac{d}{dw} f(x; w)$$

Bivariate Gradient

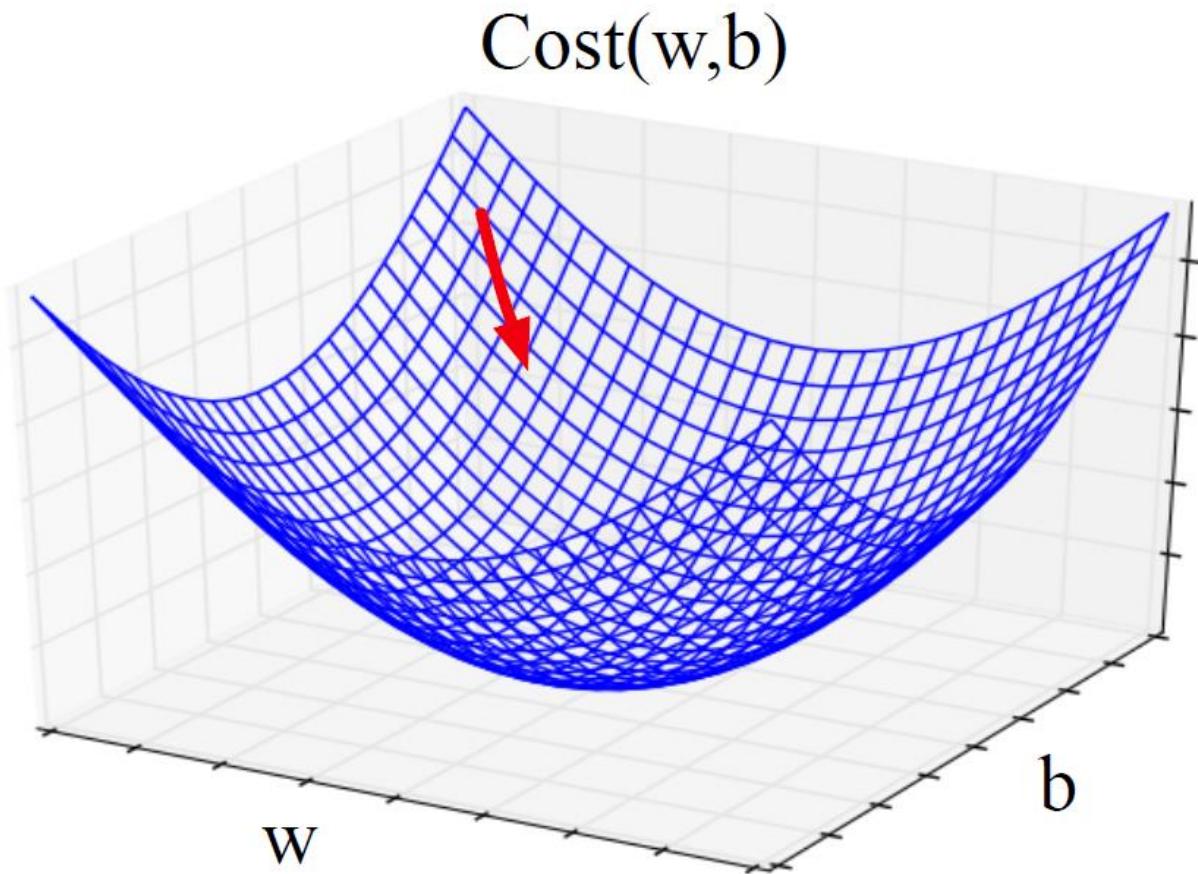


Figure 5.4 Visualization of the gradient vector in two dimensions w and b .

Using the Gradient

In an actual logistic regression, the parameter vector w is much longer than 1 or 2, since the input feature vector x can be quite long, and we need a weight w_i for each x_i . For each dimension/variable w_i in w (plus the bias b), the gradient will have a component that tells us the slope with respect to that variable. Essentially we're asking: "How much would a small change in that variable w_i influence the total loss function L ?"

In each dimension w_i , we express the slope as a partial derivative $\frac{\partial}{\partial w_i}$ of the loss function. The gradient is then defined as a vector of these partials. We'll represent \hat{y} as $f(x; \theta)$ to make the dependence on θ more obvious:

Updating using the Gradient

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

The final equation for updating θ based on the gradient is thus

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

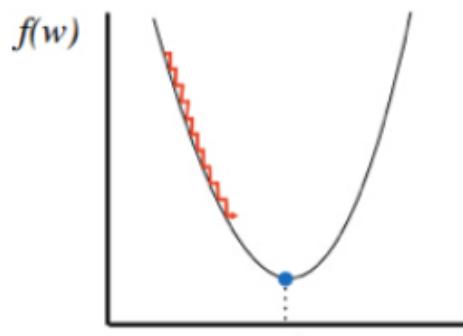
Gradient for Logistic Regression

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

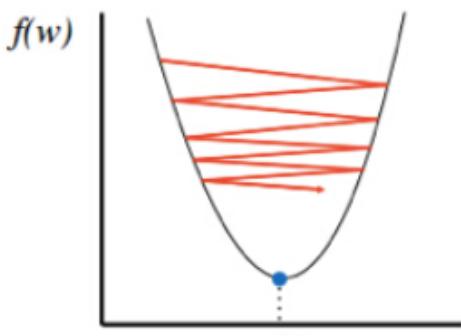
$$\frac{\partial L_{CE}(w, b)}{\partial w_j} = [\sigma(w \cdot x + b) - y]x_j$$

Learning Rate

- Updates: $\theta^{t+1} = \theta^t - \eta \frac{d}{d\theta} f(x; \theta)$
- Magnitude of movement along gradient
- Higher/faster learning rate = larger updates to parameters



Too small: converge
very slowly



Too big: overshoot and
even diverge

Stochastic Gradient Descent

```
function STOCHASTIC GRADIENT DESCENT( $L()$ ,  $f()$ ,  $x$ ,  $y$ ) returns  $\theta$ 
    # where: L is the loss function
    # f is a function parameterized by  $\theta$ 
    # x is the set of training inputs  $x^{(1)}, x^{(2)}, \dots, x^{(n)}$ 
    # y is the set of training outputs (labels)  $y^{(1)}, y^{(2)}, \dots, y^{(n)}$ 

     $\theta \leftarrow 0$ 
    repeat til done    # see caption
        For each training tuple  $(x^{(i)}, y^{(i)})$  (in random order)
            1. Optional (for reporting):      # How are we doing on this tuple?
                Compute  $\hat{y}^{(i)} = f(x^{(i)}; \theta)$  # What is our estimated output  $\hat{y}$ ?
                Compute the loss  $L(\hat{y}^{(i)}, y^{(i)})$  # How far off is  $\hat{y}^{(i)}$  from the true output  $y^{(i)}$ ?
            2.  $g \leftarrow \nabla_{\theta} L(f(x^{(i)}; \theta), y^{(i)})$  # How should we move  $\theta$  to maximize loss?
            3.  $\theta \leftarrow \theta - \eta g$            # Go the other way instead
    return  $\theta$ 
```

Figure 5.5 The stochastic gradient descent algorithm. Step 1 (computing the loss) is used to report how well we are doing on the current tuple. The algorithm can terminate when it converges (or when the gradient $< \epsilon$), or when progress halts (for example when the loss starts going up on a held-out set).

Let's walk though a single step of the gradient descent algorithm. We'll use a simplified version of the example in Fig. 5.2 as it sees a single observation x , whose correct value is $y = 1$ (this is a positive review), and with only two features:

$$x_1 = 3 \quad (\text{count of positive lexicon words})$$

$$x_2 = 2 \quad (\text{count of negative lexicon words})$$

Let's assume the initial weights and bias in θ^0 are all set to 0, and the initial learning rate η is 0.1:

$$w_1 = w_2 = b = 0$$

$$\eta = 0.1$$

The single update step requires that we compute the gradient, multiplied by the learning rate

$$\theta^{t+1} = \theta^t - \eta \nabla_{\theta} L(f(x^{(i)}; \theta), y^{(i)})$$

Example (Cont'd)

In our mini example there are three parameters, so the gradient vector has 3 dimensions, for w_1 , w_2 , and b . We can compute the first gradient as follows:

$$\nabla_{w,b} = \begin{bmatrix} \frac{\partial L_{CE}(w,b)}{\partial w_1} \\ \frac{\partial L_{CE}(w,b)}{\partial w_2} \\ \frac{\partial L_{CE}(w,b)}{\partial b} \end{bmatrix} = \begin{bmatrix} (\sigma(w \cdot x + b) - y)x_1 \\ (\sigma(w \cdot x + b) - y)x_2 \\ \sigma(w \cdot x + b) - y \end{bmatrix} = \begin{bmatrix} (\sigma(0) - 1)x_1 \\ (\sigma(0) - 1)x_2 \\ \sigma(0) - 1 \end{bmatrix} = \begin{bmatrix} -0.5x_1 \\ -0.5x_2 \\ -0.5 \end{bmatrix} = \begin{bmatrix} -1.5 \\ -1.0 \\ -0.5 \end{bmatrix}$$

Example (Cont'd)

Now that we have a gradient, we compute the new parameter vector θ^1 by moving θ^0 in the opposite direction from the gradient:

$$\theta^2 = \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix} - \eta \begin{bmatrix} -1.5 \\ -1.0 \\ -0.5 \end{bmatrix} = \begin{bmatrix} .15 \\ .1 \\ .05 \end{bmatrix}$$

So after one step of gradient descent, the weights have shifted to be: $w_1 = .15$, $w_2 = .1$, and $b = .05$.

Regularization

To avoid overfitting, a new **regularization** term $R(\theta)$ is added to the objective function in Eq. 5.12, resulting in the following objective for a batch of m examples (slightly rewritten from Eq. 5.12 to be maximizing log probability rather than minimizing loss, and removing the $\frac{1}{m}$ term which doesn't affect the argmax):

$$\hat{\theta} = \operatorname{argmax}_{\theta} \sum_{i=1}^m \log P(y^{(i)}|x^{(i)}) - \alpha R(\theta) \quad (5.23)$$

L2 Regularization

compute this regularization term $R(\theta)$. **L2 regularization** is a quadratic function of the weight values, named because it uses the (square of the) L2 norm of the weight values. The L2 norm, $\|\theta\|_2$, is the same as the **Euclidean distance** of the vector θ from the origin. If θ consists of n weights, then:

$$R(\theta) = \|\theta\|_2^2 = \sum_{j=1}^n \theta_j^2 \quad (5.24)$$

The L2 regularized objective function becomes:

$$\hat{\theta} = \operatorname{argmax}_{\theta} \left[\sum_{i=1}^m \log P(y^{(i)} | x^{(i)}) \right] - \alpha \sum_{j=1}^n \theta_j^2$$

L1 Regularization

L1 regularization is a linear function of the weight values, named after the L1 norm $\|W\|_1$, the sum of the absolute values of the weights, or **Manhattan distance** (the Manhattan distance is the distance you'd have to walk between two points in a city with a street grid like New York):

$$R(\theta) = \|\theta\|_1 = \sum_{i=1}^n |\theta_i| \quad (5.26)$$

The L1 regularized objective function becomes:

$$\hat{\theta} = \operatorname{argmax}_{\theta} \left[\sum_{i=1}^m \log P(y^{(i)} | x^{(i)}) \right] - \alpha \sum_{j=1}^n |\theta_j|$$

Multinomial LR

- More than two classes

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

$$\text{softmax}(z) = \left[\frac{e^{z_1}}{\sum_{i=1}^k e^{z_i}}, \frac{e^{z_2}}{\sum_{i=1}^k e^{z_i}}, \dots, \frac{e^{z_k}}{\sum_{i=1}^k e^{z_i}} \right]$$

Using Softmax

The denominator $\sum_{i=1}^k e^{z_i}$ is used to normalize all the values into probabilities.
Thus for example given a vector:

$$z = [0.6, 1.1, -1.5, 1.2, 3.2, -1.1]$$

the result softmax(z) is

$$[0.055, 0.090, 0.0067, 0.10, 0.74, 0.010]$$

Sigmoid vs. Softmax

Again like the sigmoid, the input to the softmax will be the dot product between a weight vector w and an input vector x (plus a bias). But now we'll need separate weight vectors (and bias) for each of the K classes.

$$p(y = c|x) = \frac{e^{w_c \cdot x + b_c}}{\sum_{j=1}^k e^{w_j \cdot x + b_j}} \quad (5.33)$$

Derivative of the Sigmoid

$$g_{\text{logistic}}(z) = \frac{1}{1+e^{-z}}$$

$$\begin{aligned} g'_{\text{logistic}}(z) &= \frac{\partial}{\partial z} \left(\frac{1}{1+e^{-z}} \right) \\ &= \frac{e^{-z}}{(1+e^{-z})^2} \text{ (by chain rule)} && \text{(double negation)} \\ &= \frac{1+e^{-z}-1}{(1+e^{-z})^2} \\ &= \frac{1+e^{-z}}{(1+e^{-z})^2} - \left(\frac{1}{1+e^{-z}} \right)^2 \\ &= \frac{1}{(1+e^{-z})} - \left(\frac{1}{1+e^{-z}} \right)^2 \\ &= g_{\text{logistic}}(z) - g_{\text{logistic}}(z)^2 \\ &= g_{\text{logistic}}(z)(1 - g_{\text{logistic}}(z)) \end{aligned}$$

Derivative of tanh

$$\begin{aligned}g_{\tanh}(z) &= \frac{\sinh(z)}{\cosh(z)} \\&= \frac{e^z - e^{-z}}{e^z + e^{-z}} \\g'_{\tanh}(z) &= \frac{\partial}{\partial z} \frac{\sinh(z)}{\cosh(z)} \\&= \frac{\frac{\partial}{\partial z} \sinh(z) \times \cosh(z) - \frac{\partial}{\partial z} \cosh(z) \times \sinh(z)}{\cosh^2(z)} \\&= \frac{\cosh^2(z) - \sinh^2(z)}{\cosh^2(z)} \\&= 1 - \frac{\sinh^2(z)}{\cosh^2(z)} \\&= 1 - \tanh^2(z)\end{aligned}$$

Logistic Regression - Summary

- ▶ Model

$$P(y = +|x) = \frac{\exp(\sum_{i=1}^n w_i x_i)}{1 + \exp(\sum_{i=1}^n w_i x_i)}$$

- ▶ Inference

$\operatorname{argmax}_y P(y|x)$ fundamentally same as Naive Bayes

$$P(y = 1|x) \geq 0.5 \Leftrightarrow w^\top x \geq 0$$

- ▶ Learning: gradient ascent on the (regularized) discriminative log-likelihood

Summary

- Logistic regression is a supervised machine learning classifier that extracts real-valued features from the input, multiplies each by a weight, sums them, and passes the sum through a **sigmoid** function to generate a probability. A threshold is used to make a decision.
- Logistic regression can be used with two classes (e.g., positive and negative sentiment) or with multiple classes (**multinomial logistic regression**, for example for n-ary text classification, part-of-speech labeling, etc.).
- Multinomial logistic regression uses the **softmax** function to compute probabilities.
- The weights (vector w and bias b) are learned from a labeled training set via a loss function, such as the **cross-entropy loss**, that must be minimized.
- Minimizing this loss function is a **convex optimization** problem, and iterative algorithms like **gradient descent** are used to find the optimal weights.
- **Regularization** is used to avoid overfitting.
- Logistic regression is also one of the most useful analytic tools, because of its ability to transparently study the importance of individual features.

Summary of Learning Methods

It is natural to ask which learning algorithm is best, but the answer depends on what characteristics are important to the problem you are trying to solve.

Naïve Bayes *Pros:* easy to implement; estimation is fast, requiring only a single pass over the data; assigns probabilities to predicted labels; controls overfitting with smoothing parameter. *Cons:* often has poor accuracy, especially with correlated features.

Perceptron *Pros:* easy to implement; online; error-driven learning means that accuracy is typically high, especially after averaging. *Cons:* not probabilistic; hard to know when to stop learning; lack of margin can lead to overfitting.

Support vector machine *Pros:* optimizes an error-based metric, usually resulting in high accuracy; overfitting is controlled by a regularization parameter. *Cons:* not probabilistic.

Logistic regression *Pros:* error-driven and probabilistic; overfitting is controlled by a regularization parameter. *Cons:* batch learning requires black-box optimization; logistic loss can “overtrain” on correctly labeled examples.

Non-linear Learning Methods

- **Kernel methods** are generalizations of the **nearest-neighbor** classification rule, which classifies each instance by the label of the most similar example in the training set. The application of the **kernel support vector machine** to information extraction is described in chapter 17.
- **Decision trees** classify instances by checking a set of conditions. Scaling decision trees to bag-of-words inputs is difficult, but decision trees have been successful in problems such as coreference resolution (chapter 15), where more compact feature sets can be constructed (Soon et al., 2001).
- **Boosting** and related **ensemble methods** work by combining the predictions of several “weak” classifiers, each of which may consider only a small subset of features. Boosting has been successfully applied to text classification (Schapire and Singer, 2000) and syntactic analysis (Abney et al., 1999), and remains one of the most successful methods on machine learning competition sites such as Kaggle (Chen and Guestrin, 2016).