

# Introduction to Machine Learning

## Linear Classifiers

Lisbon Machine Learning School, 2015

Shay Cohen

School of Informatics, University of Edinburgh

E-mail: `scohen@inf.ed.ac.uk`

Slides heavily based on Ryan McDonald's slides from 2014

# Linear Classifiers

- ▶ Go onto ACL Anthology
- ▶ Search for: “Naive Bayes”, “Maximum Entropy”, “Logistic Regression”, “SVM”, “Perceptron”
- ▶ Do the same on Google Scholar
  - ▶ “Maximum Entropy” & “NLP” 11,000 hits, 240 before 2000
  - ▶ “SVM” & “NLP” 15,000 hits, 556 before 2000
  - ▶ “Perceptron” & “NLP”, 4,000 hits, 147 before 2000
- ▶ All are examples of linear classifiers
- ▶ All have become tools in any NLP/CL researchers tool-box in past 15 years
  - ▶ One the most important tools

# Experiment

- ▶ Document 1 – label: 0; words: ★ ◇ ○
- ▶ Document 2 – label: 0; words: ★ ♥ △
- ▶ Document 3 – label: 1; words: ★ △ ♠
- ▶ Document 4 – label: 1; words: ◇ △ ○
  
- ▶ New document – words: ★ ◇ ○; label ?
- ▶ New document – words: ★ ◇ ♥; label ?
- ▶ New document – words: ★ ◇ ♠; label ?
- ▶ New document – words: ★ △ ○; label ?

Why and how can we do this?

# Experiment

- ▶ Document 1 – label: 0; words: ★ ◇ ○
- ▶ Document 2 – label: 0; words: ★ ♥ △
- ▶ Document 3 – label: 1; words: ★ △ ♠
- ▶ Document 4 – label: 1; words: ◇ △ ○
- ▶ New document – words: ★ △ ○; label ?

**Label 0**

**Label 1**

$$P(0|\star) = \frac{\text{count}(\star \text{ and } 0)}{\text{count}(\star)} = \frac{2}{3} = \mathbf{0.67} \text{ vs. } P(1|\star) = \frac{\text{count}(\star \text{ and } 1)}{\text{count}(\star)} = \frac{1}{3} = \mathbf{0.33}$$

$$P(0|\triangle) = \frac{\text{count}(\triangle \text{ and } 0)}{\text{count}(\triangle)} = \frac{1}{3} = \mathbf{0.33} \text{ vs. } P(1|\triangle) = \frac{\text{count}(\triangle \text{ and } 1)}{\text{count}(\triangle)} = \frac{2}{3} = \mathbf{0.67}$$

$$P(0|\circ) = \frac{\text{count}(\circ \text{ and } 0)}{\text{count}(\circ)} = \frac{1}{2} = 0.5 \text{ vs. } P(1|\circ) = \frac{\text{count}(\circ \text{ and } 1)}{\text{count}(\circ)} = \frac{1}{2} = 0.5$$

# Machine Learning

- ▶ Machine learning is **well-motivated counting**
- ▶ Typically, machine learning models
  1. Define a model/distribution of interest
  2. Make some assumptions if needed
  3. Count!!
- ▶ Model:  $P(\text{label}|\text{doc}) = P(\text{label}|\text{word}_1, \dots, \text{word}_n)$ 
  - ▶ Prediction for new doc =  $\arg \max_{\text{label}} P(\text{label}|\text{doc})$
- ▶ Assumption:  $P(\text{label}|\text{word}_1, \dots, \text{word}_n) = \frac{1}{n} \sum_i P(\text{label}|\text{word}_i)$
- ▶ Count (as in example)

# Lecture Outline

- ▶ Preliminaries
  - ▶ Data: input/output, assumptions
  - ▶ Feature representations
  - ▶ Linear classifiers and decision boundaries
- ▶ Classifiers
  - ▶ Naive Bayes
  - ▶ Generative versus discriminative
  - ▶ Logistic-regression
  - ▶ Perceptron
  - ▶ Large-Margin Classifiers (SVMs)
- ▶ Regularization
- ▶ Online learning
- ▶ Non-linear classifiers

# Inputs and Outputs

- ▶ Input:  $x \in \mathcal{X}$ 
  - ▶ e.g., document or sentence with some words  $x = w_1 \dots w_n$ , or a series of previous actions
- ▶ Output:  $y \in \mathcal{Y}$ 
  - ▶ e.g., parse tree, document class, part-of-speech tags, word-sense
- ▶ Input/Output pair:  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ 
  - ▶ e.g., a document  $x$  and its label  $y$
  - ▶ Sometimes  $x$  is explicit in  $y$ , e.g., a parse tree  $y$  will contain the sentence  $x$

# General Goal

When given a new input  $x$  predict the correct output  $y$

But we need to formulate this computationally!



# Feature Representations

- ▶ We assume a mapping from input  $x$  to a high dimensional **feature vector**
  - ▶  $\phi(x) : \mathcal{X} \rightarrow \mathbb{R}^m$
- ▶ For many cases, more convenient to have mapping from input-output pairs  $(x, y)$ 
  - ▶  $\phi(x, y) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^m$
- ▶ Under certain assumptions, these are equivalent
- ▶ Most papers in NLP use  $\phi(x, y)$
- ▶ (Was?) not so common in NLP:  $\phi \in \mathbb{R}^m$  (but see word embeddings)
- ▶ More common:  $\phi_i \in \{1, \dots, F_i\}$ ,  $F_i \in \mathbb{N}^+$  (categorical)
- ▶ Very common:  $\phi \in \{0, 1\}^m$  (binary)
- ▶ For any vector  $\mathbf{v} \in \mathbb{R}^m$ , let  $\mathbf{v}_j$  be the  $j^{\text{th}}$  value

# Examples

- ▶  $x$  is a document and  $y$  is a label

$$\phi_j(x, y) = \begin{cases} 1 & \text{if } x \text{ contains the word "interest"} \\ & \text{and } y = \text{"financial"} \\ 0 & \text{otherwise} \end{cases}$$

We expect this feature to have a positive weight, “interest” is a positive indicator for the label “financial”

# Examples

- ▶  $x$  is a document and  $y$  is a label

$$\phi_j(x, y) = \begin{cases} 1 & \text{if } x \text{ contains the word "president"} \\ & \text{and } y = \text{"sports"} \\ 0 & \text{otherwise} \end{cases}$$

We expect this feature to have a negative weight?

# Examples

$\phi_j(\mathbf{x}, \mathbf{y}) = \%$  of words in  $\mathbf{x}$  containing punctuation and  $\mathbf{y} = \text{"scientific"}$

Punctuation symbols - positive indicator or negative indicator for scientific articles?

# Examples

- ▶  $x$  is a word and  $y$  is a part-of-speech tag

$$\phi_j(x, y) = \begin{cases} 1 & \text{if } x = \text{"bank"} \text{ and } y = \text{Verb} \\ 0 & \text{otherwise} \end{cases}$$

What weight would it get?

## Example 2

- $x$  is a name,  $y$  is a label classifying the name

$$\phi_0(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "George"} \\ & \text{and } y = \text{"Person"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_4(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "George"} \\ & \text{and } y = \text{"Object"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_1(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "Washington"} \\ & \text{and } y = \text{"Person"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_5(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "Washington"} \\ & \text{and } y = \text{"Object"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_2(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "Bridge"} \\ & \text{and } y = \text{"Person"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_6(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "Bridge"} \\ & \text{and } y = \text{"Object"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_3(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "General"} \\ & \text{and } y = \text{"Person"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_7(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "General"} \\ & \text{and } y = \text{"Object"} \\ 0 & \text{otherwise} \end{cases}$$

- $x=\text{General George Washington}, y=\text{Person} \rightarrow \phi(x, y) = [1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$
- $x=\text{George Washington Bridge}, y=\text{Object} \rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 0]$
- $x=\text{George Washington George}, y=\text{Object} \rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0]$

# Block Feature Vectors

- ▶  $x$ =General George Washington,  $y$ =Person  $\rightarrow \phi(x, y) = [1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$
  - ▶  $x$ =General George Washington,  $y$ =Object  $\rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 1]$
  - ▶  $x$ =George Washington Bridge,  $y$ =Object  $\rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 0]$
  - ▶  $x$ =George Washington George,  $y$ =Object  $\rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0]$
- 
- ▶ Each equal size block of the feature vector corresponds to one label
  - ▶ Non-zero values allowed only in one block

## Feature Representations - $\phi(x)$

- ▶ Instead of  $\phi(x, y) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^m$  over input/outputs  $(x, y)$
- ▶ Let  $\phi(x) : \mathcal{X} \rightarrow \mathbb{R}^{m'}$  (e.g.,  $m' = m/|\mathcal{Y}|$ )
  - ▶ i.e., feature representation only over inputs  $x$
- ▶ Equivalent when  $\phi(x, y)$  includes  $y$  as a non-decomposable object
- ▶ Disadvantages to  $\phi(x)$  formulation: no complex features over properties of labels
- ▶ Advantages: can make math cleaner, especially with binary classification



# Feature Representations - $\phi(x)$ vs. $\phi(x, y)$

- ▶  $\phi(x, y)$

- ▶  $x$ =General George Washington,  $y$ =Person  $\rightarrow \phi(x, y) = [1\ 1\ 0\ 1\ 0\ 0\ 0\ 0]$
- ▶  $x$ =General George Washington,  $y$ =Object  $\rightarrow \phi(x, y) = [0\ 0\ 0\ 0\ 1\ 1\ 0\ 1]$

- ▶  $\phi(x)$

- ▶  $x$ =General George Washington  $\rightarrow \phi(x) = [1\ 1\ 0\ 1]$

- ▶ Different ways of representing same thing
- ▶ In this case, can deterministically map from  $\phi(x)$  to  $\phi(x, y)$  given  $y$

# Linear Classifiers

- ▶ **Linear classifier:** **score** (or probability) of a particular classification is based on a linear combination of features and their **weights**
- ▶ Let  $\omega \in \mathbb{R}^m$  be a high dimensional weight vector
- ▶ Assume that  $\omega$  is known
  - ▶ **Multiclass Classification:**  $\mathcal{Y} = \{0, 1, \dots, N\}$

$$\begin{aligned} y &= \arg \max_y \omega \cdot \phi(x, y) \\ &= \arg \max_y \sum_{j=0}^m \omega_j \times \phi_j(x, y) \end{aligned}$$

- ▶ **Binary Classification** just a special case of multiclass

# Linear Classifiers – $\phi(x)$

- ▶ Define  $|\mathcal{Y}|$  parameter vectors  $\omega_y \in \mathbb{R}^{m'}$ 
  - ▶ I.e., one parameter vector per output class  $y$

## ▶ Classification

$$y = \arg \max_y \omega_y \cdot \phi(x)$$

- ▶  $\phi(x, y)$ 
  - ▶  $x=\text{General George Washington}, y=\text{Person} \rightarrow \phi(x, y) = [1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$
  - ▶  $x=\text{General George Washington}, y=\text{Object} \rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 1]$
  - ▶ Single  $\omega \in \mathbb{R}^8$
- ▶  $\phi(x)$ 
  - ▶  $x=\text{General George Washington} \rightarrow \phi(x) = [1 \ 1 \ 0 \ 1]$
  - ▶ Two parameter vectors  $\omega_0 \in \mathbb{R}^4, \omega_1 \in \mathbb{R}^4$

# Linear Classifiers - Bias Terms

- ▶ Often linear classifiers presented as

$$\text{argmax}_y \sum_{j=0}^m \omega_j \times \phi_j(x, y) + b_y$$

- ▶ Where  $b$  is a bias or offset term
- ▶ Sometimes this is folded into  $\phi$

$x$ =General George Washington,  $y$ =Person  $\rightarrow \phi(x, y) = [1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0]$

$x$ =General George Washington,  $y$ =Object  $\rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 1]$

$$\phi_4(x, y) = \begin{cases} 1 & y = \text{"Person"} \\ 0 & \text{otherwise} \end{cases}$$

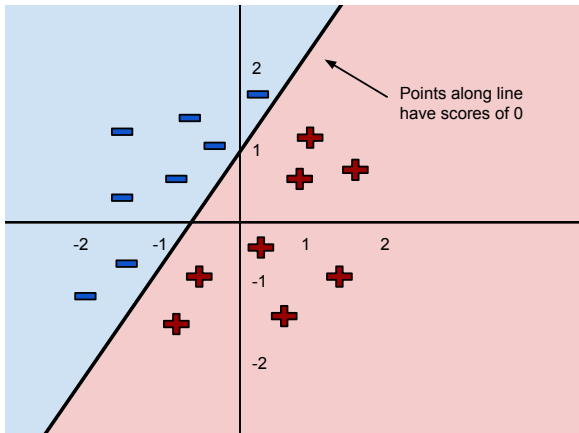
$$\phi_9(x, y) = \begin{cases} 1 & y = \text{"Object"} \\ 0 & \text{otherwise} \end{cases}$$

- ▶  $\omega_4$  and  $\omega_9$  are now the bias terms for the labels

# Binary Linear Classifier

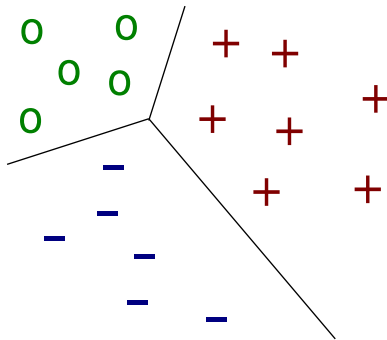
Let's say  $\omega = (1, -1)$  and  $b_y = 1, \forall y$

Then  $\omega$  is a line (generally a hyperplane) that divides all points:



# Multiclass Linear Classifier

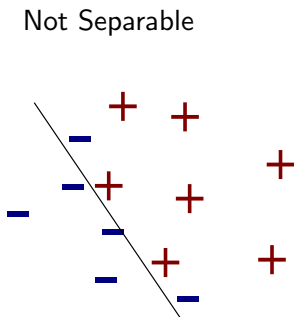
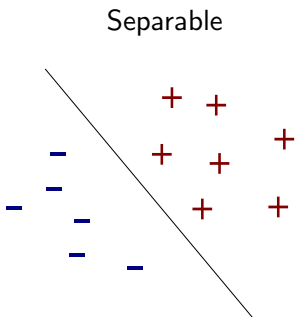
Defines regions of space. Visualization difficult.



- i.e.,  $+$  are all points  $(x, y)$  where  $+$  =  $\arg \max_y \omega \cdot \phi(x, y)$

# Separability

- ▶ A set of points is separable, if there exists a  $\omega$  such that classification is perfect



- ▶ This can also be defined mathematically (and we will do that shortly)

# Machine Learning – finding $\omega$

We now have a way to make decisions... If we have a  $\omega$ . But where do we get this  $\omega$ ?

- ▶ Supervised Learning
- ▶ Input: training examples  $\mathcal{T} = \{(x_t, y_t)\}_{t=1}^{|\mathcal{T}|}$
- ▶ Input: feature representation  $\phi$
- ▶ Output:  $\omega$  that maximizes some **important function** on the training set
  - ▶  $\omega = \arg \max \mathcal{L}(\mathcal{T}; \omega)$
- ▶ Equivalently minimize:  $\omega = \arg \min -\mathcal{L}(\mathcal{T}; \omega)$



# Objective Functions

- ▶  $\mathcal{L}(\cdot)$  is called the **objective function**
- ▶ Usually we can decompose  $\mathcal{L}$  by training pairs  $(x, y)$ 
  - ▶  $\mathcal{L}(\mathcal{T}; \omega) \propto \sum_{(x, y) \in \mathcal{T}} \text{loss}((x, y); \omega)$
  - ▶  $\text{loss}$  is a function that measures some value correlated with errors of parameters  $\omega$  on instance  $(x, y)$
- ▶ Defining  $\mathcal{L}(\cdot)$  and  $\text{loss}$  is core of linear classifiers in machine learning
- ▶ Example:  $y \in \{1, -1\}$ ,  $f(x|w)$  is the prediction we make for  $x$  using  $w$
- ▶ Loss is:

# Supervised Learning – Assumptions

- ▶ Assumption:  $(\mathbf{x}_t, \mathbf{y}_t)$  are sampled i.i.d.
  - ▶ i.i.d. = independent and identically distributed
  - ▶ independent = each sample independent of the other
  - ▶ identically = each sample from same probability distribution
- ▶ Sometimes assumption: The training data is separable
  - ▶ Needed to prove convergence for Perceptron
  - ▶ Not needed in practice

# Naive Bayes

# Probabilistic Models

- ▶ Let's put aside linear classifiers for a moment
- ▶ Here is another approach to decision making
  - ▶ Probabilistically model  $P(\mathbf{y}|\mathbf{x})$
  - ▶ If we can define this distribution, then classification becomes
    - ▶  $\arg \max_{\mathbf{y}} P(\mathbf{y}|\mathbf{x})$

# Bayes Rule

- ▶ One way to model  $P(y|x)$  is through **Bayes Rule**:

$$P(y|x) = \frac{P(y)P(x|y)}{P(x)}$$

$$\arg \max_y P(y|x) \propto \arg \max_y P(y)P(x|y)$$

- ▶ Since  $x$  is fixed
- ▶  $P(y)P(x|y) = P(x, y)$ : a joint probability
- ▶ Modeling the joint input-output distribution is at the core of **generative models**
  - ▶ Because we model a distribution that can randomly generate outputs and inputs, not just outputs
  - ▶ More on this later

# Naive Bayes (NB)

- ▶ We need to decide on the structure of  $P(x, y)$
- ▶  $P(x|y) = P(\phi(x)|y) = P(\phi_1(x), \dots, \phi_m(x)|y)$

Naive Bayes Assumption

*(conditional independence)*

$$\boxed{P(\phi_1(x), \dots, \phi_m(x)|y) = \prod_i P(\phi_i(x)|y)}$$

$$P(x, y) = P(y)P(\phi_1(x), \dots, \phi_m(x)|y) = P(y) \prod_{i=1}^m P(\phi_i(x)|y)$$

# Naive Bayes – Learning

- ▶ Input:  $\mathcal{T} = \{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^{|\mathcal{T}|}$
- ▶ Let  $\phi_i(x) \in \{1, \dots, F_i\}$  – categorical; common in NLP
- ▶ Parameters  $\mathcal{P} = \{P(\mathbf{y}), P(\phi_i(x)|\mathbf{y})\}$ 
  - ▶ Both  $P(\mathbf{y})$  and  $P(\phi_i(x)|\mathbf{y})$  are multinomials

It is important that we assume that features are independent GIVEN  $\mathbf{y}$ ... they are typically not independent a-priori e.g. We assume that "basketball", "football" | SPORT are independent given SPORT, but obviously "basketball", "football" is dependent (via sport)

# Maximum Likelihood Estimation

- ▶ What's left? Defining an objective  $\mathcal{L}(\mathcal{T})$
- ▶  $P$  plays the role of  $w$
- ▶ What objective to use?
- ▶ **Objective: Maximum Likelihood Estimation (MLE)**

$$\mathcal{L}(\mathcal{T}) = \prod_{t=1}^{|\mathcal{T}|} P(\mathbf{x}_t, \mathbf{y}_t) = \prod_{t=1}^{|\mathcal{T}|} \left( P(\mathbf{y}_t) \prod_{i=1}^m P(\phi_i(\mathbf{x}_t) | \mathbf{y}_t) \right)$$

$$\mathcal{P} = \arg \max_{\mathcal{P}} \prod_{t=1}^{|\mathcal{T}|} \left( P(\mathbf{y}_t) \prod_{i=1}^m P(\phi_i(\mathbf{x}_t) | \mathbf{y}_t) \right)$$



# Naive Bayes – Learning

MLE has **closed form solution!!** (more later) – count and normalize

$$\mathcal{P} = \arg \max_{\mathcal{P}} \prod_{t=1}^{|\mathcal{T}|} \left( P(\mathbf{y}_t) \prod_{i=1}^m P(\phi_i(\mathbf{x}_t) | \mathbf{y}_t) \right)$$

$$P(\mathbf{y}) = \frac{\sum_{t=1}^{|\mathcal{T}|} [[\mathbf{y}_t = \mathbf{y}]]}{|\mathcal{T}|}$$

$$P(\phi_i(\mathbf{x}) | \mathbf{y}) = \frac{\sum_{t=1}^{|\mathcal{T}|} [[\phi_i(\mathbf{x}_t) = \phi_i(\mathbf{x}) \text{ and } \mathbf{y}_t = \mathbf{y}]]}{\sum_{t=1}^{|\mathcal{T}|} [[\mathbf{y}_t = \mathbf{y}]]}$$

$[[X]]$  is the identity function for property  $X$

Thus, these are just normalized counts over events in  $\mathcal{T}$

**Intuitively makes sense!**

# Naive Bayes Example

- ▶  $\phi_i(x) \in 0, 1, \forall i$
- ▶ doc 1:  $y_1 = 0, \phi_0(x_1) = 1, \phi_1(x_1) = 1$
- ▶ doc 2:  $y_2 = 0, \phi_0(x_2) = 0, \phi_1(x_2) = 1$
- ▶ doc 3:  $y_3 = 1, \phi_0(x_3) = 1, \phi_1(x_3) = 0$
- ▶ Two label parameters  $P(y = 0), P(y = 1)$
- ▶ Eight feature parameters
  - ▶ 2 (labels) \* 2 (features) \* 2 (feature values)
  - ▶ E.g.,  $y = 0$  and  $\phi_0(x) = 1$ :  $P(\phi_0(x) = 1|y = 0)$
- ▶ We really have one label parameter and  $2 * 2 * (2 - 1)$  feature parameters
- ▶  $P(y = 0) = 2/3, P(y = 1) = 1/3$
- ▶  $P(\phi_0(x) = 1|y = 0) = 1/2, P(\phi_1(x) = 0|y = 1) = 1/1$

# Naive Bayes Document Classification

- ▶ doc 1:  $y_1 = \text{sports}$ , “hockey is fast”
- ▶ doc 2:  $y_2 = \text{politics}$ , “politicians talk fast”
- ▶ doc 3:  $y_3 = \text{politics}$ , “washington is sleazy”
  
- ▶  $\phi_0(x) = 1$  iff doc has word ‘hockey’, 0 o.w.
- ▶  $\phi_1(x) = 1$  iff doc has word ‘is’, 0 o.w.
- ▶  $\phi_2(x) = 1$  iff doc has word ‘fast’, 0 o.w.
- ▶  $\phi_3(x) = 1$  iff doc has word ‘politicians’, 0 o.w.
- ▶  $\phi_4(x) = 1$  iff doc has word ‘talk’, 0 o.w.
- ▶  $\phi_5(x) = 1$  iff doc has word ‘washington’, 0 o.w.
- ▶  $\phi_6(x) = 1$  iff doc has word ‘sleazy’, 0 o.w.

Your turn? What is  $P(\text{sports})$ ? What is  $P(\phi_0(0) = 1 | \text{politics})$ ?

## Deriving MLE

$$\mathcal{P} = \arg \max_{\mathcal{P}} \prod_{t=1}^{|\mathcal{T}|} \left( P(\mathbf{y}_t) \prod_{i=1}^m P(\phi_i(\mathbf{x}_t) | \mathbf{y}_t) \right)$$

# Deriving MLE (for handout)

$$\begin{aligned}
 \mathcal{P} &= \arg \max_{\mathcal{P}} \prod_{t=1}^{|\mathcal{T}|} \left( P(\mathbf{y}_t) \prod_{i=1}^m P(\phi_i(\mathbf{x}_t) | \mathbf{y}_t) \right) \\
 &= \arg \max_{\mathcal{P}} \sum_{t=1}^{|\mathcal{T}|} \left( \log P(\mathbf{y}_t) + \sum_{i=1}^m \log P(\phi_i(\mathbf{x}_t) | \mathbf{y}_t) \right) \\
 &= \arg \max_{P(\mathbf{y})} \sum_{t=1}^{|\mathcal{T}|} \log P(\mathbf{y}_t) + \arg \max_{P(\phi_i(\mathbf{x}) | \mathbf{y})} \sum_{t=1}^{|\mathcal{T}|} \sum_{i=1}^m \log P(\phi_i(\mathbf{x}_t) | \mathbf{y}_t)
 \end{aligned}$$

such that  $\sum_{\mathbf{y}} P(\mathbf{y}) = 1$ ,  $\sum_{j=1}^{F_i} P(\phi_i(\mathbf{x}) = j | \mathbf{y}) = 1$ ,  $P(\cdot) \geq 0$

# Deriving MLE

$$\mathcal{P} = \arg \max_{P(\mathbf{y})} \sum_{t=1}^{|\mathcal{T}|} \log P(\mathbf{y}_t) + \arg \max_{P(\phi_i(\mathbf{x})|\mathbf{y})} \sum_{t=1}^{|\mathcal{T}|} \sum_{i=1}^m \log P(\phi_i(\mathbf{x}_t)|\mathbf{y}_t)$$

Both optimizations are of the form

$$\arg \max_P \sum_v \text{count}(v) \log P(v), \text{ s.t., } \sum_v P(v) = 1, P(v) \geq 0$$

For example:

$$\arg \max_{P(\mathbf{y})} \sum_{t=1}^{|\mathcal{T}|} \log P(\mathbf{y}_t) = \arg \max_{P(\mathbf{y})} \sum_{\mathbf{y}} \text{count}(\mathbf{y}, \mathcal{T}) \log P(\mathbf{y})$$

$$\text{such that } \sum_{\mathbf{y}} P(\mathbf{y}) = 1, P(\mathbf{y}) \geq 0$$

# Deriving MLE

$$\begin{aligned} \arg \max_P \sum_v \text{count}(v) \log P(v) \\ \text{s.t.}, \sum_v P(v) = 1, P(v) \geq 0 \end{aligned}$$

Introduce Lagrangian multiplier  $\lambda$ , optimization becomes

~~$$\arg \max_{P, \lambda} \sum_v \text{count}(v) \log P(v) - \lambda (\sum_v P(v) - 1)$$~~

To do derivative

Derivative:

Set to zero:

Final solution:

## Deriving MLE (for handout)

$$\begin{aligned} \arg \max_P \sum_v \text{count}(v) \log P(v) \\ \text{s.t., } \sum_v P(v) = 1, P(v) \geq 0 \end{aligned}$$

Introduce Lagrangian multiplier  $\lambda$ , optimization becomes

$$\arg \max_{P, \lambda} \sum_v \text{count}(v) \log P(v) - \lambda (\sum_v P(v) - 1)$$

Derivative w.r.t  $P(v)$  is  $\frac{\text{count}(v)}{P(v)} - \lambda$

Setting this to zero  $P(v) = \frac{\text{count}(v)}{\lambda}$

Combine with  $\sum_v P(v) = 1$ .  $P(v) \geq 0$ , then  $P(v) = \frac{\text{count}(v)}{\sum_{v'} \text{count}(v')}$



## Put it together

$$\begin{aligned}\mathcal{P} &= \arg \max_{\mathcal{P}} \prod_{t=1}^{|\mathcal{T}|} \left( P(\mathbf{y}_t) \prod_{i=1}^m P(\phi_i(\mathbf{x}_t) | \mathbf{y}_t) \right) \\ &= \arg \max_{P(\mathbf{y})} \sum_{t=1}^{|\mathcal{T}|} \log P(\mathbf{y}_t) + \arg \max_{P(\phi_i(\mathbf{x}) | \mathbf{y})} \sum_{t=1}^{|\mathcal{T}|} \sum_{i=1}^m \log P(\phi_i(\mathbf{x}_t) | \mathbf{y}_t)\end{aligned}$$

$$P(\mathbf{y}) = \frac{\sum_{t=1}^{|\mathcal{T}|} [[\mathbf{y}_t = \mathbf{y}]]}{|\mathcal{T}|}$$

$$P(\phi_i(\mathbf{x}) | \mathbf{y}) = \frac{\sum_{t=1}^{|\mathcal{T}|} [[\phi_i(\mathbf{x}_t) = \phi_i(\mathbf{x}) \text{ and } \mathbf{y}_t = \mathbf{y}]]}{\sum_{t=1}^{|\mathcal{T}|} [[\mathbf{y}_t = \mathbf{y}]]}$$

# NB is a linear classifier

- ▶ Let  $\omega_y = \log P(y)$ ,  $\forall y \in \mathcal{Y}$
- ▶ Let  $\omega_{\phi_i(x), y} = \log P(\phi_i(x)|y)$ ,  $\forall y \in \mathcal{Y}, \phi_i(x) \in \{1, \dots, F_i\}$
- ▶ Let  $\omega$  be set of all  $\omega_*$  and  $\omega_{*,*}$

$$\arg \max_y P(y|\phi(x)) \propto \arg \max_y P(\phi(x), y) = \arg \max_y P(y) \prod_{i=1}^m P(\phi_i(x)|y) =$$

where  $\psi_* \in \{0, 1\}$ ,  $\psi_{i,j}(x) = [[\phi_i(x) = j]]$ ,  $\psi_{y'}(y) = [[y = y']]$

# NB is a linear classifier (for handout)

- ▶ Let  $\omega_{\mathbf{y}} = \log P(\mathbf{y}), \forall \mathbf{y} \in \mathcal{Y}$
- ▶ Let  $\omega_{\phi_i(\mathbf{x}), \mathbf{y}} = \log P(\phi_i(\mathbf{x})|\mathbf{y}), \forall \mathbf{y} \in \mathcal{Y}, \phi_i(\mathbf{x}) \in \{1, \dots, F_i\}$
- ▶ Let  $\omega$  be set of all  $\omega_*$  and  $\omega_{*,*}$

$$\arg \max_{\mathbf{y}} P(\mathbf{y}|\phi(\mathbf{x})) \propto \arg \max_{\mathbf{y}} P(\phi(\mathbf{x}), \mathbf{y}) = \arg \max_{\mathbf{y}} P(\mathbf{y}) \prod_{i=1}^m P(\phi_i(\mathbf{x})|\mathbf{y})$$

it is linear classifier:  
show it by  
concatenating the  
"sums over y' and j"

$$\arg \max_{\mathbf{y}} \log P(\mathbf{y}) + \sum_{i=1}^m \log P(\phi_i(\mathbf{x})|\mathbf{y})$$

$$\arg \max_{\mathbf{y}} \omega_{\mathbf{y}} + \sum_{i=1}^m \omega_{\phi_i(\mathbf{x}), \mathbf{y}}$$

$$= \arg \max_{\mathbf{y}} \sum_{\mathbf{y}'} \omega_{\mathbf{y}} \psi_{\mathbf{y}'}(\mathbf{y}) + \sum_{i=1}^m \sum_{j=1}^{F_i} \omega_{\phi_i(\mathbf{x}), \mathbf{y}} \psi_{i,j}(\mathbf{x})$$

where  $\psi_* \in \{0, 1\}$ ,  $\psi_{i,j}(\mathbf{x}) = [[\phi_i(\mathbf{x}) = j]]$ ,  $\psi_{\mathbf{y}'}(\mathbf{y}) = [[\mathbf{y} = \mathbf{y}']]$

# Smoothing

- ▶ doc 1:  $y_1 = \text{sports, "hockey is fast"}$
- ▶ doc 2:  $y_2 = \text{politics, "politicians talk fast"}$
- ▶ doc 3:  $y_3 = \text{politics, "washington is sleazy"}$
  
- ▶ New doc: "washington hockey is fast"
- ▶ Both 'sports' and 'politics' have probabilities of 0
  
- ▶ Smoothing aims to assign a small amount of probability to unseen events
- ▶ E.g., Additive/Laplacian smoothing

$$P(v) = \frac{\text{count}(v)}{\sum_{v'} \text{count}(v')} \implies P(v) = \frac{\text{count}(v) + \alpha}{\sum_{v'} (\text{count}(v') + \alpha)}$$

# Discriminative versus Generative

- ▶ Generative models attempt to model inputs and outputs
  - ▶ e.g., NB = MLE of joint distribution  $P(\mathbf{x}, \mathbf{y})$
  - ▶ Statistical model must explain generation of input
- ▶ Occam's Razor: why model input?
- ▶ Discriminative models
  - ▶ Use  $\mathcal{L}$  that directly optimizes  $P(\mathbf{y}|\mathbf{x})$  (or something related)
  - ▶ Logistic Regression – MLE of  $P(\mathbf{y}|\mathbf{x})$
  - ▶ Perceptron and SVMs – minimize classification error
- ▶ Generative and discriminative models use  $P(\mathbf{y}|\mathbf{x})$  for prediction
- ▶ Differ only on what distribution they use to set  $\omega$

# Logistic Regression

# Logistic Regression

exponential is based  
on max-entropy  
models: TODO How?

Define a conditional probability:

$$P(\mathbf{y}|\mathbf{x}) = \frac{e^{\omega \cdot \phi(\mathbf{x}, \mathbf{y})}}{Z_{\mathbf{x}}}, \quad \text{where } Z_{\mathbf{x}} = \sum_{\mathbf{y}' \in \mathcal{Y}} e^{\omega \cdot \phi(\mathbf{x}, \mathbf{y}')}$$


Note: still a linear classifier

$$\begin{aligned} \arg \max_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}) &= \arg \max_{\mathbf{y}} \frac{e^{\omega \cdot \phi(\mathbf{x}, \mathbf{y})}}{Z_{\mathbf{x}}} \\ &= \arg \max_{\mathbf{y}} e^{\omega \cdot \phi(\mathbf{x}, \mathbf{y})} \\ &= \arg \max_{\mathbf{y}} \omega \cdot \phi(\mathbf{x}, \mathbf{y}) \end{aligned}$$

# Logistic Regression

$$P(y|x) = \frac{e^{\omega \cdot \phi(x,y)}}{Z_x}$$

conditional  
loglikelihood



- ▶ Q: How do we learn weights  $\omega$
- ▶ A: Set weights to maximize log-likelihood of training data:

$$\begin{aligned}\omega &= \arg \max_{\omega} \mathcal{L}(\mathcal{T}; \omega) \\ &= \arg \max_{\omega} \prod_{t=1}^{|\mathcal{T}|} P(y_t|x_t) = \arg \max_{\omega} \sum_{t=1}^{|\mathcal{T}|} \log P(y_t|x_t)\end{aligned}$$

- ▶ In a nutshell we set the weights  $\omega$  so that we assign as much probability to the correct label  $y$  for each  $x$  in the training set



# Logistic Regression

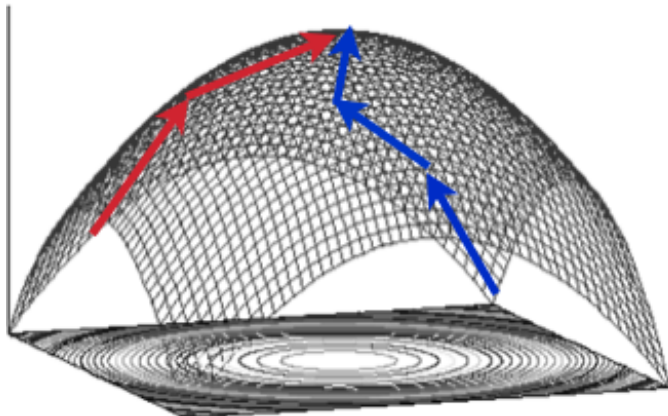
$$P(y|x) = \frac{e^{\omega \cdot \phi(x,y)}}{Z_x}, \quad \text{where } Z_x = \sum_{y' \in \mathcal{Y}} e^{\omega \cdot \phi(x,y')}$$

$$\omega = \arg \max_{\omega} \sum_{t=1}^{|\mathcal{T}|} \log P(y_t|x_t) (*)$$

- ▶ The objective function (\*) is concave (take the 2nd derivative)
- ▶ Therefore there is a global maximum
- ▶ No closed form solution, but lots of numerical techniques
  - ▶ Gradient methods (gradient ascent, conjugate gradient, iterative scaling)
  - ▶ Newton methods (limited-memory quasi-newton)

why?

# Gradient Ascent



# Gradient Ascent

- ▶ Let  $\mathcal{L}(\mathcal{T}; \omega) = \sum_{t=1}^{|\mathcal{T}|} \log(e^{\omega \cdot \phi(x_t, y_t)} / Z_x)$
- ▶ Want to find  $\arg \max_{\omega} \mathcal{L}(\mathcal{T}; \omega)$ 
  - ▶ Set  $\omega^0 = 0^m$
  - ▶ Iterate until convergence

$$\omega^i = \omega^{i-1} + \alpha \nabla \mathcal{L}(\mathcal{T}; \omega^{i-1})$$

- ▶  $\alpha > 0$  and set so that  $\mathcal{L}(\mathcal{T}; \omega^i) > \mathcal{L}(\mathcal{T}; \omega^{i-1})$
- ▶  $\nabla \mathcal{L}(\mathcal{T}; \omega)$  is gradient of  $\mathcal{L}$  w.r.t.  $\omega$ 
  - ▶ A gradient is all partial derivatives over variables  $w_i$
  - ▶ i.e.,  $\nabla \mathcal{L}(\mathcal{T}; \omega) = (\frac{\partial}{\partial \omega_0} \mathcal{L}(\mathcal{T}; \omega), \frac{\partial}{\partial \omega_1} \mathcal{L}(\mathcal{T}; \omega), \dots, \frac{\partial}{\partial \omega_m} \mathcal{L}(\mathcal{T}; \omega))$
- ▶ Gradient ascent will always find  $\omega$  to maximize  $\mathcal{L}$

# Gradient Descent

- ▶ Let  $\mathcal{L}(\mathcal{T}; \omega) = - \sum_{t=1}^{|\mathcal{T}|} \log (e^{\omega \cdot \phi(x_t, y_t)} / Z_x)$
- ▶ Want to find  $\arg \min_{\omega} \mathcal{L}(\mathcal{T}; \omega)$ 
  - ▶ Set  $\omega^0 = O^m$
  - ▶ Iterate until convergence

$$\omega^i = \omega^{i-1} - \alpha \nabla \mathcal{L}(\mathcal{T}; \omega^{i-1})$$

- ▶  $\alpha > 0$  and set so that  $\mathcal{L}(\mathcal{T}; \omega^i) < \mathcal{L}(\mathcal{T}; \omega^{i-1})$
- ▶  $\nabla \mathcal{L}(\mathcal{T}; \omega)$  is gradient of  $\mathcal{L}$  w.r.t.  $\omega$ 
  - ▶ A gradient is all partial derivatives over variables  $w_i$
  - ▶ i.e.,  $\nabla \mathcal{L}(\mathcal{T}; \omega) = (\frac{\partial}{\partial \omega_0} \mathcal{L}(\mathcal{T}; \omega), \frac{\partial}{\partial \omega_1} \mathcal{L}(\mathcal{T}; \omega), \dots, \frac{\partial}{\partial \omega_m} \mathcal{L}(\mathcal{T}; \omega))$
- ▶ Gradient descent will always find  $\omega$  to minimize  $\mathcal{L}$

# The partial derivatives

- ▶ Need to find all partial derivatives  $\frac{\partial}{\partial \omega_i} \mathcal{L}(\mathcal{T}; \omega)$

$$\begin{aligned}\mathcal{L}(\mathcal{T}; \omega) &= \sum_t \log P(\mathbf{y}_t | \mathbf{x}_t) \\ &= \sum_t \log \frac{e^{\omega \cdot \phi(\mathbf{x}_t, \mathbf{y}_t)}}{\sum_{\mathbf{y}' \in \mathcal{Y}} e^{\omega \cdot \phi(\mathbf{x}_t, \mathbf{y}')}} \\ &= \sum_t \log \frac{e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)}}{Z_{\mathbf{x}_t}}\end{aligned}$$

## Partial derivatives - some reminders

$$1. \frac{\partial}{\partial x} \log F = \frac{1}{F} \frac{\partial}{\partial x} F$$

► We always assume log is the natural logarithm  $\log_e$

$$2. \frac{\partial}{\partial x} e^F = e^F \frac{\partial}{\partial x} F$$

$$3. \frac{\partial}{\partial x} \sum_t F_t = \sum_t \frac{\partial}{\partial x} F_t$$

$$4. \frac{\partial}{\partial x} \frac{F}{G} = \frac{G \frac{\partial}{\partial x} F - F \frac{\partial}{\partial x} G}{G^2}$$

# The partial derivatives

$$\frac{\partial}{\partial \omega_i} \mathcal{L}(\mathcal{T}; \omega) =$$

# The partial derivatives 1 (for handout)

$$\begin{aligned}
 \frac{\partial}{\partial \omega_i} \mathcal{L}(\mathcal{T}; \omega) &= \frac{\partial}{\partial \omega_i} \sum_t \log \frac{e^{\sum_j \omega_j \times \phi_j(x_t, y_t)}}{Z_{x_t}} \\
 &= \sum_t \frac{\partial}{\partial \omega_i} \log \frac{e^{\sum_j \omega_j \times \phi_j(x_t, y_t)}}{Z_{x_t}} \\
 &= \sum_t \left( \frac{Z_{x_t}}{e^{\sum_j \omega_j \times \phi_j(x_t, y_t)}} \right) \left( \frac{\partial}{\partial \omega_i} \frac{e^{\sum_j \omega_j \times \phi_j(x_t, y_t)}}{Z_{x_t}} \right)
 \end{aligned}$$



# The partial derivatives

Now,  $\frac{\partial}{\partial \omega_j} \frac{e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)}}{Z_{\mathbf{x}_t}} =$

## The partial derivatives 2 (for handout)

Now,

$$\begin{aligned}
 \frac{\partial}{\partial \omega_i} \frac{e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)}}{Z_{\mathbf{x}_t}} &= \frac{Z_{\mathbf{x}_t} \frac{\partial}{\partial \omega_i} e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)} - e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)} \frac{\partial}{\partial \omega_i} Z_{\mathbf{x}_t}}{Z_{\mathbf{x}_t}^2} \\
 &= \frac{Z_{\mathbf{x}_t} e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)} \phi_i(\mathbf{x}_t, \mathbf{y}_t) - e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)} \frac{\partial}{\partial \omega_i} Z_{\mathbf{x}_t}}{Z_{\mathbf{x}_t}^2} \\
 &= \frac{e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)}}{Z_{\mathbf{x}_t}^2} (Z_{\mathbf{x}_t} \phi_i(\mathbf{x}_t, \mathbf{y}_t) - \frac{\partial}{\partial \omega_i} Z_{\mathbf{x}_t}) \\
 &= \frac{e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)}}{Z_{\mathbf{x}_t}^2} (Z_{\mathbf{x}_t} \phi_i(\mathbf{x}_t, \mathbf{y}_t) \\
 &\quad - \sum_{\mathbf{y}' \in \mathcal{Y}} e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}')} \phi_i(\mathbf{x}_t, \mathbf{y}'))
 \end{aligned}$$

because

$$\frac{\partial}{\partial \omega_i} Z_{\mathbf{x}_t} = \frac{\partial}{\partial \omega_i} \sum_{\mathbf{y}' \in \mathcal{Y}} e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}')} = \sum_{\mathbf{y}' \in \mathcal{Y}} e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}')} \phi_i(\mathbf{x}_t, \mathbf{y}')$$

# The partial derivatives

## The partial derivatives 3 (for handout)

From before,

$$\frac{\partial}{\partial \omega_i} \frac{e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)}}{Z_{\mathbf{x}_t}} = \frac{e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)}}{Z_{\mathbf{x}_t}^2} (Z_{\mathbf{x}_t} \phi_i(\mathbf{x}_t, \mathbf{y}_t) - \sum_{\mathbf{y}' \in \mathcal{Y}} e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}')} \phi_i(\mathbf{x}_t, \mathbf{y}'))$$

Sub this in,

$$\begin{aligned} \frac{\partial}{\partial \omega_i} \mathcal{L}(\mathcal{T}; \omega) &= \sum_t \left( \frac{Z_{\mathbf{x}_t}}{e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)}} \right) \left( \frac{\partial}{\partial \omega_i} \frac{e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}_t)}}{Z_{\mathbf{x}_t}} \right) \\ &= \sum_t \frac{1}{Z_{\mathbf{x}_t}} (Z_{\mathbf{x}_t} \phi_i(\mathbf{x}_t, \mathbf{y}_t) - \sum_{\mathbf{y}' \in \mathcal{Y}} e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}')} \phi_i(\mathbf{x}_t, \mathbf{y}')) \\ &= \sum_t \phi_i(\mathbf{x}_t, \mathbf{y}_t) - \sum_t \sum_{\mathbf{y}' \in \mathcal{Y}} \frac{e^{\sum_j \omega_j \times \phi_j(\mathbf{x}_t, \mathbf{y}')}}{Z_{\mathbf{x}_t}} \phi_i(\mathbf{x}_t, \mathbf{y}') \\ &= \sum_t \phi_i(\mathbf{x}_t, \mathbf{y}_t) - \sum_t \sum_{\mathbf{y}' \in \mathcal{Y}} P(\mathbf{y}' | \mathbf{x}_t) \phi_i(\mathbf{x}_t, \mathbf{y}') \end{aligned}$$

# FINALLY!!!

- ▶ After all that,

$$\frac{\partial}{\partial \omega_i} \mathcal{L}(\mathcal{T}; \omega) = \sum_t \phi_i(x_t, y_t) - \sum_t \sum_{y' \in \mathcal{Y}} P(y'|x_t) \phi_i(x_t, y')$$

- ▶ And the gradient is:

$$\nabla \mathcal{L}(\mathcal{T}; \omega) = \left( \frac{\partial}{\partial \omega_0} \mathcal{L}(\mathcal{T}; \omega), \frac{\partial}{\partial \omega_1} \mathcal{L}(\mathcal{T}; \omega), \dots, \frac{\partial}{\partial \omega_m} \mathcal{L}(\mathcal{T}; \omega) \right)$$

- ▶ So we can now use gradient ascent to find  $\omega$ !!

# Logistic Regression Summary

- ▶ Define conditional probability

$$P(\mathbf{y}|\mathbf{x}) = \frac{e^{\boldsymbol{\omega} \cdot \boldsymbol{\phi}(\mathbf{x}, \mathbf{y})}}{Z_{\mathbf{x}}}$$

- ▶ Set weights to maximize log-likelihood of training data:

$$\boldsymbol{\omega} = \arg \max_{\boldsymbol{\omega}} \sum_t \log P(\mathbf{y}_t | \mathbf{x}_t)$$

- ▶ Can find the gradient and run gradient ascent (or any gradient-based optimization algorithm)

$$\frac{\partial}{\partial \omega_i} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) = \sum_t \phi_i(\mathbf{x}_t, \mathbf{y}_t) - \sum_t \sum_{\mathbf{y}' \in \mathcal{Y}} P(\mathbf{y}' | \mathbf{x}_t) \phi_i(\mathbf{x}_t, \mathbf{y}')$$

# Logistic Regression = Maximum Entropy

More here: <https://purenanya.wordpress.com/2012/09/05/logistic-regression-and-the-maximum-entropy-modeling->

- ▶ Max Ent: maximize entropy subject to constraints on features:  $P = \arg \max_P H(P)$  under constraints
  - ▶ Empirical feature counts must equal expected counts
- ▶ Quick intuition

Matching  
1st  
moments  
(only  
average)  
other  
according  
occams  
razor

Partial derivative in logistic regression

$$\frac{\partial}{\partial \omega_i} \mathcal{L}(\mathcal{T}; \omega) = \sum_t \phi_i(x_t, y_t) - \sum_t \sum_{y' \in \mathcal{Y}} P(y'|x_t) \phi_i(x_t, y')$$

First term is empirical feature counts and second term is expected counts

Derivative set to zero maximizes function

Therefore when both counts are equivalent, we optimize the logistic regression objective!

# Perceptron



# Perceptron

- ▶ Choose a  $\omega$  that minimizes error

$$\mathcal{L}(\mathcal{T}; \omega) = \sum_{t=1}^{|\mathcal{T}|} 1 - [[y_t = \arg \max_y \omega \cdot \phi(x_t, y)]]$$

$$\omega = \arg \min_{\omega} \sum_{t=1}^{|\mathcal{T}|} 1 - [[y_t = \arg \max_y \omega \cdot \phi(x_t, y)]]$$

$$[[p]] = \begin{cases} 1 & p \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

- ▶ This is a 0-1 loss function
  - ▶ When minimizing error people tend to use **hinge-loss**
  - ▶ We'll get back to this

## Aside: Min error versus max log-likelihood

- ▶ Highly related but not identical
- ▶ Example: consider a training set  $\mathcal{T}$  with 1001 points

$$\begin{aligned}1000 \times (\mathbf{x}_i, \mathbf{y} = 0) &= [-1, 1, 0, 0] \quad \text{for } i = 1 \dots 1000 \\1 \times (\mathbf{x}_{1001}, \mathbf{y} = 1) &= [0, 0, 3, 1]\end{aligned}$$

- ▶ Now consider  $\boldsymbol{\omega} = [-1, 0, 1, 0]$
- ▶ Error in this case is 0 – so  $\boldsymbol{\omega}$  minimizes error

$$[-1, 0, 1, 0] \cdot [-1, 1, 0, 0] = 1 > [-1, 0, 1, 0] \cdot [0, 0, -1, 1] = -1$$

$$[-1, 0, 1, 0] \cdot [0, 0, 3, 1] = 3 > [-1, 0, 1, 0] \cdot [3, 1, 0, 0] = -3$$

- ▶ However, log-likelihood = -126.9 (omit calculation)

## Aside: Min error versus max log-likelihood

- ▶ Highly related but not identical
- ▶ Example: consider a training set  $\mathcal{T}$  with 1001 points

$$1000 \times (\mathbf{x}_i, \mathbf{y} = 0) = [-1, 1, 0, 0] \quad \text{for } i = 1 \dots 1000$$
$$1 \times (\mathbf{x}_{1001}, \mathbf{y} = 1) = [0, 0, 3, 1]$$

- ▶ Now consider  $\boldsymbol{\omega} = [-1, 7, 1, 0]$
- ▶ Error in this case is 1 – so  $\boldsymbol{\omega}$  does not minimize error

$$[-1, 7, 1, 0] \cdot [-1, 1, 0, 0] = 8 > [-1, 7, 1, 0] \cdot [-1, 1, 0, 0] = -1$$

$$[-1, 7, 1, 0] \cdot [0, 0, 3, 1] = 3 < [-1, 7, 1, 0] \cdot [3, 1, 0, 0] = 4$$

- ▶ However, log-likelihood = -1.4
- ▶ Better log-likelihood and worse error

## Aside: Min error versus max log-likelihood

- ▶ Max likelihood  $\neq$  min error
- ▶ Max likelihood pushes as much probability on correct labeling of training instance
  - ▶ Even at the cost of mislabeling a few examples
- ▶ Min error forces all training instances to be correctly classified
  - ▶ Often not possible
  - ▶ Ways of regularizing model to allow sacrificing some errors for better predictions on more examples

# Perceptron Learning Algorithm

Training data:  $\mathcal{T} = \{(x_t, y_t)\}_{t=1}^{|\mathcal{T}|}$

1.  $\omega^{(0)} = 0; i = 0$
2. for  $n : 1..N$
3.   for  $t : 1..T$
4.     Let  $y' = \arg \max_{y'} \omega^{(i)} \cdot \phi(x_t, y')$
5.     if  $y' \neq y_t$
6.        $\omega^{(i+1)} = \omega^{(i)} + \phi(x_t, y_t) - \phi(x_t, y')$
7.        $i = i + 1$
8. return  $\omega^i$

# Perceptron: Separability and Margin

- ▶ Given an training instance  $(x_t, y_t)$ , define:
  - ▶  $\bar{\mathcal{Y}}_t = \mathcal{Y} - \{y_t\}$
  - ▶ i.e.,  $\bar{\mathcal{Y}}_t$  is the set of incorrect labels for  $x_t$
- ▶ A training set  $\mathcal{T}$  is separable with margin  $\gamma > 0$  if there exists a vector  $\mathbf{u}$  with  $\|\mathbf{u}\| = 1$  such that:

$$\mathbf{u} \cdot \phi(x_t, y_t) - \mathbf{u} \cdot \phi(x_t, y') \geq \gamma$$

for all  $y' \in \bar{\mathcal{Y}}_t$  and  $\|\mathbf{u}\| = \sqrt{\sum_j \mathbf{u}_j^2}$

- ▶ **Assumption:** the training set is separable with margin  $\gamma$

# Perceptron: Main Theorem

- ▶ **Theorem:** For any training set separable with a margin of  $\gamma$ , the following holds for the perceptron algorithm:

$$\text{mistakes made during training} \leq \frac{R^2}{\gamma^2}$$

where  $R \geq \|\phi(\mathbf{x}_t, \mathbf{y}_t) - \phi(\mathbf{x}_t, \mathbf{y}')\|$  for all  $(\mathbf{x}_t, \mathbf{y}_t) \in \mathcal{T}$  and  $\mathbf{y}' \in \bar{\mathcal{Y}}_t$

- ▶ Thus, after a finite number of training iterations, the error on the training set will converge to zero
- ▶ **Let's prove it!** (proof taken from Collins '02)

# Perceptron Learning Algorithm

Training data:  $\mathcal{T} = \{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^{|\mathcal{T}|}$

1.  $\omega^{(0)} = 0; i = 0$
2. for  $n : 1..N$
3.   for  $t : 1..T$
4.     Let  $\mathbf{y}' = \arg \max_{\mathbf{y}'} \omega^{(i)} \cdot \phi(\mathbf{x}_t, \mathbf{y}')$
5.     if  $\mathbf{y}' \neq \mathbf{y}_t$
6.        $\omega^{(i+1)} = \omega^{(i)} + \phi(\mathbf{x}_t, \mathbf{y}_t) - \phi(\mathbf{x}_t, \mathbf{y}')$
7.        $i = i + 1$
8. return  $\omega^i$

- ▶  $\omega^{(k-1)}$  are the weights before  $k^{th}$  mistake
- ▶ Suppose  $k^{th}$  mistake made at the  $t^{th}$  example,  $(\mathbf{x}_t, \mathbf{y}_t)$
- ▶  $\mathbf{y}' = \arg \max_{\mathbf{y}'} \omega^{(k-1)} \cdot \phi(\mathbf{x}_t, \mathbf{y}')$
- ▶  $\mathbf{y}' \neq \mathbf{y}_t$
- ▶  $\omega^{(k)} = \omega^{(k-1)} + \phi(\mathbf{x}_t, \mathbf{y}_t) - \phi(\mathbf{x}_t, \mathbf{y}')$





# Perceptron Learning Algorithm (for handout)

Training data:  $\mathcal{T} = \{(x_t, y_t)\}_{t=1}^{|\mathcal{T}|}$

1.  $\omega^{(0)} = 0; i = 0$
2. for  $n : 1..N$
3.   for  $t : 1..T$
4.     Let  $y' = \arg \max_{y'} \omega^{(i)} \cdot \phi(x_t, y')$
5.     if  $y' \neq y_t$
6.          $\omega^{(i+1)} = \omega^{(i)} + \phi(x_t, y_t) - \phi(x_t, y')$
7.          $i = i + 1$
8. return  $\omega^i$

►  $\omega^{(k-1)}$  are the weights before  $k^{th}$  mistake

► Suppose  $k^{th}$  mistake made at the  $t^{th}$  example,  $(x_t, y_t)$

►  $y' = \arg \max_{y'} \omega^{(k-1)} \cdot \phi(x_t, y')$

►  $y' \neq y_t$

►  $\omega^{(k)} = \omega^{(k-1)} + \phi(x_t, y_t) - \phi(x_t, y')$

► Now:  $\mathbf{u} \cdot \omega^{(k)} = \mathbf{u} \cdot \omega^{(k-1)} + \mathbf{u} \cdot (\phi(x_t, y_t) - \phi(x_t, y')) \geq \mathbf{u} \cdot \omega^{(k-1)} + \gamma$

► Now:  $\omega^{(0)} = 0$  and  $\mathbf{u} \cdot \omega^{(0)} = 0$ , by induction on  $k$ ,  $\mathbf{u} \cdot \omega^{(k)} \geq k\gamma$

► Now: since  $\mathbf{u} \cdot \omega^{(k)} \leq \|\mathbf{u}\| \times \|\omega^{(k)}\|$  and  $\|\mathbf{u}\| = 1$  then  $\|\omega^{(k)}\| \geq k\gamma$

► Now:

$$\|\omega^{(k)}\|^2 = \|\omega^{(k-1)}\|^2 + \|\phi(x_t, y_t) - \phi(x_t, y')\|^2 + 2\omega^{(k-1)} \cdot (\phi(x_t, y_t) - \phi(x_t, y'))$$

$$\|\omega^{(k)}\|^2 \leq \|\omega^{(k-1)}\|^2 + R^2$$

(since  $R \geq \|\phi(x_t, y_t) - \phi(x_t, y')\|$

and  $\omega^{(k-1)} \cdot \phi(x_t, y_t) - \omega^{(k-1)} \cdot \phi(x_t, y') \leq 0$ )

# Perceptron Learning Algorithm

- ▶ We have just shown that  $\|\omega^{(k)}\| \geq k\gamma$  and  $\|\omega^{(k)}\|^2 \leq \|\omega^{(k-1)}\|^2 + R^2$
- ▶ By induction on  $k$  and since  $\omega^{(0)} = 0$  and  $\|\omega^{(0)}\|^2 = 0$
- ▶ Therefore,
- ▶ and solving for  $k$
- ▶ Therefore the number of errors is bounded!

# Perceptron Learning Algorithm (for handout)

- ▶ We have just shown that  $\|\omega^{(k)}\| \geq k\gamma$  and  $\|\omega^{(k)}\|^2 \leq \|\omega^{(k-1)}\|^2 + R^2$
- ▶ By induction on  $k$  and since  $\omega^{(0)} = 0$  and  $\|\omega^{(0)}\|^2 = 0$

$$\|\omega^{(k)}\|^2 \leq kR^2$$

- ▶ Therefore,

$$k^2\gamma^2 \leq \|\omega^{(k)}\|^2 \leq kR^2$$

- ▶ and solving for  $k$

$$k \leq \frac{R^2}{\gamma^2}$$

- ▶ Therefore the number of errors is bounded!

# Perceptron Summary

- ▶ Learns a linear classifier that minimizes error
- ▶ Guaranteed to find a  $\omega$  in a finite amount of time
- ▶ Perceptron is an example of an **Online Learning Algorithm**
  - ▶  $\omega$  is updated based on a single training instance in isolation

$$\omega^{(i+1)} = \omega^{(i)} + \phi(x_t, y_t) - \phi(x_t, y')$$

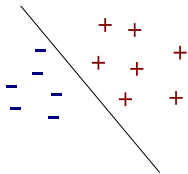
# Averaged Perceptron

Training data:  $\mathcal{T} = \{(x_t, y_t)\}_{t=1}^{|\mathcal{T}|}$

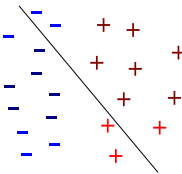
1.  $\omega^{(0)} = 0; i = 0$
2. for  $n : 1..N$
3.   for  $t : 1..T$
4.     Let  $y' = \arg \max_{y'} \omega^{(i)} \cdot \phi(x_t, y')$
5.     if  $y' \neq y_t$
6.        $\omega^{(i+1)} = \omega^{(i)} + \phi(x_t, y_t) - \phi(x_t, y')$
7.     else
6.        $\omega^{(i+1)} = \omega^{(i)}$
7.      $i = i + 1$
8. return  $(\sum_i \omega^{(i)}) / (N \times T)$

# Margin

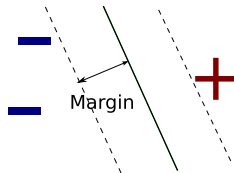
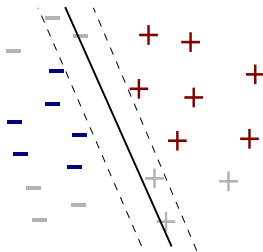
Training



Testing



Denote the value of the margin by  $\gamma$



# Maximizing Margin

- ▶ For a training set  $\mathcal{T}$
- ▶ Margin of a weight vector  $\omega$  is smallest  $\gamma$  such that

$$\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') \geq \gamma$$

- ▶ for every training instance  $(x_t, y_t) \in \mathcal{T}$ ,  $y' \in \bar{\mathcal{Y}}_t$

# Maximizing Margin

- ▶ Intuitively maximizing margin makes sense
- ▶ More importantly, generalization error to unseen test data is proportional to the inverse of the margin

$$\epsilon \propto \frac{R^2}{\gamma^2 \times |\mathcal{T}|}$$

- ▶ **Perceptron:** we have shown that:
  - ▶ If a training set is separable by some margin, the perceptron will find a  $\omega$  that separates the data
  - ▶ However, the perceptron does not pick  $\omega$  to maximize the margin!



# Support Vector Machines (SVMs)

# Maximizing Margin

Let  $\gamma > 0$

$$\max_{\|\omega\|=1} \gamma$$

such that:

$$\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') \geq \gamma$$

$$\forall (x_t, y_t) \in \mathcal{T}$$

$$\text{and } y' \in \bar{\mathcal{Y}}_t$$

- ▶ Note: algorithm still **minimizes error** if data is separable
- ▶  $\|\omega\|$  is bound since scaling trivially produces larger margin

$$\beta(\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y')) \geq \beta\gamma, \text{ for some } \beta \geq 1$$

# Max Margin = Min Norm

Let  $\gamma > 0$

**Max Margin:**

$$\max_{\|\omega\|=1} \gamma$$

such that:

$$\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') \geq \gamma$$

$$\forall (x_t, y_t) \in \mathcal{T}$$

$$\text{and } y' \in \bar{\mathcal{Y}}_t$$

Change of variable:  $u = \frac{\omega}{\gamma}$ ?

$$\|\omega\| = 1 \text{ iff } \|u\| = 1/\gamma$$

**Min Norm (step 1):**

$$\max_{\|u\|=1/\gamma} \gamma$$

such that:

$$\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') \geq \gamma$$

$$\forall (x_t, y_t) \in \mathcal{T}$$

$$\text{and } y' \in \bar{\mathcal{Y}}_t$$

# Max Margin = Min Norm

Let  $\gamma > 0$

**Max Margin:**

$$\max_{\|\omega\|=1} \gamma$$

such that:

$$\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') \geq \gamma$$

$$\forall (x_t, y_t) \in \mathcal{T}$$

$$\text{and } y' \in \bar{\mathcal{Y}}_t$$

Change variables:  $u = \frac{\omega}{\gamma}$ ?

$$\|\omega\| = 1 \text{ iff } \|u\| = \gamma$$

**Min Norm (step 2):**

$$\max_{\|u\|=1/\gamma} \gamma$$

such that:

$$\gamma u \cdot \phi(x_t, y_t) - \gamma u \cdot \phi(x_t, y') \geq \gamma$$

$$\forall (x_t, y_t) \in \mathcal{T}$$

$$\text{and } y' \in \bar{\mathcal{Y}}_t$$

# Max Margin = Min Norm

Let  $\gamma > 0$

**Max Margin:**

$$\max_{\|\omega\|=1} \gamma$$

such that:

$$\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') \geq \gamma$$

$$\forall (x_t, y_t) \in \mathcal{T}$$

$$\text{and } y' \in \bar{\mathcal{Y}}_t$$

Change variables:  $u = \frac{w}{\gamma}$ ?

$$\|\omega\| = 1 \text{ iff } \|\mathbf{u}\| = \gamma$$

**Min Norm (step 3):**

$$\max_{\|\mathbf{u}\|=1/\gamma} \gamma$$

such that:

$$\mathbf{u} \cdot \phi(x_t, y_t) - \mathbf{u} \cdot \phi(x_t, y') \geq 1$$

$$\forall (x_t, y_t) \in \mathcal{T}$$

$$\text{and } y' \in \bar{\mathcal{Y}}_t$$

But  $\gamma$  is really not constrained!

# Max Margin = Min Norm

Let  $\gamma > 0$

**Max Margin:**

$$\max_{\|\omega\|=1} \gamma$$

such that:

$$\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') \geq \gamma$$

$$\forall (x_t, y_t) \in \mathcal{T}$$

$$\text{and } y' \in \bar{\mathcal{Y}}_t$$

Change variables:  $u = \frac{\omega}{\gamma}$ ?

$$\|\omega\| = 1 \text{ iff } \|u\| = \gamma$$

**Min Norm (step 4):**

$$\max_u \frac{1}{\|u\|} = \min_u \|u\|$$

such that:

$$u \cdot \phi(x_t, y_t) - u \cdot \phi(x_t, y') \geq 1$$

$$\forall (x_t, y_t) \in \mathcal{T}$$

$$\text{and } y' \in \bar{\mathcal{Y}}_t$$

But  $\gamma$  is really not constrained!

# Max Margin = Min Norm

Let  $\gamma > 0$

**Max Margin:**

$$\max_{\|\omega\|=1} \gamma$$

such that:

$$\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') \geq \gamma$$

$$\forall (x_t, y_t) \in \mathcal{T}$$

$$\text{and } y' \in \bar{\mathcal{Y}}_t$$

**Min Norm:**

$$\min_{\mathbf{u}} \frac{1}{2} \|\mathbf{u}\|^2$$

such that:

$$\mathbf{u} \cdot \phi(x_t, y_t) - \mathbf{u} \cdot \phi(x_t, y') \geq 1$$

$$\forall (x_t, y_t) \in \mathcal{T}$$

$$\text{and } y' \in \bar{\mathcal{Y}}_t$$

- Intuition: Instead of fixing  $\|\omega\|$  we fix the margin  $\gamma = 1$

# Support Vector Machines

$$\omega = \arg \min_{\omega} \frac{1}{2} \|\omega\|^2$$

such that:

$$\begin{aligned} \omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') &\geq 1 \\ \forall (x_t, y_t) \in \mathcal{T} \text{ and } y' \in \bar{\mathcal{Y}}_t \end{aligned}$$

- ▶ Quadratic programming problem – a well-known convex optimization problem
- ▶ Can be solved with many techniques [Nocedal and Wright 1999]



# Support Vector Machines

What if data is not separable? (Original problem: will not satisfy the constraints!)

$$\omega = \arg \min_{\omega, \xi} \frac{1}{2} \|\omega\|^2 + \textcolor{red}{C} \sum_{t=1}^{|\mathcal{T}|} \xi_t$$

such that:

$$\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') \geq \textcolor{red}{1} - \xi_t \text{ and } \xi_t \geq 0$$

$$\forall (x_t, y_t) \in \mathcal{T} \text{ and } y' \in \bar{\mathcal{Y}}_t$$

$\xi_t$ : trade-off between margin per example and  $\|\omega\|$

Larger  $\textcolor{red}{C}$  = more examples correctly classified

If data is separable, optimal solution has  $\xi_i = 0, \forall i$

# Support Vector Machines

$$\omega = \arg \min_{\omega, \xi} \frac{\lambda}{2} \|\omega\|^2 + \sum_{t=1}^{|\mathcal{T}|} \xi_t \quad \lambda = \frac{1}{C}$$

such that:

$$\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') \geq 1 - \xi_t$$

Can we have a more compact representation of this objective function?

$$\omega \cdot \phi(x_t, y_t) - \max_{y' \neq y_t} \omega \cdot \phi(x_t, y') \geq 1 - \xi_t$$

$$\xi_t \geq 1 + \underbrace{\max_{y' \neq y_t} \omega \cdot \phi(x_t, y') - \omega \cdot \phi(x_t, y_t)}_{\text{negated margin for example}}$$

# Support Vector Machines

$$\xi_t \geq 1 + \underbrace{\max_{\mathbf{y}' \neq \mathbf{y}_t} \omega \cdot \phi(\mathbf{x}_t, \mathbf{y}') - \omega \cdot \phi(\mathbf{x}_t, \mathbf{y}_t)}_{\text{negated margin for example}}$$

- ▶ If  $\|\omega\|$  classifies  $(\mathbf{x}_t, \mathbf{y}_t)$  with margin 1, penalty  $\xi_t = 0$
- ▶ (Objective wants to keep  $\xi_t$  small and  $\xi_t = 0$  satisfies the constraint)
- ▶ Otherwise:  $\xi_t = 1 + \max_{\mathbf{y}' \neq \mathbf{y}_t} \omega \cdot \phi(\mathbf{x}_t, \mathbf{y}') - \omega \cdot \phi(\mathbf{x}_t, \mathbf{y}_t)$
- ▶ (Again, because that's the minimal  $\xi_t$  that satisfies the constraint, and we want  $\xi_t$  smallest as possible)
- ▶ That means that in the end  $\xi_t$  will be:

$$\xi_t = \max\{0, 1 + \max_{\mathbf{y}' \neq \mathbf{y}_t} \omega \cdot \phi(\mathbf{x}_t, \mathbf{y}') - \omega \cdot \phi(\mathbf{x}_t, \mathbf{y}_t)\}$$

(If an example is classified correctly,  $\xi_t = 0$  and the second term in the max is negative.)

# Support Vector Machines

$$\omega = \arg \min_{\omega, \xi} \frac{\lambda}{2} \|\omega\|^2 + \sum_{t=1}^{|\mathcal{T}|} \xi_t$$

such that:

$$\xi_t \geq 1 + \max_{y' \neq y_t} \omega \cdot \phi(x_t, y') - \omega \cdot \phi(x_t, y_t)$$

Hinge loss equivalent

$$\begin{aligned} \omega &= \arg \min_{\omega} \mathcal{L}(\mathcal{T}; \omega) = \arg \min_{\omega} \sum_{t=1}^{|\mathcal{T}|} \text{loss}((x_t, y_t); \omega) + \frac{\lambda}{2} \|\omega\|^2 \\ &= \arg \min_{\omega} \left( \sum_{t=1}^{|\mathcal{T}|} \max \left( 0, 1 + \max_{y' \neq y_t} \omega \cdot \phi(x_t, y') - \omega \cdot \phi(x_t, y_t) \right) \right) + \frac{\lambda}{2} \|\omega\|^2 \end{aligned}$$

# Summary

## What we have covered

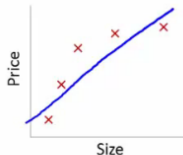
- ▶ Linear Classifiers
  - ▶ Naive Bayes
  - ▶ Logistic Regression
  - ▶ Perceptron
  - ▶ Support Vector Machines

## What is next

- ▶ Regularization
- ▶ Online learning
- ▶ Non-linear classifiers

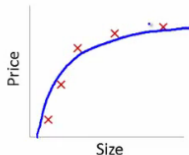
# Regularization

# Fit of a Model



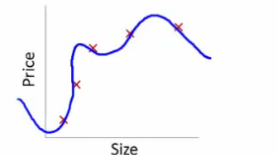
$$\theta_0 + \theta_1 x$$

High bias  
(underfit)



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

"Just right"



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

High variance  
(overfit)

- ▶ Two sources of error:
  - ▶ Bias error, measures how well the hypothesis class fits the space we are trying to model
  - ▶ Variance error, measures sensitivity to training set selection
  - ▶ Want to balance these two things

# Overfitting

- ▶ Early in lecture we made assumption data was i.i.d.
- ▶ Rarely is this true
  - ▶ E.g., syntactic analyzers typically trained on 40,000 sentences from early 1990s WSJ news text
- ▶ Even more common:  $\mathcal{T}$  is very small
- ▶ This leads to **overfitting**
- ▶ E.g.: 'fake' is never a verb in WSJ treebank (only adjective)
  - ▶ High weight on " $\phi(x, y) = 1$  if  $x$ =fake and  $y$ =adjective"
  - ▶ Of course: leads to high log-likelihood / low error
- ▶ Other features might be more indicative
- ▶ Adjacent word identities: 'He wants to X his death'  $\rightarrow$  X=verb



# Regularization

- ▶ In practice, we **regularize** models to prevent overfitting

$$\arg \max_{\omega} \mathcal{L}(\mathcal{T}; \omega) - \lambda \mathcal{R}(\omega)$$

- ▶ Where  $\mathcal{R}(\omega)$  is the regularization function
- ▶  $\lambda$  controls how much to regularize
- ▶ Common functions
  - ▶ L2:  $\mathcal{R}(\omega) \propto \|\omega\|_2 = \|\omega\| = \sqrt{\sum_i \omega_i^2}$  – smaller weights desired
  - ▶ L0:  $\mathcal{R}(\omega) \propto \|\omega\|_0 = \sum_i [[\omega_i > 0]]$  – zero weights desired
    - ▶ Non-convex
    - ▶ Approximate with L1:  $\mathcal{R}(\omega) \propto \|\omega\|_1 = \sum_i |\omega_i|$

# Logistic Regression with L2 Regularization

- ▶ Perhaps most common classifier in NLP

$$\mathcal{L}(\mathcal{T}; \omega) - \lambda \mathcal{R}(\omega) = \sum_{t=1}^{|\mathcal{T}|} \log \left( e^{\omega \cdot \phi(x_t, y_t)} / Z_x \right) - \frac{\lambda}{2} \|\omega\|^2$$

- ▶ What are the new partial derivatives?

$$\frac{\partial}{\partial w_i} \mathcal{L}(\mathcal{T}; \omega) - \frac{\partial}{\partial w_i} \lambda \mathcal{R}(\omega)$$

- ▶ We know  $\frac{\partial}{\partial w_i} \mathcal{L}(\mathcal{T}; \omega)$

- ▶ Just need  $\frac{\partial}{\partial w_i} \frac{\lambda}{2} \|\omega\|^2 = \frac{\partial}{\partial w_i} \frac{\lambda}{2} \left( \sqrt{\sum_i \omega_i^2} \right)^2 = \frac{\partial}{\partial w_i} \frac{\lambda}{2} \sum_i \omega_i^2 = \lambda \omega_i$

# Support Vector Machines

Hinge-loss formulation: L2 regularization already happening!

$$\begin{aligned}
 \omega &= \arg \min_{\omega} \mathcal{L}(\mathcal{T}; \omega) + \lambda \mathcal{R}(\omega) \\
 &= \arg \min_{\omega} \sum_{t=1}^{|\mathcal{T}|} \text{loss}((x_t, y_t); \omega) + \lambda \mathcal{R}(\omega) \\
 &= \arg \min_{\omega} \sum_{t=1}^{|\mathcal{T}|} \max(0, 1 + \max_{y \neq y_t} \omega \cdot \phi(x_t, y) - \omega \cdot \phi(x_t, y_t)) + \lambda \mathcal{R}(\omega) \\
 &= \arg \min_{\omega} \sum_{t=1}^{|\mathcal{T}|} \max(0, 1 + \max_{y \neq y_t} \omega \cdot \phi(x_t, y) - \omega \cdot \phi(x_t, y_t)) + \frac{\lambda}{2} \|\omega\|^2
 \end{aligned}$$

↑ SVM optimization ↑

# SVMs vs. Logistic Regression

$$\begin{aligned}\omega &= \arg \min_{\omega} \mathcal{L}(\mathcal{T}; \omega) + \lambda \mathcal{R}(\omega) \\ &= \arg \min_{\omega} \sum_{t=1}^{|\mathcal{T}|} \text{loss}(\mathbf{x}_t, \mathbf{y}_t; \omega) + \lambda \mathcal{R}(\omega)\end{aligned}$$

SVMs/hinge-loss:  $\max(0, 1 + \max_{\mathbf{y} \neq \mathbf{y}_t} (\omega \cdot \phi(\mathbf{x}_t, \mathbf{y}) - \omega \cdot \phi(\mathbf{x}_t, \mathbf{y}_t)))$

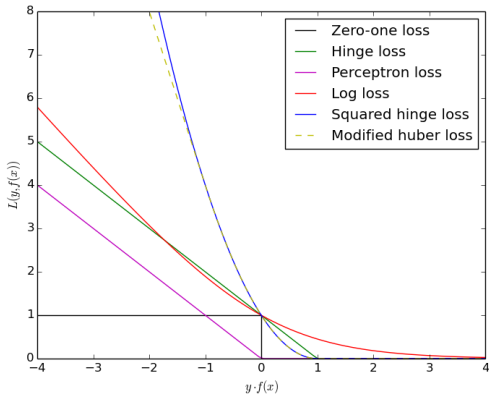
$$\omega = \arg \min_{\omega} \sum_{t=1}^{|\mathcal{T}|} \max(0, 1 + \max_{\mathbf{y} \neq \mathbf{y}_t} \omega \cdot \phi(\mathbf{x}_t, \mathbf{y}) - \omega \cdot \phi(\mathbf{x}_t, \mathbf{y}_t)) + \frac{\lambda}{2} \|\omega\|^2$$

Logistic Regression/**log-loss**:  $-\log(e^{\omega \cdot \phi(\mathbf{x}_t, \mathbf{y}_t)} / Z_{\mathbf{x}})$

$$\omega = \arg \min_{\omega} \sum_{t=1}^{|\mathcal{T}|} -\log(e^{\omega \cdot \phi(\mathbf{x}_t, \mathbf{y}_t)} / Z_{\mathbf{x}}) + \frac{\lambda}{2} \|\omega\|^2$$

# Generalized Linear Classifiers

$$\omega = \arg \min_{\omega} \mathcal{L}(\mathcal{T}; \omega) + \lambda \mathcal{R}(\omega) = \arg \min_{\omega} \sum_{t=1}^{|\mathcal{T}|} \text{loss}((x_t, y_t); \omega) + \lambda \mathcal{R}(\omega)$$



# Which Classifier to Use?

- ▶ Trial and error
- ▶ Training time available
- ▶ Choice of features is often more important

# Online Learning

# Online vs. Batch Learning

Batch( $\mathcal{T}$ );

- ▶ for 1 ...  $N$ 
  - ▶  $\omega \leftarrow \text{update}(\mathcal{T}; \omega)$
- ▶ return  $\omega$

E.g., SVMs, logistic regression, NB

Online( $\mathcal{T}$ );

- ▶ for 1 ...  $N$ 
  - ▶ for  $(x_t, y_t) \in \mathcal{T}$ 
    - ▶  $\omega \leftarrow \text{update}((x_t, y_t); \omega)$
  - ▶ end for
- ▶ end for
- ▶ return  $\omega$

E.g., Perceptron

$$\omega = \omega + \phi(x_t, y_t) - \phi(x_t, y)$$



# Online vs. Batch Learning

- ▶ Online algorithms
  - ▶ Tend to converge more quickly
  - ▶ Often easier to implement
  - ▶ Require more hyperparameter tuning (exception Perceptron)
  - ▶ More unstable convergence
- ▶ Batch algorithms
  - ▶ Tend to converge more slowly
  - ▶ Implementation more complex (quad prog, LBFGs)
  - ▶ Typically more robust to hyperparameters
  - ▶ More stable convergence

# Gradient Descent Reminder

- ▶ Let  $\mathcal{L}(\mathcal{T}; \omega) = \sum_{t=1}^{|\mathcal{T}|} \text{loss}((x_t, y_t); \omega)$ 
  - ▶ Set  $\omega^0 = O^m$
  - ▶ Iterate until convergence

$$\omega^i = \omega^{i-1} - \alpha \nabla \mathcal{L}(\mathcal{T}; \omega^{i-1}) = \omega^{i-1} - \sum_{t=1}^{|\mathcal{T}|} \alpha \nabla \text{loss}((x_t, y_t); \omega^{i-1})$$

- ▶  $\alpha > 0$  and set so that  $\mathcal{L}(\mathcal{T}; \omega^i) < \mathcal{L}(\mathcal{T}; \omega^{i-1})$
- ▶ **Stochastic Gradient Descent (SGD)**
  - ▶ Approximate  $\nabla \mathcal{L}(\mathcal{T}; \omega)$  with single  $\nabla \text{loss}((x_t, y_t); \omega)$

# Stochastic Gradient Descent

- ▶ Let  $\mathcal{L}(\mathcal{T}; \omega) = \sum_{t=1}^{|\mathcal{T}|} \text{loss}((x_t, y_t); \omega)$
- ▶ Set  $\omega^0 = O^m$
- ▶ iterate until convergence
  - ▶ sample  $(x_t, y_t) \in \mathcal{T}$  // “stochastic”
    - ▶  $\omega^i = \omega^{i-1} - \alpha \nabla \text{loss}((x_t, y_t); \omega)$
- ▶ return  $\omega$

In practice

Need to solve  $\nabla \text{loss}((x_t, y_t); \omega)$

- ▶ Set  $\omega^0 = O^m$
- ▶ for  $1 \dots N$ 
  - ▶ for  $(x_t, y_t) \in \mathcal{T}$ 
    - ▶  $\omega^i = \omega^{i-1} - \alpha \nabla \text{loss}((x_t, y_t); \omega)$
- ▶ return  $\omega$

# Online Logistic Regression

- ▶ Stochastic Gradient Descent (SGD)
- ▶  $loss((x_t, y_t); \omega) = \text{log-loss}$
- ▶  $\nabla loss((x_t, y_t); \omega) = \nabla \left( -\log \left( e^{\omega \cdot \phi(x_t, y_t)} / Z_{x_t} \right) \right)$
- ▶ From logistic regression section:

$$\nabla \left( -\log \left( e^{\omega \cdot \phi(x_t, y_t)} / Z_{x_t} \right) \right) = - \left( \phi(x_t, y_t) - \sum_y P(y|x) \phi(x_t, y) \right)$$

- ▶ Plus regularization term (if part of model)

# Online SVMs

- ▶ Stochastic Gradient Descent (SGD)
- ▶  $loss((x_t, y_t); \omega) = \text{hinge-loss}$

$$\nabla loss((x_t, y_t); \omega) = \nabla \left( \max(0, 1 + \max_{y \neq y_t} \omega \cdot \phi(x_t, y) - \omega \cdot \phi(x_t, y_t)) \right)$$

- ▶ Subgradient is:

$$\begin{aligned} & \nabla \left( \max(0, 1 + \max_{y \neq y_t} \omega \cdot \phi(x_t, y) - \omega \cdot \phi(x_t, y_t)) \right) \\ &= \begin{cases} 0, & \text{if } \omega \cdot \phi(x_t, y_t) - \max_y \omega \cdot \phi(x_t, y) \geq 1 \\ \phi(x_t, y) - \phi(x_t, y_t), & \text{otherwise, where } y = \max_y \omega \cdot \phi(x_t, y) \end{cases} \end{aligned}$$

- ▶ Plus regularization term (required for SVMs)

# Perceptron and Hinge-Loss

SVM subgradient update looks like perceptron update

$$\omega^i = \omega^{i-1} - \alpha \begin{cases} 0, & \text{if } \omega \cdot \phi(x_t, y_t) - \max_y \omega \cdot \phi(x_t, y) \geq 1 \\ \phi(x_t, y) - \phi(x_t, y_t), & \text{otherwise, where } y = \max_y \omega \cdot \phi(x_t, y) \end{cases}$$

Perceptron

$$\omega^i = \omega^{i-1} - \alpha \begin{cases} 0, & \text{if } \omega \cdot \phi(x_t, y_t) - \max_y \omega \cdot \phi(x_t, y) \geq 0 \\ \phi(x_t, y) - \phi(x_t, y_t), & \text{otherwise, where } y = \max_y \omega \cdot \phi(x_t, y) \end{cases}$$

where  $\alpha = 1$ , note  $\phi(x_t, y) - \phi(x_t, y_t)$  not  $\phi(x_t, y_t) - \phi(x_t, y)$  since '-' (descent)

Perceptron = SGD with no-margin hinge-loss

$$\max (0, 1 + \max_{y \neq y_t} \omega \cdot \phi(x_t, y) - \omega \cdot \phi(x_t, y_t))$$

# Margin Infused Relaxed Algorithm (MIRA)

Batch (SVMs):

$$\min \frac{1}{2} \|\omega\|^2$$

such that:

$$\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') \geq 1$$

$$\forall (x_t, y_t) \in \mathcal{T} \text{ and } y' \in \bar{\mathcal{Y}}_t$$

Online (MIRA):

Training data:  $\mathcal{T} = \{(x_t, y_t)\}_{t=1}^{|\mathcal{T}|}$

1.  $\omega^{(0)} = 0; i = 0$
2. for  $n : 1..N$
3.   for  $t : 1..T$
4.      $\omega^{(i+1)} = \arg \min_{\omega^*} \|\omega^* - \omega^{(i)}\|$   
       such that:  
        $\omega \cdot \phi(x_t, y_t) - \omega \cdot \phi(x_t, y') \geq 1$   
        $\forall y' \in \bar{\mathcal{Y}}_t$
5.      $i = i + 1$
6. return  $\omega^i$

- MIRA has much smaller optimizations with only  $|\bar{\mathcal{Y}}_t|$  constraints

# Quick Summary



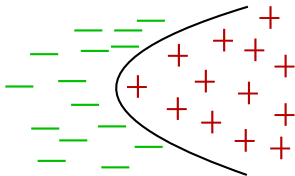
# Linear Classifiers

- ▶ Naive Bayes, Perceptron, Logistic Regression and SVMs
- ▶ Generative vs. Discriminative
- ▶ Objective functions and loss functions
  - ▶ Log-loss, min error and hinge loss
  - ▶ Generalized linear classifiers
- ▶ Regularization
- ▶ Online vs. Batch learning

# Non-linear Classifiers

# Non-Linear Classifiers

- ▶ Some data sets require more than a linear classifier to be correctly modeled
- ▶ Decision boundary is no longer a hyperplane in the feature space
- ▶ A lot of models out there
  - ▶ K-Nearest Neighbours
  - ▶ Decision Trees
  - ▶ Neural Networks
  - ▶ **Kernels**



# Kernels

- ▶ A kernel is a similarity function between two points that is symmetric and positive semi-definite, which we denote by:

$$K(\mathbf{x}_t, \mathbf{x}_r) \in \mathbb{R}$$

- ▶ Let  $M$  be a  $n \times n$  matrix such that ...

$$M_{t,r} = K(\mathbf{x}_t, \mathbf{x}_r)$$

- ▶ ... for any  $n$  points. Called the **Gram matrix**.
- ▶ Symmetric:

$$K(\mathbf{x}_t, \mathbf{x}_r) = K(\mathbf{x}_r, \mathbf{x}_t)$$

- ▶ Positive definite: for all non-zero  $\mathbf{v}$  and any set of  $\mathbf{x}$ s that define a Gram matrix:

$$\mathbf{v}M\mathbf{v}^T \geq 0$$

# Kernels

- ▶ **Mercer's Theorem:** for any kernel  $K$ , there exists an  $\phi$ , in some  $\mathbb{R}^d$ , such that:

$$K(\mathbf{x}_t, \mathbf{x}_r) = \phi(\mathbf{x}_t) \cdot \phi(\mathbf{x}_r)$$

- ▶ Since our features are over pairs  $(\mathbf{x}, \mathbf{y})$ , we will write kernels over pairs

$$K((\mathbf{x}_t, \mathbf{y}_t), (\mathbf{x}_r, \mathbf{y}_r)) = \phi(\mathbf{x}_t, \mathbf{y}_t) \cdot \phi(\mathbf{x}_r, \mathbf{y}_r)$$

# Kernel Trick: General Overview

- ▶ Define a kernel, and do not explicitly use dot product between vectors, only kernel calculations
- ▶ In some high-dimensional space, this corresponds to dot product
- ▶ In that space, the decision boundary is linear, but in the original space, we now have a non-linear decision boundary
- ▶ Let's do it for the Perceptron!

# Kernel Trick – Perceptron Algorithm

Training data:  $\mathcal{T} = \{(x_t, y_t)\}_{t=1}^{|\mathcal{T}|}$

1.  $\omega^{(0)} = 0$ ;  $i = 0$
2. for  $n : 1..N$
3.   for  $t : 1..T$
4.     Let  $y = \arg \max_y \omega^{(i)} \cdot \phi(x_t, y)$
5.     if  $y \neq y_t$
6.        $\omega^{(i+1)} = \omega^{(i)} + \phi(x_t, y_t) - \phi(x_t, y)$
7.        $i = i + 1$
8. return  $\omega^i$

- ▶ Each feature function  $\phi(x_t, y_t)$  is added and  $\phi(x_t, y)$  is subtracted to  $\omega$  say  $\alpha_{y,t}$  times
  - ▶  $\alpha_{y,t}$  is the # of times during learning label  $y$  is predicted for example  $t$
- ▶ Thus,

$$\omega = \sum_{t,y} \alpha_{y,t} [\phi(x_t, y_t) - \phi(x_t, y)]$$

# Kernel Trick – Perceptron Algorithm

- ▶ We can re-write the argmax function as:

$$\mathbf{y}^* = \arg \max_{\mathbf{y}^*} \omega^{(i)} \cdot \phi(\mathbf{x}, \mathbf{y}^*)$$

$$=$$

$$=$$

$$=$$

- ▶ We can then re-write the perceptron algorithm strictly with kernels



# Kernel Trick – Perceptron Algorithm (for handout)

- ▶ We can re-write the argmax function as:

$$\begin{aligned}
 \mathbf{y}^* &= \arg \max_{\mathbf{y}^*} \omega^{(i)} \cdot \phi(\mathbf{x}, \mathbf{y}^*) \\
 &= \arg \max_{\mathbf{y}^*} \sum_{t, \mathbf{y}} \alpha_{\mathbf{y}, t} [\phi(\mathbf{x}_t, \mathbf{y}_t) - \phi(\mathbf{x}_t, \mathbf{y})] \cdot \phi(\mathbf{x}, \mathbf{y}^*) \\
 &= \arg \max_{\mathbf{y}^*} \sum_{t, \mathbf{y}} \alpha_{\mathbf{y}, t} [\phi(\mathbf{x}_t, \mathbf{y}_t) \cdot \phi(\mathbf{x}_t, \mathbf{y}^*) - \phi(\mathbf{x}_t, \mathbf{y}) \cdot \phi(\mathbf{x}, \mathbf{y}^*)] \\
 &= \arg \max_{\mathbf{y}^*} \sum_{t, \mathbf{y}} \alpha_{\mathbf{y}, t} [\textcolor{red}{K}((\mathbf{x}_t, \mathbf{y}_t), (\mathbf{x}_t, \mathbf{y}^*)) - \textcolor{red}{K}((\mathbf{x}_t, \mathbf{y}), (\mathbf{x}, \mathbf{y}^*))]
 \end{aligned}$$

- ▶ We can then re-write the perceptron algorithm strictly with kernels

# Kernel Trick – Perceptron Algorithm

Training data:  $\mathcal{T} = \{(x_t, y_t)\}_{t=1}^{|\mathcal{T}|}$

1.  $\forall y, t$  set  $\alpha_{y,t} = 0$
2. for  $n : 1..N$
3.   for  $t : 1..T$
4.     Let  $y^* = \arg \max_{y^*} \sum_{t,y} \alpha_{y,t} [K((x_t, y_t), (x_t, y^*)) - K((x_t, y), (x_t, y^*))]$
5.     if  $y^* \neq y_t$
6.        $\alpha_{y^*,t} = \alpha_{y^*,t} + 1$

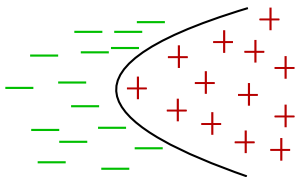
- Given a new instance  $x$

$$y^* = \arg \max_{y^*} \sum_{t,y} \alpha_{y,t} [K((x_t, y_t), (x, y^*)) - K((x_t, y), (x, y^*))]$$

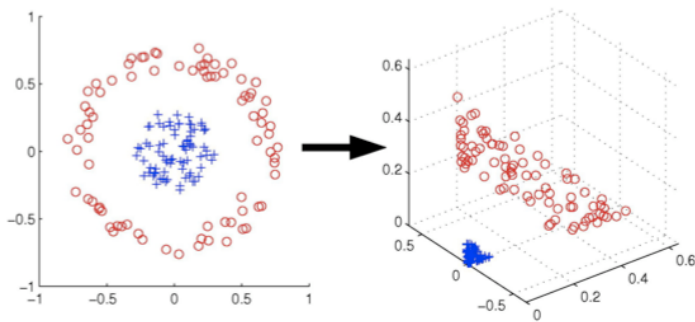
- But it seems like we have just complicated things???

# Kernels = Tractable Non-Linearity

- ▶ A linear classifier in a higher dimensional feature space is a non-linear classifier in the original space
- ▶ Computing a non-linear kernel is often better computationally than calculating the corresponding dot product in the high dimension feature space
- ▶ Thus, kernels allow us to efficiently learn non-linear classifiers



# Linear Classifiers in High Dimension



$$\mathbb{R}^2 \longrightarrow \mathbb{R}^3$$

$$(x_1, x_2) \longmapsto (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

## Example: Polynomial Kernel

- ▶  $\phi(x) \in \mathbb{R}^M$ ,  $d \geq 2$
- ▶  $K(x_t, x_s) = (\phi(x_t) \cdot \phi(x_s) + 1)^d$ 
  - ▶  $O(M)$  to calculate for any  $d!!$
- ▶ But in the original feature space (primal space)
  - ▶ Consider  $d = 2$ ,  $M = 2$ , and  $\phi(x_t) = [x_{t,1}, x_{t,2}]$

$$\begin{aligned}
 (\phi(x_t) \cdot \phi(x_s) + 1)^2 &= ([x_{t,1}, x_{t,2}] \cdot [x_{s,1}, x_{s,2}] + 1)^2 \\
 &= (x_{t,1}x_{s,1} + x_{t,2}x_{s,2} + 1)^2 \\
 &= (x_{t,1}x_{s,1})^2 + (x_{t,2}x_{s,2})^2 + 2(x_{t,1}x_{s,1}) + 2(x_{t,2}x_{s,2}) \\
 &\quad + 2(x_{t,1}x_{t,2}x_{s,1}x_{s,2}) + (1)^2
 \end{aligned}$$

which equals:

$$\underbrace{[(x_{t,1})^2, (x_{t,2})^2, \sqrt{2}x_{t,1}, \sqrt{2}x_{t,2}, \sqrt{2}x_{t,1}x_{t,2}, 1]}_{\text{feature vector in high-dimensional space}} \cdot \underbrace{[(x_{s,1})^2, (x_{s,2})^2, \sqrt{2}x_{s,1}, \sqrt{2}x_{s,2}, \sqrt{2}x_{s,1}x_{s,2}, 1]}_{\text{feature vector in high-dimensional space}}$$

# Popular Kernels

- ▶ Polynomial kernel

$$K(\mathbf{x}_t, \mathbf{x}_s) = (\phi(\mathbf{x}_t) \cdot \phi(\mathbf{x}_s) + 1)^d$$

- ▶ Gaussian radial basis kernel (infinite feature space representation!)

$$K(\mathbf{x}_t, \mathbf{x}_s) = \exp\left(\frac{-\|\phi(\mathbf{x}_t) - \phi(\mathbf{x}_s)\|^2}{2\sigma}\right)$$

- ▶ String kernels [Lodhi et al. 2002, Collins and Duffy 2002]
- ▶ Tree kernels [Collins and Duffy 2002]

# Kernels Summary

- ▶ Can turn a linear classifier into a non-linear classifier
- ▶ Kernels project feature space to higher dimensions
  - ▶ Sometimes exponentially larger
  - ▶ Sometimes an infinite space!
- ▶ Can “kernelize” algorithms to make them non-linear
- ▶ (e.g. support vector machines)

Wrap up and time for questions



# Summary

Basic principles of machine learning:

- ▶ To do learning, we set up an objective function that tells the fit of the model to the data
- ▶ We optimize with respect to the model (weights, probability model, etc.)
- ▶ Can do it in a batch or online fashion

What model to use?

- ▶ One example of a model: linear classifiers
- ▶ Can kernelize these models to get non-linear classification

## References and Further Reading

- ▶ A. L. Berger, S. A. Della Pietra, and V. J. Della Pietra. 1996.  
A maximum entropy approach to natural language processing. Computational Linguistics, 22(1).
- ▶ C.T. Chu, S.K. Kim, Y.A. Lin, Y.Y. Yu, G. Bradski, A.Y. Ng, and K. Olukotun. 2007.  
Map-Reduce for machine learning on multicore. In Advances in Neural Information Processing Systems.
- ▶ M. Collins and N. Duffy. 2002.  
New ranking algorithms for parsing and tagging: Kernels over discrete structures, and the voted perceptron. In Proc. ACL.
- ▶ M. Collins. 2002.  
Discriminative training methods for hidden Markov models: Theory and experiments with perceptron algorithms. In Proc. EMNLP.
- ▶ K. Crammer and Y. Singer. 2001.  
On the algorithmic implementation of multiclass kernel based vector machines. JMLR.
- ▶ K. Crammer and Y. Singer. 2003.  
Ultraconservative online algorithms for multiclass problems. JMLR.

- ▶ K. Crammer, O. Dekel, S. Shalev-Shwartz, and Y. Singer. 2003.  
Online passive aggressive algorithms. In Proc. NIPS.
- ▶ K. Crammer, O. Dekel, J. Keshat, S. Shalev-Shwartz, and Y. Singer. 2006.  
Online passive aggressive algorithms. JMLR.
- ▶ Y. Freund and R.E. Schapire. 1999.  
Large margin classification using the perceptron algorithm. Machine Learning, 37(3):277–296.
- ▶ T. Joachims. 2002.  
Learning to Classify Text using Support Vector Machines. Kluwer.
- ▶ J. Lafferty, A. McCallum, and F. Pereira. 2001.  
Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In Proc. ICML.
- ▶ H. Lodhi, C. Saunders, J. Shawe-Taylor, and N. Cristianini. 2002.  
Classification with string kernels. Journal of Machine Learning Research.
- ▶ G. Mann, R. McDonald, M. Mohri, N. Silberman, and D. Walker. 2009.  
Efficient large-scale distributed training of conditional maximum entropy models. In Advances in Neural Information Processing Systems.
- ▶ A. McCallum, D. Freitag, and F. Pereira. 2000.

Maximum entropy Markov models for information extraction and segmentation. In Proc. ICML.

- ▶ R. McDonald, K. Crammer, and F. Pereira. 2005.  
Online large-margin training of dependency parsers. In Proc. ACL.
- ▶ K.R. Müller, S. Mika, G. Rätsch, K. Tsuda, and B. Schölkopf. 2001.  
An introduction to kernel-based learning algorithms. IEEE Neural Networks, 12(2):181–201.
- ▶ J Nocedal and SJ Wright. 1999.  
Numerical optimization, volume 2. Springer New York.
- ▶ F. Sha and F. Pereira. 2003.  
Shallow parsing with conditional random fields. In Proc. HLT/NAACL, pages 213–220.
- ▶ C. Sutton and A. McCallum. 2006.  
An introduction to conditional random fields for relational learning. In L. Getoor and B. Taskar, editors, Introduction to Statistical Relational Learning. MIT Press.
- ▶ B. Taskar, C. Guestrin, and D. Koller. 2003.  
Max-margin Markov networks. In Proc. NIPS.
- ▶ B. Taskar. 2004.

Learning Structured Prediction Models: A Large Margin Approach. Ph.D. thesis, Stanford.

- ▶ I. Tsochantaridis, T. Hofmann, T. Joachims, and Y. Altun. 2004. Support vector learning for interdependent and structured output spaces. In Proc. ICML.
- ▶ T. Zhang. 2004. Solving large scale linear prediction problems using stochastic gradient descent algorithms. In Proceedings of the twenty-first international conference on Machine learning.