# Introduction to Machine Learning: Linear Learners<sup>1</sup>

Lisbon Machine Learning School, 2016

Stefan Riezler

Computational Linguistics & Scientific Computing Heidelberg University, Germany riezler@cl.uni-heidelberg.de

<sup>1</sup>Including (lots of) material from Shay Cohen and Ryan McDonald

# Modeling the Frog's Perceptual System



# Modeling the Frog's Perceptual System

- [Lettvin et al. 1959] show that the frog's perceptual system constructs reality by four separate operations:
  - contrast detection: presence of sharp boundary?
  - convexity detection: how curved and how big is object?
  - movement detection: is object moving?
  - dimming speed: how fast does object obstruct light?
- The frog's goal: Capture any object of the size of an insect or worm providing it moves like one.

# Modeling the Frog's Perceptual System

- [Lettvin et al. 1959] show that the frog's perceptual system constructs reality by four separate operations:
  - contrast detection: presence of sharp boundary?
  - convexity detection: how curved and how big is object?
  - movement detection: is object moving?
  - dimming speed: how fast does object obstruct light?
- The frog's goal: Capture any object of the size of an insect or worm providing it moves like one.
- Can we build a model of this perceptual system and learn to capture the right objects?

► Assume training data of edible (+) and inedible (-) objects

| convex | speed  | label | convex | speed  | label |
|--------|--------|-------|--------|--------|-------|
| small  | small  | -     | small  | large  | +     |
| small  | medium | -     | medium | large  | +     |
| small  | medium | -     | medium | large  | +     |
| medium | small  | -     | large  | small  | +     |
| large  | small  | -     | large  | large  | +     |
| small  | small  | -     | large  | medium | +     |
| small  | large  | -     |        |        |       |
| small  | medium | -     |        |        |       |

► Assume training data of edible (+) and inedible (-) objects

| convex | speed  | label | convex | speed  | label |
|--------|--------|-------|--------|--------|-------|
| small  | small  | -     | small  | large  | +     |
| small  | medium | -     | medium | large  | +     |
| small  | medium | -     | medium | large  | +     |
| medium | small  | -     | large  | small  | +     |
| large  | small  | -     | large  | large  | +     |
| small  | small  | -     | large  | medium | +     |
| small  | large  | -     |        |        |       |
| small  | medium | -     |        |        |       |

### Learning model parameters from data:

► Assume training data of edible (+) and inedible (-) objects

| convex | speed  | label | convex | speed  | label |
|--------|--------|-------|--------|--------|-------|
| small  | small  | -     | small  | large  | +     |
| small  | medium | -     | medium | large  | +     |
| small  | medium | -     | medium | large  | +     |
| medium | small  | -     | large  | small  | +     |
| large  | small  | -     | large  | large  | +     |
| small  | small  | -     | large  | medium | +     |
| small  | large  | -     |        |        |       |
| small  | medium | -     |        |        |       |

### Learning model parameters from data:

$$\begin{array}{l} \blacktriangleright & p(+)=8/14, \ p(-)=6/14 \\ \hline & p(convex = small|-) = & , \ p(convex = med|-) = & , \ p(convex = large|-) = \\ p(speed = small|-) = & , \ p(speed = med|-) = & , \ p(speed = large|-) = \\ p(speed = small|+) = & , \ p(convex = med|+) = & , \ p(speed = large|+) = \\ p(speed = small|+) = & , \ p(speed = med|+) = & , \ p(speed = large|+) = \\ \end{array}$$

Assume training data of edible (+) and inedible (-) objects

| convex | speed  | label | convex | speed  | label |
|--------|--------|-------|--------|--------|-------|
| small  | small  | -     | small  | large  | +     |
| small  | medium | -     | medium | large  | +     |
| small  | medium | -     | medium | large  | +     |
| medium | small  | -     | large  | small  | +     |
| large  | small  | -     | large  | large  | +     |
| small  | small  | -     | large  | medium | +     |
| small  | large  | -     |        |        |       |
| small  | medium | -     |        |        |       |

#### Learning model parameters from data:

$$p(+) = 8/14$$
,  $p(-) = 6/14$ 

p(convex = small|-) = 6/8, p(convex = med|-) = 1/8, p(convex = large|-) = 1/8 p(speed = small|-) = 4/8, p(speed = med|-) = 3/8, p(speed = large|-) = 1/8 p(convex = small|+) = 1/6, p(convex = med|+) = 2/6, p(convex = large|+) = 3/6 p(speed = small|+) = 1/6, p(speed = med|+) = 1/6, p(speed = large|+) = 4/6

Assume training data of edible (+) and inedible (-) objects

| convex | speed  | label | convex | speed  | label |
|--------|--------|-------|--------|--------|-------|
| small  | small  | -     | small  | large  | +     |
| small  | medium | -     | medium | large  | +     |
| small  | medium | -     | medium | large  | +     |
| medium | small  | -     | large  | small  | +     |
| large  | small  | -     | large  | large  | +     |
| small  | small  | -     | large  | medium | +     |
| small  | large  | -     |        |        |       |
| small  | medium | -     |        |        |       |

#### Learning model parameters from data:

p(convex = small|-) = 6/8, p(convex = med|-) = 1/8, p(convex = large|-) = 1/8 p(speed = small|-) = 4/8, p(speed = med|-) = 3/8, p(speed = large|-) = 1/8 p(convex = small|+) = 1/6, p(convex = med|+) = 2/6, p(convex = large|+) = 3/6 p(speed = small|+) = 1/6, p(speed = med|+) = 1/6, p(speed = large|+) = 4/6

### Predict unseen p(label = ?, convex = med, speed = med)

p(-) · p(convex = med|-) · p(speed = med|-) =
p(+) · p(convex = med|+) · p(speed = med|+) =

Assume training data of edible (+) and inedible (-) objects

| convex | speed  | label | convex | speed  | label |
|--------|--------|-------|--------|--------|-------|
| small  | small  | -     | small  | large  | +     |
| small  | medium | -     | medium | large  | +     |
| small  | medium | -     | medium | large  | +     |
| medium | small  | -     | large  | small  | +     |
| large  | small  | -     | large  | large  | +     |
| small  | small  | -     | large  | medium | +     |
| small  | large  | -     |        |        |       |
| small  | medium | -     |        |        |       |

#### Learning model parameters from data:

p(convex = small|-) = 6/8, p(convex = med|-) = 1/8, p(convex = large|-) = 1/8 p(speed = small|-) = 4/8, p(speed = med|-) = 3/8, p(speed = large|-) = 1/8 p(convex = small|+) = 1/6, p(convex = med|+) = 2/6, p(convex = large|+) = 3/6 p(speed = small|+) = 1/6, p(speed = med|+) = 1/6, p(speed = large|+) = 4/6

### Predict unseen p(label = ?, convex = med, speed = med)

$$p(-) \cdot p(\text{convex} = \text{med}|-) \cdot p(\text{speed} = \text{med}|-) = 8/14 \cdot 1/8 \cdot 3/8 = 0.027$$

▶ 
$$p(+) \cdot p(\text{convex} = \text{med}|+) \cdot p(\text{speed} = \text{med}|+) =$$

Assume training data of edible (+) and inedible (-) objects

| convex | speed  | label | convex | speed  | label |
|--------|--------|-------|--------|--------|-------|
| small  | small  | -     | small  | large  | +     |
| small  | medium | -     | medium | large  | +     |
| small  | medium | -     | medium | large  | +     |
| medium | small  | -     | large  | small  | +     |
| large  | small  | -     | large  | large  | +     |
| small  | small  | -     | large  | medium | +     |
| small  | large  | -     |        |        |       |
| small  | medium | -     |        |        |       |

#### Learning model parameters from data:

p(convex = small|-) = 6/8, p(convex = med|-) = 1/8, p(convex = large|-) = 1/8 p(speed = small|-) = 4/8, p(speed = med|-) = 3/8, p(speed = large|-) = 1/8 p(convex = small|+) = 1/6, p(convex = med|+) = 2/6, p(convex = large|+) = 3/6 p(speed = small|+) = 1/6, p(speed = med|+) = 1/6, p(speed = large|+) = 4/6

### Predict unseen p(label = ?, convex = med, speed = med)

$$\begin{array}{l} \blacktriangleright \quad p(-) \cdot p(convex = med|-) \cdot p(speed = med|-) = 8/14 \cdot 1/8 \cdot 3/8 = 0.027 \\ \hline \quad p(+) \cdot p(convex = med|+) \cdot p(speed = med|+) = 6/14 \cdot 2/6 \cdot 1/6 = 0.02 \\ \end{array}$$

Assume training data of edible (+) and inedible (-) objects

| convex | speed  | label | convex | speed  | label |
|--------|--------|-------|--------|--------|-------|
| small  | small  | -     | small  | large  | +     |
| small  | medium | -     | medium | large  | +     |
| small  | medium | -     | medium | large  | +     |
| medium | small  | -     | large  | small  | +     |
| large  | small  | -     | large  | large  | +     |
| small  | small  | -     | large  | medium | +     |
| small  | large  | -     |        |        |       |
| small  | medium | -     |        |        |       |

#### Learning model parameters from data:

▶ p(+) = 8/14, p(-) = 6/14

 $\begin{array}{l} & \mathsf{p}(\mathsf{convex} = \mathsf{small}|{-}) = 6/8, \ \mathsf{p}(\mathsf{convex} = \mathsf{med}|{-}) = 1/8, \ \mathsf{p}(\mathsf{convex} = \mathsf{large}|{-}) = 1/8 \\ & \mathsf{p}(\mathsf{speed} = \mathsf{smal}|{-}) = 4/8, \ \mathsf{p}(\mathsf{speed} = \mathsf{med}|{-}) = 3/8, \ \mathsf{p}(\mathsf{speed} = \mathsf{large}|{-}) = 1/8 \\ & \mathsf{p}(\mathsf{convex} = \mathsf{smal}||{+}) = 1/6, \ \mathsf{p}(\mathsf{sorvex} = \mathsf{med}|{+}) = 2/6, \ \mathsf{p}(\mathsf{convex} = \mathsf{large}|{+}) = 3/6 \\ & \mathsf{p}(\mathsf{speed} = \mathsf{smal}||{+}) = 1/6, \ \mathsf{p}(\mathsf{speed} = \mathsf{med}|{+}) = 1/6, \ \mathsf{p}(\mathsf{speed} = \mathsf{large}|{+}) = 4/6 \\ \end{array}$ 

### Predict unseen p(label = ?, convex = med, speed = med)

- p(-) · p(convex = med|-) · p(speed = med|-) = 8/14 · 1/8 · 3/8 = 0.027
- ▶  $p(+) \cdot p(convex = med|+) \cdot p(speed = med|+) = 6/14 \cdot 2/6 \cdot 1/6 = 0.024$
- Inedible: p(convex = med, speed = med, label = -) > p(convex = med, speed = med, label = +)!

# Machine Learning is a Frog's World

- Machine learning problems can be seen as problems of function estimation where
  - our models are based on a combined feature representation of inputs and outputs
    - similar to the frog whose world is constructed by four-dimensional feature vector based on detection operations

# Machine Learning is a Frog's World

- Machine learning problems can be seen as problems of function estimation where
  - our models are based on a combined feature representation of inputs and outputs
    - similar to the frog whose world is constructed by four-dimensional feature vector based on detection operations
  - learning of parameter weights is done by optimizing fit of model to training data
    - frog uses binary classification into edible/inedible objects as supervision signals for learning

# Machine Learning is a Frog's World

- Machine learning problems can be seen as problems of function estimation where
  - our models are based on a combined feature representation of inputs and outputs
    - similar to the frog whose world is constructed by four-dimensional feature vector based on detection operations
  - learning of parameter weights is done by optimizing fit of model to training data
    - frog uses binary classification into edible/inedible objects as supervision signals for learning
  - The model used in the frog's perception example is called Naive Bayes: It measures compatibility of inputs to outputs by a linear model and optimizes parameters by linear optimization

# Lecture Outline: Linear Learners for NLP

- Preliminaries
  - Data: input/output, assumptions
  - Feature representations
  - Linear models
- Linear learners
  - Naive Bayes
  - Generative versus discriminative
  - Logistic Regression
  - Perceptron
  - Large-Margin Learners (SVMs)
- Regularization
- Online learning
- Non-linear models

### Inputs and Outputs

▶ Input:  $x \in X$ 

▶ e.g., document or sentence with some words x = w<sub>1</sub>...w<sub>n</sub>
 ▶ Output: y ∈ Y

e.g., document class, translation, parse tree

- ▶ Input/Output pair:  $({m x},{m y})\in {\mathcal X} imes {\mathcal Y}$ 
  - e.g., a document x and its class label y,
  - a source sentence x and its translation y,
  - $\blacktriangleright$  a sentence x and its parse tree y

### **Feature Representations**

- We assume a mapping from input x to a high dimensional feature vector
  - $\phi(x):\mathcal{X} o \mathbb{R}^m$
- $\blacktriangleright$  For many cases, more convenient to have mapping from input-output pairs (x,y)

•  $\phi({m x},{m y}): {\mathcal X} imes {\mathcal Y} o {\mathbb R}^m$ 

- Under certain assumptions, these are equivalent
- Most papers in NLP use  $\phi(x,y)$

### **Feature Representations**

We assume a mapping from input x to a high dimensional feature vector

•  $\phi(x):\mathcal{X} o \mathbb{R}^m$ 

 $\blacktriangleright$  For many cases, more convenient to have mapping from input-output pairs (x,y)

•  $\phi(x,y):\mathcal{X} imes\mathcal{Y} o\mathbb{R}^m$ 

- 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)

### **Feature Representations**

We assume a mapping from input x to a high dimensional feature vector

•  $\phi(x):\mathcal{X} o \mathbb{R}^m$ 

 $\blacktriangleright$  For many cases, more convenient to have mapping from input-output pairs (x,y)

•  $\phi(x,y):\mathcal{X} imes\mathcal{Y} o\mathbb{R}^m$ 

- 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  $j^{th}$  value

#### x is a document and y is a label

$$\phi_j(oldsymbol{x},oldsymbol{y}) = \left\{egin{array}{ccc} 1 & ext{if} \ oldsymbol{x} \ ext{contains the word "interest"} \ & ext{and} \ oldsymbol{y} = ext{"financial"} \ & ext{0} \ & ext{otherwise} \end{array}
ight.$$

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

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

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

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

$$\phi_j(oldsymbol{x},oldsymbol{y}) = \left\{egin{array}{ccc} 1 & ext{if} oldsymbol{x} = ext{``bank'' and} oldsymbol{y} = ext{Verb} \ 0 & ext{otherwise} \end{array}
ight.$$

What weight would it get?

• x is a source sentence and y is translation

$$\phi_j(x,y) = \left\{egin{array}{ll} 1 & ext{if "y a-t-il" present in $x$} \ & ext{ and "are there" present in $y$} \ 0 & ext{otherwise} \end{array}
ight.$$

$$\phi_k(x,y) = \left\{egin{array}{cccc} 1 & ext{if "y a-t-il" present in $x$} \ & ext{ and "are there any" present in $y$} \ 0 & ext{otherwise} \end{array}
ight.$$

Which phrase indicator should be preferred?



#### Note: Label y includes sentence x

# Linear Models

Linear model: Defines a discriminant function that is based on linear combination of features and weights

$$egin{array}{rcl} f(m{x};m{\omega}) &=& rg\max_{m{y}\in\mathcal{Y}} \ m{\omega}\cdotm{\phi}(m{x},m{y}) \ &=& rg\max_{m{y}\in\mathcal{Y}} \ \sum_{j=0}^m m{\omega}_j imesm{\phi}_j(m{x},m{y}) \end{array}$$

# Linear Models

Linear model: Defines a discriminant function that is based on linear combination of features and weights

$$egin{array}{rcl} f(m{x};m{\omega}) &=& rgmax \ m{y}\in\mathcal{Y} &m{\omega}\cdotm{\phi}(m{x},m{y}) \ &=& rgmax \ m{y}\in\mathcal{Y} & \sum_{j=0}^mm{\omega}_j imesm{\phi}_j(m{x},m{y}) \end{array}$$

- Let  $\boldsymbol{\omega} \in \mathbb{R}^m$  be a high dimensional weight vector
- Assume that *\u03c6* is known
  - Multiclass Classification:  $\mathcal{Y} = \{0, 1, \dots, N\}$

$$oldsymbol{y} = rgmax_{oldsymbol{y}\in\mathcal{Y}} oldsymbol{\omega}\cdot \phi(oldsymbol{x},oldsymbol{y})$$

Binary Classification just a special case of multiclass

# **Binary Linear Model**

 $\omega$  defines a hyperplane (line in 2 dimensions) that divides all points:



# Multiclass Linear Model

Defines regions of space. Visualization difficult.



▶ + are all points (x, y) where + =  $rg \max_{u} \omega \cdot \phi(x, y)$ 

# Separability

A set of points is separable, if there exists a ω such that classification is perfect





Not Separable

# **Supervised Learning**

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

- Input:
  - ▶ i.i.d. (independent and identically distributed) training examples  $\mathcal{T} = \{(x_t, y_t)\}_{t=1}^{|\mathcal{T}|}$
  - feature representation  $\phi$

# Supervised Learning

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

- Input:
  - ▶ i.i.d. (independent and identically distributed) training examples  $\mathcal{T} = \{(x_t, y_t)\}_{t=1}^{|\mathcal{T}|}$
  - $\blacktriangleright$  feature representation  $\phi$
- Output: ω that maximizes an objective function on the training set
  - $\boldsymbol{\omega} = \arg \max \mathcal{L}(\mathcal{T}; \boldsymbol{\omega})$
  - Equivalently minimize:  $\boldsymbol{\omega} = \arg\min \mathcal{L}(\mathcal{T}; \boldsymbol{\omega})$

### **Objective Functions**

• Ideally we can decompose  $\mathcal{L}$  by training pairs (x, y)

▶  $\mathcal{L}(\mathcal{T}; \omega) \propto \sum_{(x,y) \in \mathcal{T}} loss((x, y); \omega)$ ▶ *loss* is a function that measures some value correlated with errors of parameters  $\omega$  on instance (x,y)

### **Objective Functions**

Ideally we can decompose L by training pairs (x,y)

- $\blacktriangleright \ \mathcal{L}(\mathcal{T}; oldsymbol{\omega}) \propto \sum_{(oldsymbol{x}, oldsymbol{y}) \in \mathcal{T}} \mathit{loss}((oldsymbol{x}, oldsymbol{y}); oldsymbol{\omega})$
- ▶ loss is a function that measures some value correlated with errors of parameters  $\omega$  on instance (x, y)
- Example:
  - ▶  $y \in \{1, -1\}, f(x; \omega)$  is the prediction we make for x using  $\omega$ ▶ 0-1 loss function:  $loss((x, y); \omega) = \begin{cases} 0 & \text{if } f(x; \omega) = y, \\ 1 & \text{else} \end{cases}$

# **Convex Optimization**



A function is convex if its graph lies on or below the line segment connecting any two points on the graph

 $f(\alpha x + \beta y) \le \alpha f(x) + \beta f(y)$  for all  $\alpha, \beta \ge 0, \ \alpha + \beta = 1$  (1)

▶ For linear models we have equality in (1)

### **Convex Optimization**

- Optimization problem is defined as problem of finding a point that minimizes our objective function (maximization is minimization of -f(x))
- Limit attention to convex (or even linear) functions
- Find point at which gradient of f is 0
  - Slope of the gradient is non-decreasing as one moves away from minimum, 0 at minimum
  - In order to find minimum, follow opposite direction of gradient
  - For convex functions, this will lead to the single global minimum
# Naive Bayes

### Naive Bayes

Probabilistic decision model:

$$rg \max_{oldsymbol{y}} oldsymbol{P}(oldsymbol{y}|oldsymbol{x}) \propto rg \max_{oldsymbol{y}} oldsymbol{P}(oldsymbol{y}) oldsymbol{P}(oldsymbol{x}|oldsymbol{y})$$

► Uses Bayes Rule:

$$egin{aligned} & P(oldsymbol{y}|oldsymbol{x}) = rac{P(oldsymbol{y})P(oldsymbol{x}|oldsymbol{y})}{P(oldsymbol{x})} ext{ for fixed }oldsymbol{x} \end{aligned}$$

- Generative model since P(y)P(x|y) = P(x, y) is a joint probability
  - Because we model a distribution that can randomly generate outputs and inputs, not just outputs

### Naivety of Naive Bayes

• We need to decide on the structure of P(x, y)

$$\blacktriangleright P(x|y) = P(\phi(x)|y) = P(\phi_1(x), \dots, \phi_m(x)|y)$$

Naive Bayes Assumption (conditional independence) $P(m{\phi}_1(m{x}),\ldots,m{\phi}_m(m{x})|m{y})=\prod_i P(m{\phi}_i(m{x})|m{y})$ 

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

#### Naive Bayes – Learning

- ▶ Input:  $\mathcal{T} = \{(x_t, y_t)\}_{t=1}^{|\mathcal{T}|}$
- Let  $\phi_i(x) \in \{1, \ldots, F_i\}$
- Parameters  $\mathcal{P} = \{P(\boldsymbol{y}), P(\phi_i(\boldsymbol{x})|\boldsymbol{y})\}$ 
  - Both P(y) and  $P(\phi_i(x)|y)$  are multinomials

## **Maximum Likelihood Estimation**

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

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

## Naive Bayes – Learning

MLE has closed form solution

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

$$egin{aligned} P(oldsymbol{y}) &= rac{\sum_{t=1}^{|\mathcal{T}|} \llbracket oldsymbol{y}_t = oldsymbol{y} 
rbrace }{|\mathcal{T}|} \ P(\phi_i(oldsymbol{x}) | oldsymbol{y}) &= rac{\sum_{t=1}^{|\mathcal{T}|} \llbracket \phi_i(oldsymbol{x}_t) = \phi_i(oldsymbol{x}) ext{ and } oldsymbol{y}_t = oldsymbol{y} 
rbrace \ \sum_{t=1}^{|\mathcal{T}|} \llbracket oldsymbol{y}_t = oldsymbol{y} 
rbrace \ \sum_{t=1}^{|\mathcal{T}|} \llbracket oldsymbol{y}_t = oldsymbol{y} 
rbrace \end{aligned}$$

where  $\llbracket p \rrbracket = \begin{cases} 1 & \text{if } p \text{ is true,} \\ 0 & \text{otherwise.} \end{cases}$ Thus, these are just normalized counts over events in  $\mathcal{T}$ 

### Naive Bayes Document Classification

- doc 1:  $y_1 =$  sports, "hockey is fast"
- doc 2: y<sub>2</sub> = politics, "politicians talk fast"
- doc 3: y<sub>3</sub> = politics, "washington is sleazy"
- $\phi_0(x) = 1$  if doc has word 'hockey', 0 else
- $\phi_1(x) = 1$  if doc has word 'is', 0 else
- $\phi_2(x) = 1$  if doc has word 'fast', 0 else
- $\phi_3(x) = 1$  if doc has word 'politicians', 0 else
- $\phi_4(x) = 1$  if doc has word 'talk', 0 else
- $\phi_5(x) = 1$  if doc has word 'washington', 0 else
- $\phi_6(x) = 1$  if doc has word 'sleazy', 0 else

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

### Naive Bayes Document Classification

- doc 1:  $y_1 =$  sports, "hockey is fast"
- doc 2: y<sub>2</sub> = politics, "politicians talk fast"
- doc 3: y<sub>3</sub> = politics, "washington is sleazy"
- $\phi_0(x) = 1$  if doc has word 'hockey', 0 else
- $\phi_1(x) = 1$  if doc has word 'is', 0 else
- $\phi_2(x) = 1$  if doc has word 'fast', 0 else
- $\phi_3(x) = 1$  if doc has word 'politicians', 0 else
- $\phi_4(x) = 1$  if doc has word 'talk', 0 else
- $\phi_5(x) = 1$  if doc has word 'washington', 0 else
- $\phi_6(x) = 1$  if doc has word 'sleazy', 0 else

Your turn? What is P(sports)? 1/3 What is  $P(\phi_0(x) = 1 | \text{politics})$ ?

### Naive Bayes Document Classification

- doc 1:  $y_1 =$  sports, "hockey is fast"
- doc 2: y<sub>2</sub> = politics, "politicians talk fast"
- doc 3: y<sub>3</sub> = politics, "washington is sleazy"
- $\phi_0(x) = 1$  if doc has word 'hockey', 0 else
- $\phi_1(x) = 1$  if doc has word 'is', 0 else
- $\phi_2(x) = 1$  if doc has word 'fast', 0 else
- $\phi_3(x) = 1$  if doc has word 'politicians', 0 else
- $\phi_4(x) = 1$  if doc has word 'talk', 0 else
- $\phi_5(x) = 1$  if doc has word 'washington', 0 else
- $\phi_6(x) = 1$  if doc has word 'sleazy', 0 else

Your turn? What is P(sports)? 1/3 What is  $P(\phi_0(x) = 1 | \text{politics})$ ? 0

Naive Bayes

# **Deriving MLE**

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

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

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

$$\mathcal{P} = \operatorname*{arg\,max}_{P(\boldsymbol{y})} \sum_{t=1}^{|\mathcal{T}|} \log P(\boldsymbol{y}_t) + \operatorname*{arg\,max}_{P(\phi_i(\boldsymbol{x})|\boldsymbol{y})} \sum_{t=1}^{|\mathcal{T}|} \sum_{i=1}^{m} \log P(\phi_i(\boldsymbol{x}_t)|\boldsymbol{y}_t)$$

Both optimizations are of the form

 $\arg \max_{P} \sum_{v} \operatorname{count}(v) \log P(v)$ , s.t.,  $\sum_{v} P(v) = 1$ ,  $P(v) \ge 0$ 

where v is event in  $\mathcal{T}$ , e.g.  $(\boldsymbol{y}_t = \boldsymbol{y})$  or  $(\phi_i(\boldsymbol{x}_t) = \phi_i(\boldsymbol{x})$  and  $\boldsymbol{y}_t = \boldsymbol{y})$ 

$$\arg \max_{P} \sum_{v} \operatorname{count}(v) \log P(v) \\ \text{s.t., } \sum_{v} P(v) = 1, \ P(v) \ge 0$$

Introduce Lagrangian multiplier  $\lambda$ , optimization becomes

$$\operatorname{arg\,max}_{P,\lambda} \sum_{v} \operatorname{count}(v) \log P(v) - \lambda (\sum_{v} P(v) - 1)$$

Derivative:

Set to zero:

Final solution:

$$\arg \max_{P} \sum_{v} \operatorname{count}(v) \log P(v) \\ \text{s.t., } \sum_{v} P(v) = 1, \ P(v) \ge 0$$

Introduce Lagrangian multiplier  $\lambda$ , optimization becomes

$$\operatorname{arg\,max}_{P,\lambda} \sum_{v} \operatorname{count}(v) \log P(v) - \lambda (\sum_{v} P(v) - 1)$$

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

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

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

Reinstantiate events v in  $\mathcal{T}$ :

$$egin{aligned} & P(oldsymbol{y}) = rac{\sum_{t=1}^{|\mathcal{T}|} \llbracket oldsymbol{y}_t = oldsymbol{y} 
rbrace \ |\mathcal{T}|}{|\mathcal{T}|} \ & P(\phi_i(oldsymbol{x}) | oldsymbol{y}) = rac{\sum_{t=1}^{|\mathcal{T}|} \llbracket \phi_i(oldsymbol{x}_t) = \phi_i(oldsymbol{x}) ext{ and } oldsymbol{y}_t = oldsymbol{y} 
rbrace \ & \sum_{t=1}^{|\mathcal{T}|} \llbracket oldsymbol{y}_t = oldsymbol{y} 
rblace \ & \sum_{t=1}^{|\mathcal{T}|} \llbracket oldsymbol{y}_t = oldsymbol{y} 
rblace$$

### Naive Bayes is a linear model

▶ Let 
$$oldsymbol{\omega}_{oldsymbol{y}} = \log P(oldsymbol{y}), \, orall oldsymbol{y} \in \mathcal{Y}$$

$$\blacktriangleright \text{ Let } \omega_{\phi_i(\boldsymbol{x}),\boldsymbol{y}} = \log P(\phi_i(\boldsymbol{x})|\boldsymbol{y}), \ \forall \boldsymbol{y} \in \mathcal{Y}, \phi_i(\boldsymbol{x}) \in \{1,\ldots,F_i\}$$

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

where 
$$oldsymbol{\psi}_*\in\{0,1\}$$
,  $oldsymbol{\psi}_{i,j}(oldsymbol{x})=[\![oldsymbol{\phi}_i(oldsymbol{x})=j]\!]$ ,  $oldsymbol{\psi}_{oldsymbol{y}'}(oldsymbol{y})=[\![oldsymbol{y}=oldsymbol{y}']\!]$ 

### Naive Bayes is a linear model

$$\begin{aligned} \arg \max_{\boldsymbol{y}} \ P(\boldsymbol{y}|\boldsymbol{\phi}(\boldsymbol{x})) &\propto \ \arg \max_{\boldsymbol{y}} \ P(\boldsymbol{\phi}(\boldsymbol{x}), \boldsymbol{y}) = \arg \max_{\boldsymbol{y}} \ P(\boldsymbol{y}) \prod_{i=1}^{m} P(\boldsymbol{\phi}_{i}(\boldsymbol{x})|\boldsymbol{y}) \\ &= \ \arg \max_{\boldsymbol{y}} \ \log P(\boldsymbol{y}) + \sum_{i=1}^{m} \log P(\boldsymbol{\phi}_{i}(\boldsymbol{x})|\boldsymbol{y}) \\ &= \ \arg \max_{\boldsymbol{y}} \ \boldsymbol{\omega}_{\boldsymbol{y}} + \sum_{i=1}^{m} \boldsymbol{\omega}_{\boldsymbol{\phi}_{i}(\boldsymbol{x}), \boldsymbol{y}} \\ &= \ \arg \max_{\boldsymbol{y}} \ \sum_{\boldsymbol{y}'} \boldsymbol{\omega}_{\boldsymbol{y}} \boldsymbol{\psi}_{\boldsymbol{y}'}(\boldsymbol{y}) + \sum_{i=1}^{m} \sum_{j=1}^{F_{i}} \boldsymbol{\omega}_{\boldsymbol{\phi}_{i}(\boldsymbol{x}), \boldsymbol{y}} \boldsymbol{\psi}_{i, j}(\boldsymbol{x}) \\ & \end{aligned}$$
where  $\boldsymbol{\psi}_{*} \in \{0, 1\}, \ \boldsymbol{\psi}_{i, j}(\boldsymbol{x}) = \llbracket \boldsymbol{\phi}_{i}(\boldsymbol{x}) = j \rrbracket, \ \boldsymbol{\psi}_{\boldsymbol{y}'}(\boldsymbol{y}) = \llbracket \boldsymbol{y} = \boldsymbol{y'} \rrbracket$ 

## Smoothing

- doc 1:  $y_1 =$  sports, "hockey is fast"
- doc 2: y<sub>2</sub> = politics, "politicians talk fast"
- doc 3:  $y_3$  = 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

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

### **Discriminative versus Generative**

- Generative models attempt to model inputs and outputs
  - e.g., Naive Bayes = MLE of joint distribution P(x, y)
  - Statistical model must explain generation of input
- Occam's Razor: "Among competing hypotheses, the one with the fewest assumptions should be selected"
- Discriminative models
  - Use  $\mathcal{L}$  that directly optimizes P(y|x) (or something related)
  - Logistic Regression MLE of P(y|x)
  - Perceptron and SVMs minimize classification error
- Generative and discriminative models use P(y|x) for prediction
- $\blacktriangleright$  Differ only on what distribution they use to set  $\omega$

Define a conditional probability:

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

Note: still a linear model

$$\underset{y}{\operatorname{arg\,max}} P(y|x) = \operatorname{arg\,max}_{y} \frac{e^{\omega \cdot \phi(x,y)}}{Z_{x}}$$

$$= \operatorname{arg\,max}_{y} e^{\omega \cdot \phi(x,y)}$$

$$= \operatorname{arg\,max}_{y} \omega \cdot \phi(x,y)$$

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

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

$$\begin{split} \boldsymbol{\omega} &= \arg \max_{\boldsymbol{\omega}} \ \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) \\ &= \arg \max_{\boldsymbol{\omega}} \ \prod_{t=1}^{|\mathcal{T}|} P(\boldsymbol{y}_t | \boldsymbol{x}_t) = \arg \max_{\boldsymbol{\omega}} \ \sum_{t=1}^{|\mathcal{T}|} \log P(\boldsymbol{y}_t | \boldsymbol{x}_t) \end{split}$$

In a nutshell we set the weights ω so that we assign as much probability to the correct label y for each x in the training set

$$egin{aligned} P(m{y}|m{x}) &= rac{e^{\omega \cdot \phi(x,y)}}{Z_x}, & ext{ where } Z_x &= \sum_{m{y}' \in \mathcal{Y}} e^{\omega \cdot \phi(x,y')} \ & \omega &= rg\max_{m{\omega}} \ \sum_{t=1}^{|\mathcal{T}|} \log P(m{y}_t|m{x}_t) \ (*) \end{aligned}$$

- The objective function (\*) is concave
- 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)

# **Gradient Ascent**



### **Gradient Ascent**

▶ Let 
$$\mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) = \sum_{t=1}^{|\mathcal{T}|} \log \left( e^{\boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t)} / Z_{\boldsymbol{x}} \right)$$

- Want to find  $\operatorname{arg\,max}_{\boldsymbol{\omega}} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega})$ 
  - Set  $\omega^0 = O^m$
  - Iterate until convergence

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

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

#### Gradient Descent

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

• Want to find  $\arg \min_{\omega} \mathcal{L}(\mathcal{T}; \omega)$ 

• Set 
$$\omega^0 = O^m$$

Iterate until convergence

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

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

# The partial derivatives

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

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

# Partial derivatives - some reminders

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

$$rac{\partial}{\partial oldsymbol{\omega}_i} \mathcal{L}(\mathcal{T};oldsymbol{\omega}) =$$

# The partial derivatives (1)

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

# The partial derivatives

Now, 
$$\frac{\partial}{\partial \boldsymbol{\omega}_i} \frac{e^{\sum_j \boldsymbol{\omega}_j \times \boldsymbol{\phi}_j(\boldsymbol{x}_t, \boldsymbol{y}_t)}}{Z_{\boldsymbol{x}_t}} =$$

# The partial derivatives (2)

Now,

$$\frac{\partial}{\partial \omega_{i}} \frac{e^{\sum_{j} \omega_{j} \times \phi_{j}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t})}}{Z_{\boldsymbol{x}_{t}}} = \frac{Z_{\boldsymbol{x}_{t}} \frac{\partial}{\partial \omega_{i}} e^{\sum_{j} \omega_{j} \times \phi_{j}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t})} - e^{\sum_{j} \omega_{j} \times \phi_{j}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t})} \frac{\partial}{\partial \omega_{i}} Z_{\boldsymbol{x}_{t}}}{Z_{\boldsymbol{x}_{t}}^{2}}$$

$$= \frac{Z_{\boldsymbol{x}_{t}} e^{\sum_{j} \omega_{j} \times \phi_{j}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t})} \phi_{i}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t}) - e^{\sum_{j} \omega_{j} \times \phi_{j}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t})} \frac{\partial}{\partial \omega_{i}} Z_{\boldsymbol{x}_{t}}}{Z_{\boldsymbol{x}_{t}}^{2}}$$

$$= \frac{e^{\sum_{j} \omega_{j} \times \phi_{j}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t})}}{Z_{\boldsymbol{x}_{t}}^{2}} (Z_{\boldsymbol{x}_{t}} \phi_{i}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t}) - \frac{\partial}{\partial \omega_{i}} Z_{\boldsymbol{x}_{t}})$$

$$= \frac{e^{\sum_{j} \omega_{j} \times \phi_{j}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t})}}{Z_{\boldsymbol{x}_{t}}^{2}} (Z_{\boldsymbol{x}_{t}} \phi_{i}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t}) - \frac{\partial}{\partial \omega_{i}} Z_{\boldsymbol{x}_{t}})$$

because

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

# The partial derivatives

## The partial derivatives (3)

From (2),  $\frac{\partial}{\partial \omega_{i}} \frac{e^{\sum_{j} \omega_{j} \times \phi_{j}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t})}}{Z_{\boldsymbol{x}_{t}}} = \frac{e^{\sum_{j} \omega_{j} \times \phi_{j}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t})}}{Z_{\boldsymbol{x}_{t}}^{2}} (Z_{\boldsymbol{x}_{t}} \phi_{i}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t}) - \sum_{\boldsymbol{y}' \in \mathcal{Y}} e^{\sum_{j} \omega_{j} \times \phi_{j}(\boldsymbol{x}_{t}, \boldsymbol{y}')} \phi_{i}(\boldsymbol{x}_{t}, \boldsymbol{y}'))$ 

Sub this in (1),

$$\begin{aligned} \frac{\partial}{\partial \omega_{i}} \mathcal{L}(\mathcal{T}; \omega) &= \sum_{t} \left( \frac{Z_{x_{t}}}{e^{\sum_{j} \omega_{j} \times \phi_{j}(x_{t}, y_{t})}} \right) \left( \frac{\partial}{\partial \omega_{i}} \frac{e^{\sum_{j} \omega_{j} \times \phi_{j}(x_{t}, y_{t})}}{Z_{x_{t}}} \right) \\ &= \sum_{t} \frac{1}{Z_{x_{t}}} \left( Z_{x_{t}} \phi_{i}(x_{t}, y_{t}) - \sum_{y' \in \mathcal{Y}} e^{\sum_{j} \omega_{j} \times \phi_{j}(x_{t}, y')} \phi_{i}(x_{t}, y') \right) \right) \\ &= \sum_{t} \phi_{i}(x_{t}, y_{t}) - \sum_{t} \sum_{y' \in \mathcal{Y}} \frac{e^{\sum_{j} \omega_{j} \times \phi_{j}(x_{t}, y')}}{Z_{x_{t}}} \phi_{i}(x_{t}, y') \\ &= \sum_{t} \phi_{i}(x_{t}, y_{t}) - \sum_{t} \sum_{y' \in \mathcal{Y}} P(y'|x_{t}) \phi_{i}(x_{t}, y') \end{aligned}$$

# FINALLY!!!

After all that,

$$rac{\partial}{\partial oldsymbol{\omega}_i} \mathcal{L}(\mathcal{T};oldsymbol{\omega}) \hspace{0.2cm} = \hspace{0.2cm} \sum_t \phi_i(oldsymbol{x}_t,oldsymbol{y}_t) - \sum_t \sum_{oldsymbol{y}' \in \mathcal{Y}} oldsymbol{\mathcal{P}}(oldsymbol{y}'|oldsymbol{x}_t) \phi_i(oldsymbol{x}_t,oldsymbol{y}')$$

And the gradient is:

$$\nabla \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) = (\frac{\partial}{\partial \boldsymbol{\omega}_0} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}), \frac{\partial}{\partial \boldsymbol{\omega}_1} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}), \dots, \frac{\partial}{\partial \boldsymbol{\omega}_m} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}))$$

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

# Logistic Regression Summary

Define conditional probability

$$P(oldsymbol{y}|oldsymbol{x}) = rac{e^{oldsymbol{\omega}\cdot\phi(oldsymbol{x},oldsymbol{y})}}{Z_{oldsymbol{x}}}$$

Set weights to maximize log-likelihood of training data:

$$oldsymbol{\omega} = rgmax_{oldsymbol{\omega}} \sum_t \log P(oldsymbol{y}_t | oldsymbol{x}_t)$$

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

$$rac{\partial}{\partial oldsymbol{\omega}_i} \mathcal{L}(\mathcal{T};oldsymbol{\omega}) = \sum_t \phi_i(oldsymbol{x}_t,oldsymbol{y}_t) - \sum_t \sum_{oldsymbol{y}' \in \mathcal{Y}} \mathcal{P}(oldsymbol{y}'|oldsymbol{x}_t) \phi_i(oldsymbol{x}_t,oldsymbol{y}')$$
## Logistic Regression = Maximum Entropy

- Well-known equivalence
- Max Ent: maximize entropy subject to constraints on features: P = arg max<sub>P</sub> H(P) under constraints
  - Empirical feature counts must equal expected counts
- Quick intuition
  - Partial derivative in logistic regression

$$rac{\partial}{\partial oldsymbol{\omega}_i} \mathcal{L}(\mathcal{T};oldsymbol{\omega}) = \sum_t \phi_i(oldsymbol{x}_t,oldsymbol{y}_t) - \sum_t \sum_{oldsymbol{y}' \in \mathcal{Y}} \mathcal{P}(oldsymbol{y}'|oldsymbol{x}_t) \phi_i(oldsymbol{x}_t,oldsymbol{y}')$$

- 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

$$\mathcal{L}(\mathcal{T}; oldsymbol{\omega}) = \sum_{t=1}^{|\mathcal{T}|} 1 - \llbracket oldsymbol{y}_t = rgmax_{oldsymbol{y}} oldsymbol{\omega} \cdot oldsymbol{\phi}(oldsymbol{x}_t, oldsymbol{y}) 
rbrace$$

$$\boldsymbol{\omega} = \operatorname*{arg\,min}_{\boldsymbol{\omega}} \sum_{t=1}^{|\mathcal{T}|} 1 - \llbracket \boldsymbol{y}_t = \operatorname*{arg\,max}_{\boldsymbol{y}} \ \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}) \rrbracket$$
$$\llbracket \boldsymbol{p} \rrbracket = \left\{ \begin{array}{cc} 1 & \boldsymbol{p} \text{ is true} \\ 0 & \text{otherwise} \end{array} \right.$$

This is a 0-1 loss function

We will see later how to formalize the perceptron error function as hinge-loss

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:

A training set T is separable with margin γ > 0 if there exists a vector u with ||u|| = 1 such that:

$$\mathbf{u} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \mathbf{u} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}') \ge \gamma \tag{2}$$

for all  $oldsymbol{y}'\in ar{\mathcal{Y}}_t$  and  $||oldsymbol{u}||=\sqrt{\sum_joldsymbol{\mathsf{u}}_j^2}$ 

• Assumption: the training set is separable with margin  $\gamma$ 

#### Perceptron: Main Theorem

Theorem: For any training set separable with a margin of γ, the following holds for the perceptron algorithm:

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

where  $R \geq ||\phi(x_t,y_t) - \phi(x_t,y')||$  for all  $(x_t,y_t) \in \mathcal{T}$  and  $y' \in ar{\mathcal{Y}}_t$ 

- 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 2002])

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

Lower bound: 
$$\begin{split} & \boldsymbol{\omega}^{(k-1)} \text{ are weights before } k^{th} \text{ error} \\ & \text{Suppose } k^{th} \text{ error made at } (\boldsymbol{x}_t, \boldsymbol{y}_t) \\ & \boldsymbol{y}' = \arg \max_{\boldsymbol{y}'} \boldsymbol{\omega}^{(k-1)} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}') \\ & \boldsymbol{y}' \neq \boldsymbol{y}_t \\ & \boldsymbol{\omega}^{(k)} = \\ & \boldsymbol{\omega}^{(k-1)} + \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}') \end{split}$$

Training data: 
$$\mathcal{T} = \{(\boldsymbol{x}_t, \boldsymbol{y}_t)\}_{t=1}^{|\mathcal{T}|}$$
   
Lower bound:  
 $\boldsymbol{\omega}^{(k-1)}$  are we  
2. for  $n: 1..N$   
3. for  $t: 1..T$   
4. Let  $\boldsymbol{y}' = \arg \max_{\boldsymbol{y}'} \boldsymbol{\omega}^{(i)} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}')$   
5. if  $\boldsymbol{y}' \neq \boldsymbol{y}_t$   
6.  $\boldsymbol{\omega}^{(i+1)} = \boldsymbol{\omega}^{(i)} + \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}')$   
7.  $i = i+1$   
8. return  $\boldsymbol{\omega}^i$   
Lower bound:  
 $\boldsymbol{\omega}^{(k-1)}$  are we  
 $\boldsymbol{y}' = \arg \max_{\boldsymbol{y}'} \boldsymbol{y}_t$   
 $\boldsymbol{\psi}' = \arg \max_{\boldsymbol{y}'} \boldsymbol{\psi}_t$   
 $\boldsymbol{\psi}^{(k)} = \boldsymbol{\psi}^{(k-1)} + \boldsymbol{\phi}(\boldsymbol{x}_t) + \boldsymbol{\psi}(\boldsymbol{x}_t, \boldsymbol{y}_t)$ 

 $egin{aligned} & egin{aligned} & egi$ 

$$\begin{split} \mathbf{u} \cdot \boldsymbol{\omega}^{(k)} &= \mathbf{u} \cdot \boldsymbol{\omega}^{(k-1)} + \mathbf{u} \cdot (\boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}')) \geq \mathbf{u} \cdot \boldsymbol{\omega}^{(k-1)} + \boldsymbol{\gamma}, \text{ by (2)} \\ \text{Since } \boldsymbol{\omega}^{(0)} &= 0 \text{ and } \mathbf{u} \cdot \boldsymbol{\omega}^{(0)} = 0, \text{ for all } k: \mathbf{u} \cdot \boldsymbol{\omega}^{(k)} \geq k\boldsymbol{\gamma}, \text{ by induction on } k \\ \text{Since } \mathbf{u} \cdot \boldsymbol{\omega}^{(k)} \leq ||\mathbf{u}|| \times ||\boldsymbol{\omega}^{(k)}||, \text{ by the law of cosines, and } ||\mathbf{u}|| = 1, \text{ then } \\ ||\boldsymbol{\omega}^{(k)}|| \geq k\boldsymbol{\gamma} \end{split}$$

Training data: 
$$\mathcal{T} = \{(\boldsymbol{x}_{t}, \boldsymbol{y}_{t})\}_{t=1}^{|\mathcal{T}|}$$
  
1.  $\omega^{(0)} = 0; i = 0$   
2. for  $n: 1..N$   
3. for  $t: 1..T$   
4. Let  $\boldsymbol{y}' = \arg \max_{\boldsymbol{y}'} \omega^{(i)} \cdot \phi(\boldsymbol{x}_{t}, \boldsymbol{y}')$   
5. if  $\boldsymbol{y}' \neq \boldsymbol{y}_{t}$   
6.  $\omega^{(i+1)} = \omega^{(i)} + \phi(\boldsymbol{x}_{t}, \boldsymbol{y}_{t}) - \phi(\boldsymbol{x}_{t}, \boldsymbol{y}')$   
7.  $i = i+1$   
8. return  $\omega^{i}$   
Lower bound:  
 $\omega^{(k-1)}$  are weights before  $k^{th}$  error made at  $(\boldsymbol{x}_{t}, \boldsymbol{y}_{t})$   
 $\boldsymbol{y}' = \arg \max_{\boldsymbol{y}'} \omega^{(k-1)} \cdot \phi(\boldsymbol{x}_{t}, \boldsymbol{y}')$   
 $\boldsymbol{y}' \neq \boldsymbol{y}_{t}$   
 $\omega^{(k)} = \omega^{(k-1)} + \phi(\boldsymbol{x}_{t}, \boldsymbol{y}_{t}) - \phi(\boldsymbol{x}_{t}, \boldsymbol{y}')$ 

$$\begin{split} \mathbf{u} \cdot \boldsymbol{\omega}^{(k)} &= \mathbf{u} \cdot \boldsymbol{\omega}^{(k-1)} + \mathbf{u} \cdot (\boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}')) \geq \mathbf{u} \cdot \boldsymbol{\omega}^{(k-1)} + \gamma, \text{ by (2)} \\ \text{Since } \boldsymbol{\omega}^{(0)} &= 0 \text{ and } \mathbf{u} \cdot \boldsymbol{\omega}^{(0)} = 0, \text{ for all } k: \mathbf{u} \cdot \boldsymbol{\omega}^{(k)} \geq k\gamma, \text{ by induction on } k \\ \text{Since } \mathbf{u} \cdot \boldsymbol{\omega}^{(k)} \leq ||\mathbf{u}|| \times ||\boldsymbol{\omega}^{(k)}||, \text{ by the law of cosines, and } ||\mathbf{u}|| = 1, \text{ then } \\ ||\boldsymbol{\omega}^{(k)}|| \geq k\gamma \end{split}$$

Upper bound:

$$\begin{aligned} ||\omega^{(k)}||^2 &= ||\omega^{(k-1)}||^2 + ||\phi(x_t, y_t) - \phi(x_t, y')||^2 + 2\omega^{(k-1)} \cdot (\phi(x_t, y_t) - \phi(x_t, y')) \\ ||\omega^{(k)}||^2 &\leq ||\omega^{(k-1)}||^2 + R^2, \text{ since } R \geq ||\phi(x_t, y_t) - \phi(x_t, y')|| \\ &\quad \text{ and } \omega^{(k-1)} \cdot \phi(x_t, y_t) - \omega^{(k-1)} \cdot \phi(x_t, y') \leq 0 \\ &\leq kR^2 \text{ for all } k, \text{ by induction on } k \end{aligned}$$

• We have just shown that  $||\omega^{(k)}|| \ge k\gamma$  and  $||\omega^{(k)}||^2 \le kR^2$ 

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

and solving for k

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

Therefore the number of errors is bounded!

#### **Perceptron Summary**

- Learns parameters of a linear model by minimizing error
- Guaranteed to find a  $\omega$  in a finite amount of time
- Perceptron is an example of an Online Learning Algorithm
  - $\blacktriangleright \, \omega$  is updated based on a single training instance in isolation

$$oldsymbol{\omega}^{(i+1)} = oldsymbol{\omega}^{(i)} + \phi(oldsymbol{x}_t,oldsymbol{y}_t) - \phi(oldsymbol{x}_t,oldsymbol{y}')$$

#### **Averaged** Perceptron

```
Training data: \mathcal{T} = \{(\boldsymbol{x}_t, \boldsymbol{y}_t)\}_{t=1}^{|\mathcal{T}|}
  1. \omega^{(0)} = 0; i = 0
 2. for n: 1..N
 3. for t: 1..T
     Let oldsymbol{y}' = rg\max_{oldsymbol{u}'} oldsymbol{\omega}^{(i)} \cdot oldsymbol{\phi}(oldsymbol{x}_t,oldsymbol{y}')
  4.
 5.
               if u' \neq u_t
                   \omega^{(i+1)}=\omega^{(i)}+\phi(x_t,y_t)-\phi(x_t,y')
 6.
 7.
      else
      \omega^{(i+1)} = \omega^{(i)}
 6.
 7. i = i + 1
```

8. return  $\left(\sum_{i} \omega^{(i)}\right) / (N \times T)$ 

# Margin



# **Maximizing Margin**

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

$$\boldsymbol{\omega}\cdot\boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \boldsymbol{\omega}\cdot\boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}') \geq \gamma$$

ullet for every training instance  $(m{x}_t,m{y}_t)\in\mathcal{T}$  ,  $m{y}'\inar{\mathcal{Y}}_t$ 

# Maximizing Margin

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

$$\epsilon \propto rac{R^2}{\gamma^2 imes |\mathcal{T}|}$$

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

# Support Vector Machines (SVMs)

# **Maximizing Margin**

Let  $\gamma > 0$ 

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

$$egin{aligned} oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t,oldsymbol{y}_t) &= eta \cdot \phi(oldsymbol{x}_t,oldsymbol{y}_t) \geq \gamma \ & orall (oldsymbol{x}_t,oldsymbol{y}_t) \in \mathcal{T} \ & ext{ and } oldsymbol{y}' \in ar{\mathcal{Y}}_t \end{aligned}$$

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

$$eta(oldsymbol{\omega}\cdotoldsymbol{\phi}(oldsymbol{x}_t,oldsymbol{y}_t)) \geq eta\gamma$$
 , for some  $eta\geq 1$ 

Let  $\gamma > 0$ 

#### Max Margin:

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

$$egin{aligned} oldsymbol{\omega} \cdot oldsymbol{\phi}(oldsymbol{x}_t,oldsymbol{y}_t) - oldsymbol{\omega} \cdot oldsymbol{\phi}(oldsymbol{x}_t,oldsymbol{y}_t) & \geq \gamma \ & orall (oldsymbol{x}_t,oldsymbol{y}_t) \in \mathcal{T} \ & ext{ and } oldsymbol{y}' \in ar{\mathcal{Y}}_t \end{aligned}$$

Let  $\gamma > 0$ 

#### Max Margin:

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

$$egin{aligned} & \omega{\cdot}\phi(x_t,y_t){-}\omega{\cdot}\phi(x_t,y')\geq\gamma \ & \forall (x_t,y_t)\in\mathcal{T} \ & ext{ and } y'\inar{\mathcal{Y}}_t \end{aligned}$$
 Change variables:  $\mathbf{u}=rac{\omega}{\gamma}?$   $||\omega||=1 ext{ iff } ||\mathbf{u}||=1/\gamma, \ & ext{ then } \gamma=1/||\mathbf{u}|| \end{aligned}$ 

Let  $\gamma > 0$ 

Max Margin:

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

such that:

$$egin{aligned} & \omega \cdot \phi(m{x}_t,m{y}_t) - \omega \cdot \phi(m{x}_t,m{y}') \geq \gamma \ & orall (m{x}_t,m{y}_t) \in \mathcal{T} \ & ext{ and } m{y}' \in ar{\mathcal{Y}}_t \end{aligned}$$
Change variables:  $m{u} = rac{m{\omega}}{\gamma}?$ 
 $||m{\omega}|| = 1 ext{ iff } ||m{u}|| = 1/\gamma, \ & ext{then } \gamma = 1/||m{u}|| \end{aligned}$ 

Min Norm (step 1):

$$\underset{\boldsymbol{u}}{\text{max}} \quad \frac{1}{||\boldsymbol{u}||} = \underset{\boldsymbol{u}}{\text{min}} \left||\boldsymbol{u}|\right|$$

$$egin{aligned} & \omega{\cdot}\phi(m{x}_t,m{y}_t){-}\omega{\cdot}\phi(m{x}_t,m{y}') \geq \gamma \ & orall (m{x}_t,m{y}_t) \in \mathcal{T} \ & ext{ and } m{y}' \in ar{\mathcal{Y}}_t \end{aligned}$$

Let  $\gamma > 0$ 

Max Margin:

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

such that:

$$egin{aligned} & \omega \cdot \phi(m{x}_t,m{y}_t) - \omega \cdot \phi(m{x}_t,m{y}') \geq \gamma \ & orall (m{x}_t,m{y}_t) \in \mathcal{T} \ & ext{ and } m{y}' \in ar{\mathcal{Y}}_t \end{aligned}$$
Change variables:  $m{u} = rac{m{\omega}}{\gamma}?$ 
 $& ||m{\omega}|| = 1 ext{ iff } ||m{u}|| = 1/\gamma, \end{aligned}$ 
then  $\gamma = 1/||m{u}||$ 

Min Norm (step 2):  $\min_{\mathbf{U}} ||\mathbf{u}||$ such that:  $\gamma \mathbf{u} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \gamma \mathbf{u} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}') > \gamma$ 

$$orall (oldsymbol{x}_t,oldsymbol{y}_t)\in\mathcal{T}$$
 and  $oldsymbol{y}'\inar{\mathcal{Y}}_t$ 

Let  $\gamma > 0$ 

Max Margin:

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

such that:

$$egin{aligned} & \omega{\cdot}\phi(m{x}_t,m{y}_t){-}\omega{\cdot}\phi(m{x}_t,m{y}')\geq\gamma \ & \forall(m{x}_t,m{y}_t)\in\mathcal{T} \ & ext{ and }m{y}'\inar{\mathcal{Y}}_t \end{aligned}$$
Change variables:  $m{u}=rac{\omega}{\gamma}?$ 
 $||m{\omega}||=1 ext{ iff }||m{u}||=1/\gamma, \ & ext{then }\gamma=1/||m{u}|| \end{aligned}$ 

Min Norm (step 2):  $\begin{array}{c} \min_{\mathbf{u}} ||\mathbf{u}|| \\ \text{such that:} \\ \mathbf{u} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}_t) - \mathbf{u} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}') \geq 1 \\ \forall (\boldsymbol{x}_t, \boldsymbol{y}_t) \in \mathcal{T} \\ \text{and } \boldsymbol{y}' \in \bar{\mathcal{Y}}_t \end{array}$ 

Let  $\gamma > 0$ 

Max Margin:

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

such that:

$$egin{aligned} & \omega{\cdot}\phi(m{x}_t,m{y}_t){-}\omega{\cdot}\phi(m{x}_t,m{y}')\geq\gamma \ & \forall(m{x}_t,m{y}_t)\in\mathcal{T} \ & ext{ and }m{y}'\inar{\mathcal{Y}}_t \end{aligned}$$
Change variables:  $m{u}=rac{\omega}{\gamma}?$ 
 $||\omega||=1 ext{ iff }||m{u}||=1/\gamma, \ & ext{then }\gamma=1/||m{u}|| \end{aligned}$ 

Min Norm (step 3):  $\min_{\mathbf{u}} \quad \frac{1}{2} ||\mathbf{u}||^2$ 

$$egin{aligned} \mathsf{u}{\cdot}\phi(m{x}_t,m{y}_t){-}\mathsf{u}{\cdot}\phi(m{x}_t,m{y}') \geq 1 \ & orall (m{x}_t,m{y}_t) \in \mathcal{T} \ & ext{and} \ m{y}' \in ar{\mathcal{Y}}_t \end{aligned}$$

Let  $\gamma > 0$ 

| Max Margin:   | Min Norm:   |
|---|---|
| $\max_{  \boldsymbol{\omega}  =1} \gamma$   | $\min_{\mathbf{u}}  \frac{1}{2}   \mathbf{u}  ^2$                       |
| such that:  | such that:  |
| $oldsymbol{\omega}{\cdot}\phi(oldsymbol{x}_t,oldsymbol{y}_t){-}\omega{\cdot}\phi(oldsymbol{x}_t,oldsymbol{y}')\geq\gamma$ | $\mathbf{u}{\cdot}\phi(x_t,y_t){-}\mathbf{u}{\cdot}\phi(x_t,y') \geq 1$ |
| $orall (oldsymbol{x}_t,oldsymbol{y}_t)\in\mathcal{T}$  | $orall (oldsymbol{x}_t,oldsymbol{y}_t)\in\mathcal{T}$                  |
| and $oldsymbol{y}'\in ar{\mathcal{Y}}_t$  | and $oldsymbol{y}'\in ar{\mathcal{Y}}_t$                                |

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

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

$$oldsymbol{\omega} = \operatorname*{arg\,min}_{oldsymbol{\omega},\xi} \; rac{1}{2} ||oldsymbol{\omega}||^2 + oldsymbol{C} \sum_{t=1}^{|\mathcal{T}|} oldsymbol{\xi}_t$$

· ----

such that:

$$egin{aligned} &\omega\cdot\phi(x_t,y_t)-\omega\cdot\phi(x_t,y')\geq 1-\xi_t ext{ and } \xi_t\geq 0 \ &orall (x_t,y_t)\in\mathcal{T} ext{ and } y'\inar{\mathcal{Y}}_t \end{aligned}$$

 $\xi_t$ : slack variable representing amount of constraint violation If data is separable, optimal solution has  $\xi_i = 0$ ,  $\forall i$ C balances focus on margin  $(C < \frac{1}{2})$  and on error  $(C > \frac{1}{2})$ 

$$oldsymbol{\omega} = rgmin_{oldsymbol{\omega},\xi} \; rac{1}{2} ||oldsymbol{\omega}||^2 + C \sum_{t=1}^{|\mathcal{T}|} \xi_t$$

such that:

$$oldsymbol{\omega}\cdot \phi(oldsymbol{x}_t,oldsymbol{y}_t) - oldsymbol{\omega}\cdot \phi(oldsymbol{x}_t,oldsymbol{y}') \geq 1-\xi_t$$
  
where  $\xi_t \geq 0$  and  $orall(oldsymbol{x}_t,oldsymbol{y}_t) \in \mathcal{T}$  and  $oldsymbol{y}'\in ar{\mathcal{Y}}_t$ 

 Computing the dual form results in a quadratic programming problem – a well-known convex optimization problem [Boyd and Vandenberghe 2004]

Can we have representation of this objective that allows more direct optimization?

$$\boldsymbol{\omega} = \operatorname*{arg\,min}_{\boldsymbol{\omega},\boldsymbol{\xi}} \; \frac{1}{2} ||\boldsymbol{\omega}||^2 + C \sum_{t=1}^{|\mathcal{T}|} \xi_t$$

$$egin{aligned} & \omega \cdot \phi(x_t,y_t) - \max_{oldsymbol{y} 
eq y_t} \ \omega \cdot \phi(x_t,y') \geq 1 - arepsilon_t \end{aligned}$$

$$oldsymbol{\omega} = rgmin_{oldsymbol{\omega},\xi} rac{1}{2} ||oldsymbol{\omega}||^2 + C \sum_{t=1}^{|\mathcal{T}|} \xi_t$$

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

. \_\_\_\_

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

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

- If  $\|\omega\|$  classifies  $(x_t, y_t)$  with margin 1, penalty  $\xi_t = 0$
- ► Otherwise:  $\xi_t = 1 + \max_{y' \neq y_t} \omega \cdot \phi(x_t, y') \omega \cdot \phi(x_t, y_t)$
- That means that in the end  $\xi_t$  will be:

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

$$oldsymbol{\omega} = rgmin_{oldsymbol{\omega},\xi} \; rac{\lambda}{2} ||oldsymbol{\omega}||^2 + \sum_{t=1}^{|\mathcal{T}|} \xi_t \; ext{s.t.} \; \xi_t \geq 1 + \max_{oldsymbol{y}' 
eq oldsymbol{y}_t} \; oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t,oldsymbol{y}') - oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t,oldsymbol{y}_t)$$

#### Hinge loss

$$egin{aligned} & \omega = rgmin_{oldsymbol{\omega}} \ \mathcal{L}(\mathcal{T};oldsymbol{\omega}) = rgmin_{oldsymbol{\omega}} \ \sum_{t=1}^{|\mathcal{T}|} \mathit{loss}((oldsymbol{x}_t,oldsymbol{y}_t);oldsymbol{\omega}) \ + \ rac{\lambda}{2}||oldsymbol{\omega}||^2 \ = rgmin_{oldsymbol{\omega}} \ \left(\sum_{t=1}^{|\mathcal{T}|} \max{(0,1+\max_{oldsymbol{y}
eq oldsymbol{y}_t,oldsymbol{y}_t)} \ \omega\cdot\phi(oldsymbol{x}_t,oldsymbol{y}') - \omega\cdot\phi(oldsymbol{x}_t,oldsymbol{y}_t))}
ight) + rac{\lambda}{2}||oldsymbol{\omega}||^2 \end{aligned}$$

We will see later how to do efficient optimization of hinge loss

# Summary

#### What we have covered

- Linear Learners
  - Naive Bayes
  - Logistic Regression
  - Perceptron
  - Support Vector Machines

#### What is next

- Regularization
- Online learning
- Non-linear models

# Regularization

#### Fit of a Model



- 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

$$\underset{\boldsymbol{\omega}}{\operatorname{arg\,max}} \ \mathcal{L}(\mathcal{T};\boldsymbol{\omega}) - \lambda \mathcal{R}(\boldsymbol{\omega})$$

- 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
  - $\blacktriangleright$  L0:  $\mathcal{R}(\boldsymbol{\omega}) \propto \|\boldsymbol{\omega}\|_0 = \sum_i \llbracket \boldsymbol{\omega}_i > 0 \rrbracket$  zero weights desired
    - Non-convex
    - Approximate with L1:  $\mathcal{R}(\boldsymbol{\omega}) \propto \|\boldsymbol{\omega}\|_1 = \sum_i |\boldsymbol{\omega}_i|$
# Logistic Regression with L2 Regularization

Perhaps most common learner in NLP

$$\mathcal{L}(\mathcal{T}; oldsymbol{\omega}) - \lambda \mathcal{R}(oldsymbol{\omega}) = \sum_{t=1}^{|\mathcal{T}|} \log \left( e^{oldsymbol{\omega} \cdot oldsymbol{\phi}(oldsymbol{x}_t, oldsymbol{y}_t)} / Z_{oldsymbol{x}} 
ight) - rac{\lambda}{2} \|oldsymbol{\omega}\|^2$$

What are the new partial derivatives?

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

► We know 
$$\frac{\partial}{\partial w_i} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega})$$
  
► Just need  $\frac{\partial}{\partial w_i} \frac{\lambda}{2} \|\boldsymbol{\omega}\|^2 = \frac{\partial}{\partial w_i} \frac{\lambda}{2} \left(\sqrt{\sum_i \boldsymbol{\omega}_i^2}\right)^2 = \frac{\partial}{\partial w_i} \frac{\lambda}{2} \sum_i \boldsymbol{\omega}_i^2 = \lambda \boldsymbol{\omega}_i$ 

# **Support Vector Machines**

 SVM in hinge-loss formulation: L2 regularization corresponds to margin maximization!

$$\begin{split} \boldsymbol{\omega} &= \operatorname*{arg\,min}_{\boldsymbol{\omega}} \ \mathcal{L}(\mathcal{T};\boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega}) \\ &= \operatorname{arg\,min}_{\boldsymbol{\omega}} \ \sum_{t=1}^{|\mathcal{T}|} loss((\boldsymbol{x}_t,\boldsymbol{y}_t);\boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega}) \\ &= \operatorname{arg\,min}_{\boldsymbol{\omega}} \ \sum_{t=1}^{|\mathcal{T}|} \max\left(0, 1 + \max_{\boldsymbol{y} \neq \boldsymbol{y}_t} \ \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t,\boldsymbol{y}) - \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t,\boldsymbol{y}_t)\right) + \lambda \mathcal{R}(\boldsymbol{\omega}) \\ &= \operatorname{arg\,min}_{\boldsymbol{\omega}} \ \sum_{t=1}^{|\mathcal{T}|} \max\left(0, 1 + \max_{\boldsymbol{y} \neq \boldsymbol{y}_t} \ \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t,\boldsymbol{y}) - \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t,\boldsymbol{y}_t)\right) + \lambda \mathcal{R}(\boldsymbol{\omega}) \end{split}$$

# SVMs vs. Logistic Regression

$$\begin{split} \boldsymbol{\omega} &= \mathop{\arg\min}_{\boldsymbol{\omega}} \ \mathcal{L}(\mathcal{T};\boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega}) \\ &= \mathop{\arg\min}_{\boldsymbol{\omega}} \ \sum_{t=1}^{|\mathcal{T}|} loss((\boldsymbol{x}_t,\boldsymbol{y}_t);\boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega}) \end{split}$$

# SVMs vs. Logistic Regression

$$\begin{aligned} \boldsymbol{\omega} &= \operatorname*{arg\,min}_{\boldsymbol{\omega}} \ \mathcal{L}(\mathcal{T};\boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega}) \\ &= \operatorname{arg\,min}_{\boldsymbol{\omega}} \ \sum_{t=1}^{|\mathcal{T}|} \mathit{loss}((\boldsymbol{x}_t,\boldsymbol{y}_t);\boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega}) \end{aligned}$$

 $\mathsf{SVMs}/\mathsf{hinge-loss:}\;\max\left(0,1+\mathsf{max}_{\bm{y}\neq\bm{y}_t}\;\left(\bm{\omega}\cdot\bm{\phi}(\bm{x}_t,\bm{y})-\bm{\omega}\cdot\bm{\phi}(\bm{x}_t,\bm{y}_t)\right)\right)$ 

$$oldsymbol{\omega} = rgmin_{oldsymbol{\omega}} \sum_{t=1}^{|\mathcal{T}|} \max \left( 0, 1 + \max_{oldsymbol{y}
eq oldsymbol{y}_t} oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t,oldsymbol{y}) - oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t,oldsymbol{y}_t) + rac{\lambda}{2} \|oldsymbol{\omega}\|^2$$

# SVMs vs. Logistic Regression

$$\begin{aligned} \boldsymbol{\omega} &= \operatorname*{arg\,min}_{\boldsymbol{\omega}} \ \mathcal{L}(\mathcal{T};\boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega}) \\ &= \operatorname{arg\,min}_{\boldsymbol{\omega}} \ \sum_{t=1}^{|\mathcal{T}|} \mathit{loss}((\boldsymbol{x}_t,\boldsymbol{y}_t);\boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega}) \end{aligned}$$

 $\mathsf{SVMs}/\mathsf{hinge-loss:}\;\max\left(0,1+\mathsf{max}_{\bm{y}\neq\bm{y}_t}\;(\bm{\omega}\cdot\bm{\phi}(\bm{x}_t,\bm{y})-\bm{\omega}\cdot\bm{\phi}(\bm{x}_t,\bm{y}_t))\right)$ 

$$oldsymbol{\omega} = rgmin_{oldsymbol{\omega}} \; \sum_{t=1}^{|\mathcal{T}|} \max \left( 0, 1 + \max_{oldsymbol{y} 
eq oldsymbol{y}_t} \; oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t,oldsymbol{y}) - oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t,oldsymbol{y}_t) + rac{\lambda}{2} \|oldsymbol{\omega}\|^2$$

 $\text{Logistic Regression}/\text{log-loss:} - \text{log} \; \left( e^{\boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t)} / Z_{\boldsymbol{x}} \right)$ 

$$oldsymbol{\omega} = rgmin_{oldsymbol{\omega}} \sum_{t=1}^{|\mathcal{T}|} - \log \left( e^{oldsymbol{\omega} \cdot oldsymbol{\phi}(oldsymbol{x}_t, oldsymbol{y}_t)} / Z_{oldsymbol{x}} 
ight) + rac{\lambda}{2} \|oldsymbol{\omega}\|^2$$

#### **Generalized Linear Learners**

$$oldsymbol{\omega} = rgmin_{oldsymbol{\omega}} \mathcal{L}(\mathcal{T};oldsymbol{\omega}) + \lambda \mathcal{R}(oldsymbol{\omega}) = rgmin_{oldsymbol{\omega}} \sum_{t=1}^{|\mathcal{T}|} loss((oldsymbol{x}_t,oldsymbol{y}_t);oldsymbol{\omega}) + \lambda \mathcal{R}(oldsymbol{\omega})$$



# Which Learner to Use?

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

# **Online Learning**

# **Online vs. Batch Learning**

 $\mathsf{Batch}(\mathcal{T});$ 

▶ for 1 ... N

 $\blacktriangleright \ \omega \leftarrow \mathsf{update}(\mathcal{T}; \omega)$ 

 $\blacktriangleright$  return  $\omega$ 

Online( $\mathcal{T}$ );

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

- E.g., SVMs, logistic regression, Naive Bayes
- E.g., Perceptron $\omega = \omega + \phi(x_t, y_t) \phi(x_t, y)$

#### **Batch Gradient Descent**

► Let 
$$\mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) = \sum_{t=1}^{|\mathcal{T}|} loss((\boldsymbol{x}_t, \boldsymbol{y}_t); \boldsymbol{\omega})$$
  
► Set  $\boldsymbol{\omega}^0 = O^m$ 

Iterate until convergence

$$egin{array}{rcl} oldsymbol{\omega}^{i} &=& oldsymbol{\omega}^{i-1} - lpha 
abla \mathcal{L}(\mathcal{T};oldsymbol{\omega}^{i-1}) \ &=& oldsymbol{\omega}^{i-1} - \sum_{t=1}^{|\mathcal{T}|} lpha 
abla oldsymbol{loss}((oldsymbol{x}_t,oldsymbol{y}_t);oldsymbol{\omega}^{i-1}) \end{array}$$

•  $\alpha > 0$  and set so that  $\mathcal{L}(\mathcal{T}; \boldsymbol{\omega}^i) < \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}^{i-1})$ 

# **Stochastic Gradient Descent**

Approximate batch gradient ∇L(T; ω) with stochastic gradient ∇loss((x<sub>t</sub>, y<sub>t</sub>); ω)

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

 $\blacktriangleright$  return  $\omega$ 

# SGD over Finite Training Data

#### NLP practice: Cycling over finite data set

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

 $\blacktriangleright$  return  $\omega$ 

# **Online Logistic Regression**

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

$$abla \left( -\log \left( e^{oldsymbol{\omega} \cdot oldsymbol{\phi}(oldsymbol{x}_t,oldsymbol{y}_t)} 
ight) = -\left( \phi(oldsymbol{x}_t,oldsymbol{y}_t) - \sum_{oldsymbol{y}} oldsymbol{P}(oldsymbol{y}|oldsymbol{x}) \phi(oldsymbol{x}_t,oldsymbol{y}) 
ight)$$

Plus regularization term (if part of model)

# Online SVMs

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

$$egin{aligned} & au ext{loss}((m{x}_t,m{y}_t);m{\omega}) = igvee \left( \max\left(0,1+\max_{m{y}
eq m{y}_t} m{\omega}\cdot m{\phi}(m{x}_t,m{y}) - m{\omega}\cdot m{\phi}(m{x}_t,m{y}_t) 
ight) 
ight) \end{aligned}$$

Subgradient is:

$$egin{aligned} & 
abla \left( \mathsf{max} \; (0, 1 + \max_{oldsymbol{y} 
eq oldsymbol{y}_t} \; oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t, oldsymbol{y}) - oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t, oldsymbol{y}_t)) 
ight) \end{aligned}$$

$$= \begin{cases} 0, & \text{if } \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \max_{\boldsymbol{y}} \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}) \geq 1 \\ \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}) - \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t), & \text{otherwise, where } \boldsymbol{y} = \max_{\boldsymbol{y}} \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}) \end{cases}$$

Plus regularization term (required for SVMs)

#### Perceptron and Hinge-Loss

SVM subgradient update looks like perceptron update

$$egin{aligned} & egin{aligned} \omega^i = \omega^{i-1} - lpha iggl\{ 0, & ext{if } oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t, oldsymbol{y}_t) - ext{max}_{oldsymbol{y}} oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t, oldsymbol{y}_t) \geq \mathbf{1} \ \phi(oldsymbol{x}_t, oldsymbol{y}_t) - \phi(oldsymbol{x}_t, oldsymbol{y}_t), & ext{otherwise, where } oldsymbol{y} = ext{max}_{oldsymbol{y}} oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t, oldsymbol{y}_t) \geq \mathbf{1} \ \phi(oldsymbol{x}_t, oldsymbol{y}_t) - \phi(oldsymbol{x}_t, oldsymbol{y}_t), & ext{otherwise, where } oldsymbol{y} = ext{max}_{oldsymbol{y}} oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t, oldsymbol{y}_t) \geq \mathbf{1} \ \phi(oldsymbol{x}_t, oldsymbol{y}_t) - \phi(oldsymbol{x}_t, oldsymbol{y}_t), & ext{otherwise, where } oldsymbol{y} = ext{max}_{oldsymbol{y}} oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t, oldsymbol{y}_t) > \mathbf{1} \ \phi(oldsymbol{x}_t, oldsymbol{y}_t) - oldsymbol{x}_t, oldsymbol{y}_t), & ext{otherwise, where } oldsymbol{y} = ext{max}_{oldsymbol{y}} oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t, oldsymbol{y}_t) > \mathbf{1} \ \phi(oldsymbol{x}_t, oldsymbol{y}_t) - oldsymbol{x}_t, oldsymbol{y}_t), & ext{otherwise, where } oldsymbol{y} = ext{max}_{oldsymbol{y}} oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t, oldsymbol{y}_t), & ext{otherwise, where } oldsymbol{y} = ext{max}_{oldsymbol{y}} oldsymbol{\omega} \cdot \phi(oldsymbol{x}_t, oldsymbol{y}_t) & ext{otherwise, where } oldsymbol{y} = ext{max}_{oldsymbol{y}} oldsymbol{x}_t, oldsymbol{y}_t) & ext{otherwise, where } oldsymbol{y} = ext{max}_{oldsymbol{y}} oldsymbol{w}_t, oldsymbol{y}_t, & ext{otherwise, where } oldsymbol{y} = ext{max}_{oldsymbol{y}} oldsymbol{x}_t, oldsymbol{y}_t, & ext{otherwise, where } oldsymbol{y} = ext{max}_{oldsymbol{w}} oldsymbol{x}_t, & ext{max}_{oldsymbol{w}} oldsymbol{x}_t, & ext{max}_{oldsymbol{w}} oldsymbol{x}_t, & ext{max}_{oldsymbol{w}} oldsymbol{x}_t, & ext{max}_{oldsymbol{w}} oldsymbol{w}_t, & ext{max}_{oldsymbol{w}} oldsymbol{w}_t, & ext{max}_{oldsymbol{w}} oldsymbol{w}_t, & ext{max}_{oldsymbol{w}} oldsymbol$$

Perceptron

$$\boldsymbol{\omega}^{i} = \boldsymbol{\omega}^{i-1} - \alpha \begin{cases} 0, & \text{if } \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t}) - \max_{\boldsymbol{y}} \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_{t}, \boldsymbol{y}) \geq \boldsymbol{0} \\ \boldsymbol{\phi}(\boldsymbol{x}_{t}, \boldsymbol{y}) - \boldsymbol{\phi}(\boldsymbol{x}_{t}, \boldsymbol{y}_{t}), & \text{otherwise, where } \boldsymbol{y} = \max_{\boldsymbol{y}} \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_{t}, \boldsymbol{y}) \end{cases}$$

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

$$\max \left(0, 1 + \max_{oldsymbol{y} 
eq oldsymbol{y}_t} oldsymbol{\omega} \cdot oldsymbol{\phi}(oldsymbol{x}_t, oldsymbol{y}) - oldsymbol{\omega} \cdot oldsymbol{\phi}(oldsymbol{x}_t, oldsymbol{y}_t) 
ight)$$

# **Online vs. Batch Learning**

#### Online algorithms

- Each update step relies only on the derivative for a single randomly chosen example
  - Computational cost of one step is  $1/\mathcal{T}$  compared to batch
  - Easier to implement
- Larger variance since each gradient is different
  - Variance slows down convergence
  - Requires fine-tuning of decaying learning rate
- Batch algorithms
  - $\blacktriangleright$  Higher cost of averaging gradients over  ${\cal T}$  for each update
    - Implementation more complex
    - Less fine-tuning, e.g., allows constant learning rates
    - Faster convergence

# Variance-Reduced Online Learning

▶ SGD update extended by velocity vector v weighted by momentum coefficient  $0 \le \mu < 1$  [Polyak 1964]:

$$\boldsymbol{\omega}^{i+1} = \boldsymbol{\omega}^i - \alpha \nabla loss((\boldsymbol{x}_t, \boldsymbol{y}_t); \boldsymbol{\omega}^i) + \mu \boldsymbol{v}^i$$

where

►

$$m{v}^i=m{\omega}^i-m{\omega}^{i-1}$$

- Momentum accelerates learning if gradients are aligned along same direction, and restricts changes when successive gradient are opposite of each other
- General direction of gradient reinforced, perpendicular directions filtered out
- Best of both worlds: Efficient and effective!

# **Online-to-Batch Conversion**

Classical online learning:

- data are given as an infinite sequence of input examples
- model makes prediction on next example in sequence
- online error is averaged over predictions after each update
- Standard NLP applications:
  - Finite set of training data, prediction on new batch of test data
  - Online learning applied by cycling over finite data
  - Online-to-batch conversion: Which model to use at test time?
    - Last model? Random model? Best model on heldout set?

# **Online-to-Batch Conversion by Averaging**

Averaged Perceptron

• 
$$\bar{\boldsymbol{\omega}} = \left(\sum_{i} \boldsymbol{\omega}^{(i)}\right) / (N \times T)$$

Use weight vector averaged over online updates for prediction

How does the perceptron mistake bound carry over to batch?

▶ Let  $M_K$  be number of mistakes made during online learning, then with probability of at least  $1 - \delta$ :

$$\mathbb{E}[\mathit{loss}((oldsymbol{x},oldsymbol{y});ar{oldsymbol{\omega}}] \leq M_k + \sqrt{rac{2}{k}\lnrac{1}{\delta}}$$

- generalization bound based on online performance [Cesa-Bianchi et al. 2004]
- can be applied to all online learners with convex losses

# Quick Summary

#### Linear Learners

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

# Non-Linear Models

#### **Non-Linear Models**

- Some data sets require more than a linear decision boundary 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:

$$K(x_t, x_r) \in \mathbb{R}$$

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

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

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

$$K(x_t, x_r) = K(x_r, x_t)$$

Positive definite: positivity on diagonal

 $\mathcal{K}({m{x}},{m{x}}) \geq 0$  forall  ${m{x}}$  with equality only for  ${m{x}}=0$ 

# Kernels

Mercer's Theorem: for any kernel K, there exists an φ, in some R<sup>d</sup>, such that:

$$\mathcal{K}(x_t,x_r)=\phi(x_t)\cdot\phi(x_r)$$

Since our features are over pairs (x, y), we will write kernels over pairs

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

# 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!

Training data:  $\mathcal{T} = \{(x_t, y_t)\}_{t=1}^{|\mathcal{T}|}$ 1.  $\omega^{(0)} = 0; i = 0$ 2. for n: 1..N3. for t: 1..T4. 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 + 18. 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
  - ▶ α<sub>y,t</sub> is the # of times during learning label y is predicted for example t

Thus,

$$oldsymbol{\omega} = \sum_{t,oldsymbol{y}} lpha_{oldsymbol{y},t} [\phi(oldsymbol{x}_t,oldsymbol{y}_t) - \phi(oldsymbol{x}_t,oldsymbol{y})]$$

• We can re-write the argmax function as:  $y* = \arg \max_{y^*} \omega^{(i)} \cdot \phi(x, y^*)$ 

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

=

=

=

We can re-write the argmax function as:

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

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

. \_ .

Training data: 
$$\mathcal{T} = \{(\boldsymbol{x}_t, \boldsymbol{y}_t)\}_{t=1}^{|\mathcal{T}|}$$
  
1.  $\forall \boldsymbol{y}, t \text{ set } \alpha_{\boldsymbol{y},t} = 0$   
2. for  $n: 1..N$   
3. for  $t: 1..T$   
4. Let  $\boldsymbol{y}^* = \arg \max_{\boldsymbol{y}^*} \sum_{t,\boldsymbol{y}} \alpha_{\boldsymbol{y},t} [\mathcal{K}((\boldsymbol{x}_t, \boldsymbol{y}_t), (\boldsymbol{x}_t, \boldsymbol{y}^*)) - \mathcal{K}((\boldsymbol{x}_t, \boldsymbol{y}), (\boldsymbol{x}_t, \boldsymbol{y}^*))]$   
5. if  $\boldsymbol{y}^* \neq \boldsymbol{y}_t$   
6.  $\alpha_{\boldsymbol{y}^*,t} = \alpha_{\boldsymbol{y}^*,t} + 1$ 

Given a new instance x

$$oldsymbol{y}^* = rgmax_{oldsymbol{y}^*} \sum_{t,oldsymbol{y}} lpha_{oldsymbol{y},t} [ \mathcal{K}((oldsymbol{x}_t,oldsymbol{y}_t),(oldsymbol{x},oldsymbol{y}^*)) - \mathcal{K}((oldsymbol{x}_t,oldsymbol{y}),(oldsymbol{x},oldsymbol{y}^*))]$$

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

#### Kernels = Tractable Non-Linearity

- A linear model in a higher dimensional feature space is a non-linear model in the original space
- Computing a non-linear kernel is often better computationally than calculating the corresponding dot product in the high dimensional feature space
- Thus, kernels allow us to efficiently learn non-linear models by convex optimization



# Linear Learners in High Dimension



 $\begin{array}{rccc} \Re^2 & \longrightarrow & \Re^3 \\ (x_1, x_2) & \longmapsto & (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1 x_2, x_2^2) \end{array}$ 

# **Example: Polynomial Kernel**

which equals:

$$[(x_{t,1})^2, (x_{t,2})^2, \sqrt{2}x_{t,1}, \sqrt{2}x_{t,2}, \sqrt{2}x_{t,1}x_{t,2}, 1] + [(x_{s,1})^2, (x_{s,2})^2, \sqrt{2}x_{s,1}, \sqrt{2}x_{s,2}, \sqrt{2}x_{s,1}x_{s,2}, 1]$$

feature vector in high-dimensional space

feature vector in high-dimensional space

#### **Popular Kernels**

Polynomial kernel

$$K(x_t,x_s)=(\phi(x_t)\cdot\phi(x_s)+1)^d$$

 Gaussian radial basis kernel (infinite feature space representation!)

$$\mathcal{K}(x_t,x_s) = exp(rac{-||\phi(x_t)-\phi(x_s)||^2}{2\sigma})$$

- String kernels
- Tree kernels

# **Kernels Summary**

- Can turn a linear model into a non-linear model
- Kernels project feature space to higher dimensions
  - Sometimes exponentially larger
  - Sometimes an infinite space!
- Can "kernelize" algorithms to make them non-linear
- Convex optimization methods still applicable to learn parameters

# Kernels for Large Training Sets

- Exact kernel methods depend polynomially on the number of training examples - infeasible for large datasets
- Alternative: Explicit randomized feature map

[Rahimi and Recht 2007]

- Shallow neural network by random Fourier transformation:
  - Random weights from input to hidden units
  - Cosine as transfer function
  - Linear learning of weights from hidden to output units


# 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 (preferred!) fashion

What model to use?

- One example of a model: linear model
- Can kernelize/randomize these models to get non-linear models
- Convex optimization applicable for both types of model

# Further Reading

## Introductory Example:

J. Y. Lettvin, H. R. Maturana, W. S. McCulloch, and W. H. Pitts. 1959. What the frog's eye tells the frog's brain. Proc. Inst. Radio Engr., 47:1940–1951.

#### Naive Bayes:

Pedro Domingos and Michael Pazzani. 1997. On the optimality of the simple bayesian classifier under zero-one loss. <u>Machine</u> <u>Learning</u>, (29):103–130.

## Logistic Regression:

Bradley Efron. 1975.

The efficiency of logistic regression compared to normal discriminant analysis. Journal of the American Statistical Association, 70(352):892–898.

- Adam L. Berger, Vincent J. Della Pietra, and Stephen A. Della Pietra. 1996. A maximum entropy approach to natural language processing. <u>Computational</u> Linguistics, 22(1):39–71.
- Stefan Riezler, Detlef Prescher, Jonas Kuhn, and Mark Johnson. 2000.

Lexicalized Stochastic Modeling of Constraint-Based Grammars using Log-Linear Measures and EM Training. In <u>Proceedings of the 38th Annual Meeting of the</u> Association for Computational Linguistics (ACL'00), Hong Kong.

#### Perceptron:

Albert B.J. Novikoff. 1962.

On convergence proofs on perceptrons. Symposium on the Mathematical Theory of Automata, 12:615–622.

Yoav Freund and Robert E. Schapire. 1999. Large margin classification using the perceptron algorithm. <u>Journal of Machine</u> Learning Research, 37:277–296.

# Michael Collins. 2002. Discriminative training methods for hidden markov models: theory and experiments with perceptron algorithms. In <u>Proceedings of the conference on Empirical</u> <u>Methods in Natural Language Processing (EMNLP'02)</u>, Philadelphia, PA.

## ► SVM:

- Vladimir N. Vapnik. 1998.
  Statistical Learning Theory. Wiley.
- Olivier Chapelle. 2007.

Training a support vector machine in the primal. <u>Neural Computation</u>, 19(5):1155–1178.

- Ben Taskar, Carlos Guestrin, and Daphne Koller. 2003. Max-margin markov networks. In <u>Advances in Neural Information Processing</u> Systems 17 (NIPS'03), Vancouver, Canada.
- Ioannis Tsochantaridis, Thomas Hofmann, Thorsten Joachims, and Yasemin Altun. 2004.

Support vector machine learning for interdependent and structured output spaces. In Proceedings of the 21st International Conference on Machine Learning (ICML'04), Banff, Canada.

#### Kernels and Regularization:

Bernhard Schölkopf and Alexander J. Smola. 2002. Learning with Kernels. Support Vector Machines, Regularization, Optimization, and Beyond. The MIT Press.

Ali Rahimi and Ben Recht. 2007. Random features for large-scale kernel machines. In <u>Advances in Neural</u> Information Processing Systems (NIPS), Vancouver, B.C., Canada. Zhiyun Lu, Dong Guo, Alireza Bagheri Garakani, Kuan Liu, Avner May, Aurelien Bellet, Linxi Fan, Michael Collins, Brian Kingsbury, Michael Picheny, and Fei Sha. 2016.

A comparison between deep neural nets and kernel acoustic models for speech recognition. In IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP).

#### Convex and Non-Convex Optimization:

- Yurii Nesterov. 2004. Introductory lectures on convex optimization: A basic course. Springer.
- Stephen Boyd and Lieven Vandenberghe. 2004. Convex Optimization. Cambridge University Press.
- Dimitri P. Bertsekas and John N. Tsitsiklis. 1996. Neuro-Dynamic Programming. Athena Scientific.

#### Online/Stochastic Optimization:

- Herbert Robbins and Sutton Monro. 1951. A stochastic approximation method. <u>Annals of Mathematical Statistics</u>, 22(3):400–407.
- Boris T. Polyak. 1964.

Some methods of speeding up the convergence of iteration methods. <u>USSR</u> Computational Mathematics and Mathematical Physics, 4(5):1 – 17.

- Nicolò Cesa-Bianchi, Alex Conconi, and Claudio Gentile. 2004. On the generalization ablility of on-line learning algorithms. <u>IEEE Transactons on</u> Information Theory, 50(9):2050–2057.
- Nicolas LeRoux, Mark Schmidt, and Francis Bach. 2012. A stochastic gradient method with an exponential convergence rate for finite training sets. In <u>Advances in Neural Information Processing Systems (NIPS)</u>, Lake Tahoe, CA.
- Ilya Sutskever, James Martens, George E. Dahl, and Geoffrey E. Hinton. 2013. On the importance of initialization and momentum in deep learning. In <u>Proceedings</u> of International Conference on Machine Learning (ICML), pages 1139–1147.