# Introduction to Machine Learning: Linear Learners

Lisbon Machine Learning School, 2017

Stefan Riezler

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

# 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?

# Learning 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 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:
  ▶ p(+) =       , p(-) =
  ▶ p(convex = small|-) =       , p(convex = med|-) =       , p(convex = large|-) =
     p(speed = small|-) =       , p(speed = med|-) =       , p(speed = large|- ) =
     p(convex = small|+) =       , p(convex = med|+) =       , p(convex = large|+) =
     p(speed = small|+) =       , p(speed = med|+) =       , p(speed = large|+ ) =

# Learning 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:
  - $p(+) = 6/14$, $p(-) = 8/14$
  - $p(convex = small|-) = \quad$, $p(convex = med|-) = \quad$, $p(convex = large|-) =$
    $p(speed = small|-) = \quad$, $p(speed = med|-) = \quad$, $p(speed = large|- ) =$
    $p(convex = small|+) = \quad$, $p(convex = med|+) = \quad$, $p(convex = large|+) =$
    $p(speed = small|+) = \quad$, $p(speed = med|+) = \quad$, $p(speed = large|+ ) =$

# Learning 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:
  ▶ $p(+) = 6/14$, $p(-) = 8/14$
  ▶ $p(\text{convex} = \text{small}|-) = 6/8$, $p(\text{convex} = \text{med}|-) = 1/8$, $p(\text{convex} = \text{large}|-) = 1/8$
    $p(\text{speed} = \text{small}|-) = 4/8$, $p(\text{speed} = \text{med}|-) = 3/8$, $p(\text{speed} = \text{large}|- ) = 1/8$
    $p(\text{convex} = \text{small}|+) = 1/6$, $p(\text{convex} = \text{med}|+) = 2/6$, $p(\text{convex} = \text{large}|+) = 3/6$
    $p(\text{speed} = \text{small}|+) = 1/6$, $p(\text{speed} = \text{med}|+) = 1/6$, $p(\text{speed} = \text{large}|+ ) = 4/6$

# Learning 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:

  ▶ p(+) = 6/14, p(-) = 8/14
  ▶ 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

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

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

# Learning 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:
  ▶ p(+) = 6/14, p(-) = 8/14
  ▶ 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

▶ Predict unseen p(label = ?, convex = med, speed = med)
  ▶ p(-) · p(convex = med|-) · p(speed = med|-) = 8/14 · 1/8 · 3/8 = 0.027
  ▶ p(+) · p(convex = med|+) · p(speed = med|+) =

# Learning 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:

▶ $p(+) = 6/14$, $p(-) = 8/14$
▶ $p(\text{convex} = \text{small}|-) = 6/8$, $p(\text{convex} = \text{med}|-) = 1/8$, $p(\text{convex} = \text{large}|-) = 1/8$
  $p(\text{speed} = \text{small}|-) = 4/8$, $p(\text{speed} = \text{med}|-) = 3/8$, $p(\text{speed} = \text{large}|- ) = 1/8$
  $p(\text{convex} = \text{small}|+) = 1/6$, $p(\text{convex} = \text{med}|+) = 2/6$, $p(\text{convex} = \text{large}|+) = 3/6$
  $p(\text{speed} = \text{small}|+) = 1/6$, $p(\text{speed} = \text{med}|+) = 1/6$, $p(\text{speed} = \text{large}|+ ) = 4/6$

▶ Predict unseen $p(\text{label} = ?, \text{convex} = \text{med}, \text{speed} = \text{med})$

▶ $p(-) \cdot p(\text{convex} = \text{med}|-) \cdot p(\text{speed} = \text{med}|-) = 8/14 \cdot 1/8 \cdot 3/8 = 0.027$
▶ $p(+) \cdot p(\text{convex} = \text{med}|+) \cdot p(\text{speed} = \text{med}|+) = 6/14 \cdot 2/6 \cdot 1/6 = 0.024$

# Learning 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:
  ▶ p(+) = 6/14, p(-) = 8/14
  ▶ 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

▶ Predict unseen p(label = ?, convex = med, speed = med)
  ▶ p(-) · p(convex = med|-) · p(speed = med|-) = 8/14 · 1/8 · 3/8 = 0.027
  ▶ p(+) · p(convex = med|+) · p(speed = med|+) = 6/14 · 2/6 · 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 **convex optimization**

## Lecture Outline

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

# Inputs and Outputs

- Input: $x \in \mathcal{X}$
  - e.g., document or sentence with some words $x = w_1 \ldots w_n$
- Output: $y \in \mathcal{Y}$
  - e.g., document class, translation, parse tree

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

# Feature Representations

- Most NLP problems can be cast as multiclass classification where we assume a high-dimensional <span style="color:red">joint feature map</span> on input-output pairs $(\boldsymbol{x}, \boldsymbol{y})$
  - $\phi(\boldsymbol{x}, \boldsymbol{y}) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^m$

# Feature Representations

- Most NLP problems can be cast as multiclass classification where we assume a high-dimensional joint feature map on input-output pairs $(\boldsymbol{x}, \boldsymbol{y})$
  - $\phi(\boldsymbol{x}, \boldsymbol{y}) : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}^m$
- Common ranges:
  - categorical (e.g., counts): $\phi_i \in \{1, \dots, F_i\}$, $F_i \in \mathbb{N}^+$
  - binary (e.g., binning): $\phi \in \{0, 1\}^m$
  - continuous (e.g., word embeddings): $\phi \in \mathbb{R}^m$

# Feature Representations

- Most NLP problems can be cast as multiclass classification where we assume a high-dimensional joint feature map on input-output pairs $(\boldsymbol{x}, \boldsymbol{y})$
  - $\phi(\boldsymbol{x}, \boldsymbol{y}) : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}^m$
- Common ranges:
  - categorical (e.g., counts): $\phi_i \in \{1, \ldots, F_i\}$, $F_i \in \mathbb{N}^+$
  - binary (e.g., binning): $\phi \in \{0, 1\}^m$
  - continuous (e.g., word embeddings): $\phi \in \mathbb{R}^m$

- For any vector $\mathbf{v} \in \mathbb{R}^m$, let $\mathbf{v}_j$ be the $j^{th}$ 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

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

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

## Examples

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

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

What weight would it get?

## Examples

- $x$ is a source sentence and $y$ is translation

$$\phi_j(x, y) = \begin{cases} 1 & \text{if "y a-t-il" present in } x \\ & \text{and "are there" present in } y \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_k(x, y) = \begin{cases} 1 & \text{if "y a-t-il" present in } x \\ & \text{and "are there any" present in } y \\ 0 & \text{otherwise} \end{cases}$$

Which phrase indicator should be preferred?

## Examples



Note: Label $y$ includes sentence $x$

## Linear Models

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

$$
\begin{aligned}
f(\boldsymbol{x}; \boldsymbol{\omega}) &= \underset{\boldsymbol{y} \in \mathcal{Y}}{\arg\max} \; \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}, \boldsymbol{y}) \\
&= \underset{\boldsymbol{y} \in \mathcal{Y}}{\arg\max} \; \sum_{j=0}^{m} \boldsymbol{\omega}_j \times \boldsymbol{\phi}_j(\boldsymbol{x}, \boldsymbol{y})
\end{aligned}
$$

## Linear Models

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

$$
\begin{aligned}
f(\boldsymbol{x}; \boldsymbol{\omega}) &= \operatorname*{arg\,max}_{\boldsymbol{y} \in \mathcal{Y}} \ \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}, \boldsymbol{y}) \\
&= \operatorname*{arg\,max}_{\boldsymbol{y} \in \mathcal{Y}} \ \sum_{j=0}^{m} \boldsymbol{\omega}_j \times \boldsymbol{\phi}_j(\boldsymbol{x}, \boldsymbol{y})
\end{aligned}
$$

▶ Let $\boldsymbol{\omega} \in \mathbb{R}^m$ be a high dimensional weight vector

▶ Assume that $\boldsymbol{\omega}$ is known

　▶ **Multiclass Classification**: $\mathcal{Y} = \{0, 1, \ldots, N\}$

$$
\boldsymbol{y} = \operatorname*{arg\,max}_{\boldsymbol{y}' \in \mathcal{Y}} \ \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}, \boldsymbol{y}')
$$

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

# Linear Models for Binary Classification

- $\omega$ defines a linear decision boundary that divides space of instances in two classes
  - 2 dimensions: line
  - 3 dimensions: plane
  - $n$ dimensions: hyperplane of $n-1$ dimensions

## Multiclass Linear Model

Defines regions of space. Visualization difficult.



- ► $+$ are all points $(\boldsymbol{x}, \boldsymbol{y})$ where $+ = \arg\max_{\boldsymbol{y}} \ \boldsymbol{\omega} \cdot \phi(\boldsymbol{x}, \boldsymbol{y})$

## Convex Optimization for Supervised Learning

How to learn weight vector $\boldsymbol{\omega}$ in order to make decisions?

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

# Convex Optimization for Supervised Learning

How to learn weight vector $\boldsymbol{\omega}$ in order to make decisions?

- Input:
  - i.i.d. (independent and identically distributed) training examples $\mathcal{T} = \{(\boldsymbol{x}_t, \boldsymbol{y}_t)\}_{t=1}^{|\mathcal{T}|}$
  - feature representation $\phi$
- Output: $\boldsymbol{\omega}$ 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 $(\boldsymbol{x}, \boldsymbol{y})$
  - ▶ $\mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) \propto \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in \mathcal{T}} loss((\boldsymbol{x}, \boldsymbol{y}); \boldsymbol{\omega})$
  - ▶ *loss* is a function that measures some value correlated with errors of parameters $\boldsymbol{\omega}$ on instance $(\boldsymbol{x}, \boldsymbol{y})$

## Objective Functions

- Ideally we can decompose $\mathcal{L}$ by training pairs $(\boldsymbol{x}, \boldsymbol{y})$
  - $\mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) \propto \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in \mathcal{T}} loss((\boldsymbol{x}, \boldsymbol{y}); \boldsymbol{\omega})$
  - *loss* is a function that measures some value correlated with errors of parameters $\boldsymbol{\omega}$ on instance $(\boldsymbol{x}, \boldsymbol{y})$

- Example:
  - $\boldsymbol{y} \in \{1, -1\}$, $f(\boldsymbol{x}; \boldsymbol{\omega})$ is the prediction we make for $\boldsymbol{x}$ using $\boldsymbol{\omega}$
  - 0-1 loss function: $loss((\boldsymbol{x}, \boldsymbol{y}); \boldsymbol{\omega}) = \begin{cases} 0 & \text{if } f(\boldsymbol{x}; \boldsymbol{\omega}) = \boldsymbol{y}, \\ 1 & \text{else} \end{cases}$

## Convexity



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

$$f(\alpha\boldsymbol{x}+\beta\boldsymbol{y}) \leq \alpha f(\boldsymbol{x})+\beta f(\boldsymbol{y}) \text{ for all } \alpha,\beta \geq 0,\ \alpha+\beta = 1 \quad (1)$$

## Gradient



- Gradient of function f is vector of partial derivatives.
  $\nabla f(x) = \left( \frac{\partial}{\partial x_1} f(x), \frac{\partial}{\partial x_2} f(x), ..., \frac{\partial}{\partial x_n} f(x) \right)$
- Rate of increase of $f$ at point $x$ in each of the axis-parallel directions.

# Convex Optimization



(a)    (b)

▶ Optimization problem is defined as problem of finding a point
  that minimizes our objective function (maximization is
  minimization of $-f(x)$)

# Convex Optimization



(a) (b)

- ▶ Optimization problem is defined as problem of finding a point that minimizes our objective function (maximization is minimization of $-f(x)$)
- ▶ In order to find minimum, follow opposite direction of gradient
- ▶ For convex (or linear) functions, global minimum at point where $\nabla f(x) = 0$

# Naive Bayes

# Naive Bayes

▶ Probabilistic decision model:

$$\arg\max_{\boldsymbol{y}} P(\boldsymbol{y}|\boldsymbol{x}) \propto \arg\max_{\boldsymbol{y}} P(\boldsymbol{y})P(\boldsymbol{x}|\boldsymbol{y})$$

  ▶ Uses Bayes Rule:

$$P(\boldsymbol{y}|\boldsymbol{x}) = \frac{P(\boldsymbol{y})P(\boldsymbol{x}|\boldsymbol{y})}{P(\boldsymbol{x})} \text{ for fixed } \boldsymbol{x}$$

▶ Generative model since $P(\boldsymbol{y})P(\boldsymbol{x}|\boldsymbol{y}) = P(\boldsymbol{x}, \boldsymbol{y})$ 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 $P(x, y)$

- $P(x|y) = P(\phi(x)|y) = P(\phi_1(x), \ldots, \phi_m(x)|y)$

<div align="center">

Naive Bayes Assumption

*(conditional independence)*

$$P(\phi_1(x), \ldots, \phi_m(x)|y) = \prod_i P(\phi_i(x)|y)$$

</div>

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

## Naive Bayes – Learning

- Input: $\mathcal{T} = \{(\boldsymbol{x}_t, \boldsymbol{y}_t)\}_{t=1}^{|\mathcal{T}|}$

- Let $\phi_i(\boldsymbol{x}) \in \{1, \ldots, F_i\}$

- Parameters $\mathcal{P} = \{P(\boldsymbol{y}), P(\phi_i(\boldsymbol{x})|\boldsymbol{y})\}$

# 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)

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

## Naive Bayes – Learning

MLE has closed form solution

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

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

$$P(\boldsymbol{\phi}_i(\boldsymbol{x})|\boldsymbol{y}) = \frac{\sum_{t=1}^{|\mathcal{T}|} [\![\boldsymbol{\phi}_i(\boldsymbol{x}_t) = \boldsymbol{\phi}_i(\boldsymbol{x}) \text{ and } \boldsymbol{y}_t = \boldsymbol{y}]\!]}{\sum_{t=1}^{|\mathcal{T}|} [\![\boldsymbol{y}_t = \boldsymbol{y}]\!]}$$

where $[\![p]\!] = \begin{cases} 1 & \text{if } p \text{ is true,} \\ 0 & \text{otherwise.} \end{cases}$

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

## Deriving MLE

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

## Deriving MLE

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

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

## Deriving MLE

$$\mathcal{P} = \underset{P(\boldsymbol{y})}{\arg\max} \sum_{t=1}^{|\mathcal{T}|} \log P(\boldsymbol{y_t}) + \underset{P(\phi_i(\boldsymbol{x})|\boldsymbol{y})}{\arg\max} \sum_{t=1}^{|\mathcal{T}|} \sum_{i=1}^{m} \log P(\phi_i(\boldsymbol{x_t})|\boldsymbol{y_t})$$

## Deriving MLE

$$\mathcal{P} = \underset{P(\boldsymbol{y})}{\arg\max} \sum_{t=1}^{|\mathcal{T}|} \log P(\boldsymbol{y_t}) + \underset{P(\phi_i(\boldsymbol{x})|\boldsymbol{y})}{\arg\max} \sum_{t=1}^{|\mathcal{T}|} \sum_{i=1}^{m} \log P(\phi_i(\boldsymbol{x_t})|\boldsymbol{y_t})$$

Both optimizations are of the form

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

where $v$ is event in $\mathcal{T}$, either $(\boldsymbol{y_t} = \boldsymbol{y})$ or $(\phi_i(\boldsymbol{x_t}) = \phi_i(\boldsymbol{x}), \boldsymbol{y_t} = \boldsymbol{y})$

## Deriving MLE

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

Introduce Lagrangian multiplier $\lambda$, optimization becomes

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

## Deriving MLE

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

Introduce Lagrangian multiplier $\lambda$, optimization becomes

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

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

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

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

## Deriving MLE

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

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

$$P(\boldsymbol{\phi}_i(\boldsymbol{x})|\boldsymbol{y}) = \frac{\sum_{t=1}^{|\mathcal{T}|} [\![\boldsymbol{\phi}_i(\boldsymbol{x}_t) = \boldsymbol{\phi}_i(\boldsymbol{x}) \text{ and } \boldsymbol{y}_t = \boldsymbol{y}]\!]}{\sum_{t=1}^{|\mathcal{T}|} [\![\boldsymbol{y}_t = \boldsymbol{y}]\!]}$$

## Naive Bayes is a linear model

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

## Naive Bayes is a linear model

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

$$
\begin{aligned}
\arg\max_y \; P(y|\phi(x)) \quad \propto \quad & \arg\max_y \; P(\phi(x), y) = \arg\max_y \; P(y) \prod_{i=1}^{m} P(\phi_i(x)|y) \\
= \quad & \arg\max_y \; \log P(y) + \sum_{i=1}^{m} \log P(\phi_i(x)|y) \\
= \quad & \arg\max_y \; \omega_y + \sum_{i=1}^{m} \omega_{\phi_i(x),y} \\
= \quad & \arg\max_y \; \sum_{y'} \omega_y \psi_{y'}(y) + \sum_{i=1}^{m} \sum_{j=1}^{F_i} \omega_{\phi_i(x),y} \psi_{i,j}(x)
\end{aligned}
$$

where $\psi_{i,j}(x) = [\![\phi_i(x) = j]\!]$, $\psi_{y'}(y) = [\![y = y']\!]$

## Discriminative versus Generative Models

▶ Generative models attempt to model inputs and outputs
  ▸ e.g., Naive Bayes = MLE of joint distribution $P(\boldsymbol{x}, \boldsymbol{y})$
  ▸ 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 $P(\boldsymbol{y}|\boldsymbol{x})$ (or something related)
  ▸ Logistic Regression – MLE of $P(\boldsymbol{y}|\boldsymbol{x})$
  ▸ Perceptron and SVMs – minimize classification error

▶ Generative and discriminative models use $P(\boldsymbol{y}|\boldsymbol{x})$ for prediction

▶ Differ only on what distribution they use to set $\boldsymbol{\omega}$

# Logistic Regression

## Logistic Regression

Define a conditional probability:

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

Note: still a linear model

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

## Logistic Regression

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

▶ Q: How do we learn weights $\boldsymbol{\omega}$

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

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

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

## Logistic Regression

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

$$\boldsymbol{\omega} = \arg\max_{\boldsymbol{\omega}} \sum_{t=1}^{|\mathcal{T}|} \log P(\boldsymbol{y}_t|\boldsymbol{x}_t) \; (*)$$

- 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 $\mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) = \sum_{t=1}^{|\mathcal{T}|} \log \left( e^{\boldsymbol{\omega} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}_t)} / Z_{\boldsymbol{x}} \right)$
- Want to find $\arg \max_{\boldsymbol{\omega}} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega})$
  - Set $\boldsymbol{\omega}^0 = O^m$
  - Iterate until convergence

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

- $\alpha > 0$ is a step size / learning rate
- $\nabla \mathcal{L}(\mathcal{T}; \boldsymbol{\omega})$ is gradient of $\mathcal{L}$ w.r.t. $\boldsymbol{\omega}$
  - A gradient is all partial derivatives over variables $w_i$
  - i.e., $\nabla \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) = (\frac{\partial}{\partial \boldsymbol{\omega}_0} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}), \frac{\partial}{\partial \boldsymbol{\omega}_1} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}), \ldots, \frac{\partial}{\partial \boldsymbol{\omega}_m} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}))$
- Gradient ascent will always find $\boldsymbol{\omega}$ to maximize $\mathcal{L}$

# Gradient Descent

- Let $\mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) = -\sum_{t=1}^{|\mathcal{T}|} \log \left( e^{\boldsymbol{\omega} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}_t)} / Z_{\boldsymbol{x}} \right)$
- Want to find $\arg\min_{\boldsymbol{\omega}} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega})$
  - Set $\boldsymbol{\omega}^0 = O^m$
  - Iterate until convergence

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

- $\alpha > 0$ is step size / learning rate
- $\nabla \mathcal{L}(\mathcal{T}; \boldsymbol{\omega})$ is gradient of $\mathcal{L}$ w.r.t. $\boldsymbol{\omega}$
  - A gradient is all partial derivatives over variables $w_i$
  - i.e., $\nabla \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) = (\frac{\partial}{\partial \boldsymbol{\omega}_0} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}), \frac{\partial}{\partial \boldsymbol{\omega}_1} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}), \ldots, \frac{\partial}{\partial \boldsymbol{\omega}_m} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}))$

- Gradient descent will always find $\boldsymbol{\omega}$ to minimize $\mathcal{L}$

## The partial derivatives

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

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

## **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

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

## The partial derivatives (1)

$$
\begin{aligned}
\frac{\partial}{\partial \boldsymbol{\omega}_i} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) &= \frac{\partial}{\partial \boldsymbol{\omega}_i} \sum_t \log \frac{e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)}}{Z_{\boldsymbol{x}_t}} \\
&= \sum_t \frac{\partial}{\partial \boldsymbol{\omega}_i} \log \frac{e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)}}{Z_{\boldsymbol{x}_t}} \\
&= \sum_t \left( \frac{Z_{\boldsymbol{x}_t}}{e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)}} \right) \left( \frac{\partial}{\partial \boldsymbol{\omega}_i} \frac{e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)}}{Z_{\boldsymbol{x}_t}} \right)
\end{aligned}
$$

## The partial derivatives

Now, $\dfrac{\partial}{\partial \omega_i} \dfrac{e^{\sum_j \omega_j \times \phi_j(x_t, y_t)}}{Z_{x_t}} =$

## The partial derivatives (2)

Now,

$$
\begin{aligned}
\frac{\partial}{\partial \boldsymbol{\omega}_i} \frac{e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)}}{Z_{\boldsymbol{x}_t}} &= \frac{Z_{\boldsymbol{x}_t} \frac{\partial}{\partial \boldsymbol{\omega}_i} e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)} - e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)} \frac{\partial}{\partial \boldsymbol{\omega}_i} Z_{\boldsymbol{x}_t}}{Z_{\boldsymbol{x}_t}^2} \\
&= \frac{Z_{\boldsymbol{x}_t} e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)} \phi_i(\boldsymbol{x}_t, \boldsymbol{y}_t) - e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)} \frac{\partial}{\partial \boldsymbol{\omega}_i} Z_{\boldsymbol{x}_t}}{Z_{\boldsymbol{x}_t}^2} \\
&= \frac{e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)}}{Z_{\boldsymbol{x}_t}^2} (Z_{\boldsymbol{x}_t} \phi_i(\boldsymbol{x}_t, \boldsymbol{y}_t) - \frac{\partial}{\partial \boldsymbol{\omega}_i} Z_{\boldsymbol{x}_t}) \\
&= \frac{e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)}}{Z_{\boldsymbol{x}_t}^2} (Z_{\boldsymbol{x}_t} \phi_i(\boldsymbol{x}_t, \boldsymbol{y}_t) \\
&\qquad\qquad - \sum_{\boldsymbol{y}' \in \mathcal{Y}} e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}')} \phi_i(\boldsymbol{x}_t, \boldsymbol{y}'))
\end{aligned}
$$

because

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

# The partial derivatives

## The partial derivatives (3)

From (2),

$$
\frac{\partial}{\partial \boldsymbol{\omega}_i} \frac{e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)}}{Z_{\boldsymbol{x}_t}} = \frac{e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}_t)}}{Z_{\boldsymbol{x}_t}^2} (Z_{\boldsymbol{x}_t} \phi_i(\boldsymbol{x}_t, \boldsymbol{y}_t) \\
- \sum_{\boldsymbol{y}' \in \mathcal{Y}} e^{\sum_j \boldsymbol{\omega}_j \times \phi_j(\boldsymbol{x}_t, \boldsymbol{y}')} \phi_i(\boldsymbol{x}_t, \boldsymbol{y}'))
$$

Sub this in (1),

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

## **FINALLY!!!**

▶ After all that,

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

▶ And the gradient is:

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

▶ So we can now use gradient ascent to find $\boldsymbol{\omega}$!!

## Logistic Regression Summary

▶ Define conditional probability

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

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

$$\boldsymbol{\omega} = \arg\max_{\boldsymbol{\omega}} \sum_{t} \log P(\boldsymbol{y_t}|\boldsymbol{x_t})$$

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

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

## Logistic Regression = Maximum Entropy

- ▶ Well-known equivalence
- ▶ 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
    - ▶ Partial derivative in logistic regression

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

    - ▶ 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 Learning Algorithm

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 $\boldsymbol{\omega}^i$

## 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:

$$\mathbf{u} \cdot \phi(x_t, y_t) - \mathbf{u} \cdot \phi(x_t, y') \geq \gamma \tag{2}$$

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

- Thus, after a finite number of training iterations, the error on the training set will converge to zero

- **Let's prove it!**

## Perceptron Learning Algorithm

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:

$\boldsymbol{\omega}^{(k-1)}$ are weights before $k^{th}$ error

Suppose $k^{th}$ error made at $(\boldsymbol{x}_t, \boldsymbol{y}_t)$

$\boldsymbol{y}' = \arg\max_{\boldsymbol{y}'} \boldsymbol{\omega}^{(k-1)} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}')$

$\boldsymbol{y}' \neq \boldsymbol{y}_t$

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

## Perceptron Learning Algorithm

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:

$\boldsymbol{\omega}^{(k-1)}$ are weights before $k^{th}$ error

Suppose $k^{th}$ error made at $(\boldsymbol{x}_t, \boldsymbol{y}_t)$

$\boldsymbol{y}' = \arg\max_{\boldsymbol{y}'} \boldsymbol{\omega}^{(k-1)} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}')$

$\boldsymbol{y}' \neq \boldsymbol{y}_t$

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

$\mathbf{u} \cdot \boldsymbol{\omega}^{(k)} = \mathbf{u} \cdot \boldsymbol{\omega}^{(k-1)} + \mathbf{u} \cdot (\boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}')) \geq \mathbf{u} \cdot \boldsymbol{\omega}^{(k-1)} + \gamma$, by (2)

Since $\boldsymbol{\omega}^{(0)} = 0$ and $\mathbf{u} \cdot \boldsymbol{\omega}^{(0)} = 0$, for all $k$: $\mathbf{u} \cdot \boldsymbol{\omega}^{(k)} \geq k\gamma$, by induction on $k$

Since $\mathbf{u} \cdot \boldsymbol{\omega}^{(k)} \leq ||\mathbf{u}|| \times ||\boldsymbol{\omega}^{(k)}||$, by the law of cosines, and $||\mathbf{u}|| = 1$, then
$||\boldsymbol{\omega}^{(k)}|| \geq k\gamma$

## Perceptron Learning Algorithm

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 $\boldsymbol{y}' = \arg\max_{\boldsymbol{y}'} \boldsymbol{\omega}^{(i)} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}')$
5.     if $\boldsymbol{y}' \neq \boldsymbol{y}_t$
6.      $\boldsymbol{\omega}^{(i+1)} = \boldsymbol{\omega}^{(i)} + \phi(\boldsymbol{x}_t, \boldsymbol{y}_t) - \phi(\boldsymbol{x}_t, \boldsymbol{y}')$
7.      $i = i + 1$
8.   return $\boldsymbol{\omega}^i$

▶ Lower bound:

$\boldsymbol{\omega}^{(k-1)}$ are weights before $k^{th}$ error

Suppose $k^{th}$ error made at $(\boldsymbol{x}_t, \boldsymbol{y}_t)$

$\boldsymbol{y}' = \arg\max_{\boldsymbol{y}'} \boldsymbol{\omega}^{(k-1)} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}')$

$\boldsymbol{y}' \neq \boldsymbol{y}_t$

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

$\mathbf{u} \cdot \boldsymbol{\omega}^{(k)} = \mathbf{u} \cdot \boldsymbol{\omega}^{(k-1)} + \mathbf{u} \cdot (\phi(\boldsymbol{x}_t, \boldsymbol{y}_t) - \phi(\boldsymbol{x}_t, \boldsymbol{y}')) \geq \mathbf{u} \cdot \boldsymbol{\omega}^{(k-1)} + \gamma$, by (2)

Since $\boldsymbol{\omega}^{(0)} = 0$ and $\mathbf{u} \cdot \boldsymbol{\omega}^{(0)} = 0$, for all $k$: $\mathbf{u} \cdot \boldsymbol{\omega}^{(k)} \geq k\gamma$, by induction on $k$

Since $\mathbf{u} \cdot \boldsymbol{\omega}^{(k)} \leq ||\mathbf{u}|| \times ||\boldsymbol{\omega}^{(k)}||$, by the law of cosines, and $||\mathbf{u}|| = 1$, then

$||\boldsymbol{\omega}^{(k)}|| \geq k\gamma$

▶ Upper bound:

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

$$||\boldsymbol{\omega}^{(k)}||^2 \ \leq \ ||\boldsymbol{\omega}^{(k-1)}||^2 + R^2, \text{ since } R \geq ||\phi(\boldsymbol{x}_t, \boldsymbol{y}_t) - \phi(\boldsymbol{x}_t, \boldsymbol{y}')||$$

$$\text{and } \boldsymbol{\omega}^{(k-1)} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}_t) - \boldsymbol{\omega}^{(k-1)} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}') \leq 0$$

$$\leq \ kR^2 \text{ for all } k, \text{ by induction on } k$$

# Perceptron Learning Algorithm

- We have just shown that $||\boldsymbol{\omega}^{(k)}|| \geq k\gamma$ and $||\boldsymbol{\omega}^{(k)}||^2 \leq kR^2$

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

- and solving for $k$
$$k \leq \frac{R^2}{\gamma^2}$$

- Therefore the number of errors is bounded!

## Perceptron Summary

- Learns parameters of a linear model by minimizing error
- Guaranteed to find a $\boldsymbol{\omega}$ in a finite amount of time
- Perceptron is an example of an <span style="color:red">Online Learning Algorithm</span>
  - $\boldsymbol{\omega}$ is updated based on a single training instance in isolation

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

# **Averaged Perceptron**

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.       else
6.         $\boldsymbol{\omega}^{(i+1)} = \boldsymbol{\omega}^{(i)}$
7.       $i = i + 1$
8.   return $\left(\sum_i \boldsymbol{\omega}^{(i)}\right) / (N \times T)$

# Margin

Training

Testing



Denote the
value of the
margin by $\gamma$

## Maximizing Margin

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

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

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

# Maximizing Margin

- Intuitively maximizing margin makes sense
- By cross-validation, the generalization error on unseen test data can be shown to be 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 $\gamma > 0$

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

such that:

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

$$\forall (\boldsymbol{x}_t, \boldsymbol{y}_t) \in \mathcal{T}$$

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

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

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

## Max Margin = Min Norm

Let $\gamma > 0$

**Max Margin**:

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

such that:

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

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

# Max Margin = Min Norm

Let $\gamma > 0$

**Max Margin**:

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

such that:

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

$$\forall (\boldsymbol{x}_t, \boldsymbol{y}_t) \in \mathcal{T}$$

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

Change variables: $\mathbf{u} = \dfrac{\boldsymbol{\omega}}{\gamma}$

$||\boldsymbol{\omega}|| = 1$ iff $||\mathbf{u}|| = 1/\gamma$,
then $\gamma = 1/||\mathbf{u}||$

## Max Margin $=$ Min Norm

Let $\gamma > 0$

**Max Margin**:

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

such that:

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

$$\forall(\boldsymbol{x}_t,\boldsymbol{y}_t) \in \mathcal{T}$$

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

Change variables: $\mathbf{u} = \dfrac{\boldsymbol{\omega}}{\gamma}$

$||\boldsymbol{\omega}|| = 1$ iff $||\mathbf{u}|| = 1/\gamma$,
then $\gamma = 1/||\mathbf{u}||$

**Min Norm (step 1)**:

$$\max_{\mathbf{u}} \quad \frac{1}{||\mathbf{u}||}$$

such that:

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

$$\forall(\boldsymbol{x}_t,\boldsymbol{y}_t) \in \mathcal{T}$$

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

## Max Margin = Min Norm

Let $\gamma > 0$

**Max Margin**:

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

such that:

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

$$\forall(\boldsymbol{x}_t,\boldsymbol{y}_t) \in \mathcal{T}$$

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

Change variables: $\mathbf{u} = \dfrac{\boldsymbol{\omega}}{\gamma}$

$||\boldsymbol{\omega}|| = 1$ iff $||\mathbf{u}|| = 1/\gamma$,
then $\gamma = 1/||\mathbf{u}||$

**Min Norm (step 1)**:

$$\min_{\mathbf{u}} ||\mathbf{u}||$$

such that:

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

$$\forall(\boldsymbol{x}_t,\boldsymbol{y}_t) \in \mathcal{T}$$

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

# Max Margin = Min Norm

Let $\gamma > 0$

**Max Margin**:

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

such that:

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

$$\forall(\boldsymbol{x}_t,\boldsymbol{y}_t) \in \mathcal{T}$$

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

Change variables: $\mathbf{u} = \dfrac{\boldsymbol{\omega}}{\gamma}$

$||\boldsymbol{\omega}|| = 1$ iff $||\mathbf{u}|| = 1/\gamma$,
then $\gamma = 1/||\mathbf{u}||$

**Min Norm (step 2)**:

$$\min_{\mathbf{u}} ||\mathbf{u}||$$

such that:

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

$$\forall(\boldsymbol{x}_t,\boldsymbol{y}_t) \in \mathcal{T}$$

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

## Max Margin $=$ Min Norm

Let $\gamma > 0$

**Max Margin**:

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

such that:

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

$$\forall(\boldsymbol{x}_t,\boldsymbol{y}_t) \in \mathcal{T}$$

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

Change variables: $\mathbf{u} = \dfrac{\boldsymbol{\omega}}{\gamma}$

$||\boldsymbol{\omega}|| = 1$ iff $||\mathbf{u}|| = 1/\gamma$,
then $\gamma = 1/||\mathbf{u}||$

**Min Norm (step 2)**:

$$\min_{\mathbf{u}} ||\mathbf{u}||$$

such that:

$$\mathbf{u}\cdot\phi(\boldsymbol{x}_t,\boldsymbol{y}_t)-\mathbf{u}\cdot\phi(\boldsymbol{x}_t,\boldsymbol{y}') \geq 1$$

$$\forall(\boldsymbol{x}_t,\boldsymbol{y}_t) \in \mathcal{T}$$

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

## Max Margin $=$ Min Norm

Let $\gamma > 0$

**Max Margin**:

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

such that:

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

$$\forall (\boldsymbol{x}_t, \boldsymbol{y}_t) \in \mathcal{T}$$

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

Change variables: $\mathbf{u} = \dfrac{\boldsymbol{\omega}}{\gamma}$

$||\boldsymbol{\omega}|| = 1$ iff $||\mathbf{u}|| = 1/\gamma$,
then $\gamma = 1/||\mathbf{u}||$

**Min Norm (step 3)**:

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

such that:

$$\mathbf{u} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \mathbf{u} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}') \geq 1$$

$$\forall (\boldsymbol{x}_t, \boldsymbol{y}_t) \in \mathcal{T}$$

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

# Max Margin = Min Norm

Let $\gamma > 0$

**Max Margin**:

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

such that:

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

$$\forall (\boldsymbol{x}_t, \boldsymbol{y}_t) \in \mathcal{T}$$

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

**Min Norm**:

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

such that:

$$\mathbf{u} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}_t) - \mathbf{u} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}') \geq 1$$

$$\forall (\boldsymbol{x}_t, \boldsymbol{y}_t) \in \mathcal{T}$$

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

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

## Support Vector Machines

▶ **Constrained Optimization Problem**

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

such that:

$$\boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}') \geq 1$$

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

▶ **Support Vectors:** Examples where

$$\boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}') = 1$$

for training instance $(\boldsymbol{x}_t, \boldsymbol{y}_t) \in \mathcal{T}$ and all $\boldsymbol{y}' \in \bar{\mathcal{Y}}_t$

# Support Vector Machines

- **What if data is not separable?**

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

such that:

$$\boldsymbol{\omega} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}_t) - \boldsymbol{\omega} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}') \geq 1 - \xi_t \text{ and } \xi_t \geq 0$$

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

- $\xi_t$: slack variable representing amount of constraint violation
- If data is separable, optimal solution has $\xi_i = 0$, $\forall i$
  C balances focus on margin and on error

# Support Vector Machines

- **What if data is not separable?**

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

such that:

$$\boldsymbol{\omega} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}_t) - \boldsymbol{\omega} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}') \geq 1 - \xi_t \text{ and } \xi_t \geq 0$$

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

- $\xi_t$: 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 ($C > \frac{1}{2}$)

# Support Vector Machines

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

such that:

$$\boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}') \geq 1 - \xi_t$$
$$\text{where } \xi_t \geq 0 \text{ and } \forall (\boldsymbol{x}_t, \boldsymbol{y}_t) \in \mathcal{T} \text{ and } \boldsymbol{y}' \in \bar{\mathcal{Y}}_t$$

- Computing the dual form results in a quadratic programming problem – a well-known convex optimization problem
- Can we have representation of this objective that allows more direct optimization?

## Support Vector Machines

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

such that:

$$\boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) - \max_{\boldsymbol{y}' \neq \boldsymbol{y}_t} \ \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}') \geq 1 - \xi_t$$

## Support Vector Machines

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

such that:

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

## Support Vector Machines

$$\boldsymbol{\omega} = \underset{\boldsymbol{\omega},\xi}{\arg\min} \ \frac{\lambda}{2}||\boldsymbol{\omega}||^2 + \sum_{t=1}^{|\mathcal{T}|} \xi_t \qquad \lambda = \frac{1}{C}$$

such that:

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

## Support Vector Machines

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

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

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

## Support Vector Machines

$$\boldsymbol{\omega} = \arg\min_{\boldsymbol{\omega}, \xi} \; \frac{\lambda}{2}||\boldsymbol{\omega}||^2 + \sum_{t=1}^{|\mathcal{T}|} \xi_t \; \text{s.t.} \; \xi_t \geq 1 + \max_{\boldsymbol{y}' \neq \boldsymbol{y}_t} \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}') - \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t)$$

<span style="color:red">Hinge loss</span>

$$\boldsymbol{\omega} = \arg\min_{\boldsymbol{\omega}} \; \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) = \arg\min_{\boldsymbol{\omega}} \; \sum_{t=1}^{|\mathcal{T}|} loss((\boldsymbol{x}_t, \boldsymbol{y}_t); \boldsymbol{\omega}) \; + \; \frac{\lambda}{2}||\boldsymbol{\omega}||^2$$

$$= \arg\min_{\boldsymbol{\omega}} \left( \sum_{t=1}^{|\mathcal{T}|} \max \left(0, 1 + \max_{\boldsymbol{y}' \neq \boldsymbol{y}_t} \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}') - \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t)\right) \right) + \frac{\lambda}{2}||\boldsymbol{\omega}||^2$$

- Hinge loss allows **unconstrained optimization** (later!)

# Summary
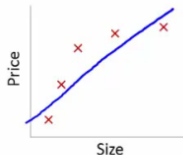
## What we have covered

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

## What is next

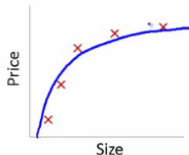- ▶ Regularization
- ▶ Online learning
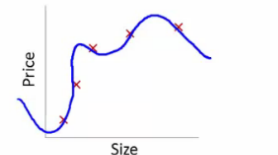- ▶ Non-linear models

# Regularization

# Fit of a Model



|  |  |  |
| --- | --- | --- |
| $\theta_0 + \theta_1 x$ | $\theta_0 + \theta_1 x + \theta_2 x^2$ | $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$ |
| High bias (underfit) | "Just right" | 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

# Fit of a Model



Degree 15 polynomial

## **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, e.g., adjacent word identities: 'He wants to X his death' $\rightarrow$ X=verb

## Regularization

▶ In practice, we regularize models to prevent overfitting

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

▶ Where $\mathcal{R}(\boldsymbol{\omega})$ is the regularization function
▶ $\lambda$ controls how much to regularize
▶ Most common regularizer
  ▹ L2: $\mathcal{R}(\boldsymbol{\omega}) \propto \|\boldsymbol{\omega}\|_2 = \|\boldsymbol{\omega}\| = \sqrt{\sum_i \omega_i^2}$ – smaller weights desired

## Logistic Regression with L2 Regularization

- Perhaps most common learner in NLP

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

- What are the new partial derivatives?

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

- 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!

$$\boldsymbol{\omega} = \arg\min_{\boldsymbol{\omega}} \ \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega})$$

## Support Vector Machines

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

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

# Support Vector Machines

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

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

# Support Vector Machines

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

$$
\begin{aligned}
\boldsymbol{\omega} &= \underset{\boldsymbol{\omega}}{\arg\min} \ \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega}) \\
&= \underset{\boldsymbol{\omega}}{\arg\min} \ \sum_{t=1}^{|\mathcal{T}|} \textit{loss}((\boldsymbol{x}_t, \boldsymbol{y}_t); \boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega}) \\
&= \underset{\boldsymbol{\omega}}{\arg\min} \ \sum_{t=1}^{|\mathcal{T}|} \max\left(0, 1 + \max_{\boldsymbol{y} \neq \boldsymbol{y}_t} \boldsymbol{\omega} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}) - \boldsymbol{\omega} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}_t)\right) + \lambda \mathcal{R}(\boldsymbol{\omega}) \\
&= \underset{\boldsymbol{\omega}}{\arg\min} \ \sum_{t=1}^{|\mathcal{T}|} \max\left(0, 1 + \max_{\boldsymbol{y} \neq \boldsymbol{y}_t} \boldsymbol{\omega} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}) - \boldsymbol{\omega} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}_t)\right) + \frac{\lambda}{2} \|\boldsymbol{\omega}\|^2
\end{aligned}
$$

## SVMs vs. Logistic Regression

$$\begin{aligned}
\boldsymbol{\omega} &= \underset{\boldsymbol{\omega}}{\arg\min} \; \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega}) \\
&= \underset{\boldsymbol{\omega}}{\arg\min} \; \sum_{t=1}^{|\mathcal{T}|} loss((\boldsymbol{x_t}, \boldsymbol{y_t}); \boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega})
\end{aligned}$$

# SVMs vs. Logistic Regression

$$\omega = \arg\min_{\omega} \mathcal{L}(\mathcal{T}; \omega) + \lambda\mathcal{R}(\omega)$$

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

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

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

## SVMs vs. Logistic Regression

$$\boldsymbol{\omega} = \arg\min_{\boldsymbol{\omega}} \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}) + \lambda \mathcal{R}(\boldsymbol{\omega})$$

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

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

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

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

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

# Summary: Loss Functions

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

# Online Learning

## Online vs. Batch Learning

Batch($\mathcal{T}$);

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

Online($\mathcal{T}$);

- ▶ for 1 ... N
    - ▶ for $(\boldsymbol{x}_t, \boldsymbol{y}_t) \in \mathcal{T}$
        - ▶ $\boldsymbol{\omega} \leftarrow$ update($(\boldsymbol{x}_t, \boldsymbol{y}_t); \boldsymbol{\omega}$)
    - ▶ end for
- ▶ end for
- ▶ return $\boldsymbol{\omega}$

E.g., SVMs, logistic regression, Naive Bayes

E.g., Perceptron
$$\boldsymbol{\omega} = \boldsymbol{\omega} + \phi(\boldsymbol{x}_t, \boldsymbol{y}_t) - \phi(\boldsymbol{x}_t, \boldsymbol{y})$$

## 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

$$
\begin{aligned}
\boldsymbol{\omega}^i &= \boldsymbol{\omega}^{i-1} - \alpha \nabla \mathcal{L}(\mathcal{T}; \boldsymbol{\omega}^{i-1}) \\
&= \boldsymbol{\omega}^{i-1} - \sum_{t=1}^{|\mathcal{T}|} \alpha \nabla loss((\boldsymbol{x}_t, \boldsymbol{y}_t); \boldsymbol{\omega}^{i-1})
\end{aligned}
$$

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

# Stochastic Gradient Descent

- Stochastic Gradient Descent (SGD)
    - Approximate batch gradient $\nabla \mathcal{L}(\mathcal{T}; \boldsymbol{\omega})$ with stochastic gradient $\nabla loss((\boldsymbol{x_t}, \boldsymbol{y_t}); \boldsymbol{\omega})$

- 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
        - sample $(\boldsymbol{x_t}, \boldsymbol{y_t}) \in \mathcal{T}$       // "stochastic"
        - $\boldsymbol{\omega}^i = \boldsymbol{\omega}^{i-1} - \alpha \nabla loss((\boldsymbol{x_t}, \boldsymbol{y_t}); \boldsymbol{\omega}^{i-1})$
    - return $\boldsymbol{\omega}$

## Online Logistic Regression

- Stochastic Gradient Descent (SGD)
- $loss((\boldsymbol{x}_t, \boldsymbol{y}_t); \boldsymbol{\omega}) = $ log-loss
- $\nabla loss((\boldsymbol{x}_t, \boldsymbol{y}_t); \boldsymbol{\omega}) = \nabla \left( -\log \left( e^{\boldsymbol{\omega} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y}_t)} / Z_{\boldsymbol{x}_t} \right) \right)$
- From logistic regression section:

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

- Plus regularization term (if part of model)

## Online SVMs

- Stochastic Gradient Descent (SGD)
- $loss((\boldsymbol{x}_t, \boldsymbol{y}_t); \boldsymbol{\omega}) = $ hinge-loss

$$\bigtriangledown loss((\boldsymbol{x}_t, \boldsymbol{y}_t); \boldsymbol{\omega}) = \bigtriangledown \left( \max \left(0, 1 + \max_{\boldsymbol{y} \neq \boldsymbol{y}_t} \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}) - \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) \right) \right)$$

- Subgradient is:

$$\bigtriangledown \left( \max \left(0, 1 + \max_{\boldsymbol{y} \neq \boldsymbol{y}_t} \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}) - \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t) \right) \right)$$

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

- Plus regularization term (L2 norm for SVMs):

$$\bigtriangledown \frac{\lambda}{2} ||\boldsymbol{\omega}||^2 = \lambda \boldsymbol{\omega}$$

# Perceptron and Hinge-Loss

SVM subgradient update looks like perceptron update

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

Perceptron

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

Perceptron = SGD optimization of no-margin hinge-loss (without regularization):

$$\max \left(0, 1 + \max_{\boldsymbol{y} \neq \boldsymbol{y}_t} \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}) - \boldsymbol{\omega} \cdot \boldsymbol{\phi}(\boldsymbol{x}_t, \boldsymbol{y}_t)\right)$$

## **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
  - ▶ Higher cost of averaging gradients over $\mathcal{T}$ 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 $0 \leq \mu < 1$ [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
  $$\boldsymbol{v}^i = \boldsymbol{\omega}^i - \boldsymbol{\omega}^{i-1}$$

  - 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
- ▶ 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 $1 - \delta$:

$$\mathbb{E}[loss((\boldsymbol{x}, \boldsymbol{y}); \bar{\boldsymbol{\omega}})] \leq M_k + \sqrt{\frac{2}{k} \ln \frac{1}{\delta}}$$

  - ▸ = 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
- ▶ Objective functions and loss functions
- ▶ Convex Optimization
- ▶ 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

# Kernel Machines = Convex Optimization for Non-Linear Models

- ▶ Projecting a linear model into a higher dimensional feature space can correspond to a non-linear model in the original space and make non-separable problems separable
- ▶ For classifiers based on similarity functions (= kernels), computing a non-linear kernel is often more efficient 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

## Monomial Features and Polynomial Kernels

- ▶ Monomial features $= d^{th}$ order products of entries $x_j$ of $\boldsymbol{x}$ s.t.
  $x_{j_1} * x_{j_2} * \cdots * x_{j_d}$ for $j_1, \ldots, j_d \in \{1 \ldots n\}$
- ▶ Ordered monomial feature map: $\phi : \mathbb{R}^2 \to \mathbb{R}^4$ s.t.
  $(x_1, x_2) \mapsto (x_1^2, x_2^2, x_1 x_2, x_2 x_1)$
- ▶ Computation of kernel from feature map:

$$
\begin{aligned}
\phi(\boldsymbol{x}) \cdot \phi(\boldsymbol{x}') &= \sum_{i=1}^{4} \phi_i(\boldsymbol{x}) \phi_i(\boldsymbol{x}') \text{ (Def. dot product)} \\
&= x_1^2 {x'_1}^2 + x_2^2 {x'_2}^2 + x_1 x_2 x'_1 x'_2 + x_2 x_1 x'_2 x'_1 \text{ (Def. } \phi) \\
&= x_1^2 {x'_1}^2 + x_2^2 {x'_2}^2 + 2 x_1 x_2 x'_1 x'_2 \\
&= \left( x_1 x'_1 + x_2 x'_2 \right)^2
\end{aligned}
$$

- ▶ Direct application of kernel: $\phi(\boldsymbol{x}) \cdot \phi(\boldsymbol{x}') = (\boldsymbol{x} \cdot \boldsymbol{x}')^2$

## Direct Application of Kernel

- Let $C_d$ be a map from $x \in \mathbb{R}^m$ to vectors $C_d(x)$ of all $d^{th}$-degree ordered products of entries of $x$.
  Then the corresponding kernel computing the dot product of vectors mapped by $C_d$ is:
  $K(x, x') = C_d(x) \cdot C_d(x') = (x \cdot x')^d$

- Alternative feature map satisfying this definition = unordered monomial: $\phi_2 : \mathbb{R}^2 \to \mathbb{R}^3$ s.t. $(x_1, x_2) \mapsto (x_1^2, x_2^2, \sqrt{2}x_1 x_2)$

## Non-Linear Feature Map



$$\phi_2 : \mathbb{R}^2 \to \mathbb{R}^3 \text{ s.t. } (x_1, x_2) \mapsto (z_1, z_2, z_3) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

# Kernel Definition

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

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

- 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:

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

- Positive definite: positivity on diagonal

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

## Mercer's Theorem

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

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

- This means that instead of mapping input data via non-lineear feature map $\phi$ and then computing dot product, we can apply kernels directly *without even knowing about $\phi$*!

## Kernel Trick

- ▶ 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

## Kernel Trick

▶ 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

▶ Note: Since our features are over pairs $(x, y)$, we will write kernels over pairs

$$K((x_t, y_t), (x_r, y_r)) = \phi(x_t, y_t) \cdot \phi(x_r, y_r)$$

▶ Let's do it for the Perceptron!

## Kernel Trick – Perceptron Algorithm

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 $\boldsymbol{y} = \arg\max_{\boldsymbol{y}} \boldsymbol{\omega}^{(i)} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y})$
5.         if $\boldsymbol{y} \neq \boldsymbol{y}_t$
6.            $\boldsymbol{\omega}^{(i+1)} = \boldsymbol{\omega}^{(i)} + \phi(\boldsymbol{x}_t, \boldsymbol{y}_t) - \phi(\boldsymbol{x}_t, \boldsymbol{y})$
7.            $i = i + 1$
8.     return $\boldsymbol{\omega}^i$

## Kernel Trick – Perceptron Algorithm

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 $\boldsymbol{y} = \arg\max_{\boldsymbol{y}} \boldsymbol{\omega}^{(i)} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y})$
5.        if $\boldsymbol{y} \neq \boldsymbol{y}_t$
6.          $\boldsymbol{\omega}^{(i+1)} = \boldsymbol{\omega}^{(i)} + \phi(\boldsymbol{x}_t, \boldsymbol{y}_t) - \phi(\boldsymbol{x}_t, \boldsymbol{y})$
7.          $i = i + 1$
8.    return $\boldsymbol{\omega}^i$

- Each feature function $\phi(\boldsymbol{x}_t, \boldsymbol{y}_t)$ is added and $\phi(\boldsymbol{x}_t, \boldsymbol{y})$ is subtracted to $\boldsymbol{\omega}$ say $\alpha_{\boldsymbol{y},t}$ times
  - $\alpha_{\boldsymbol{y},t}$ is proportional to the $\#$ of times during learning label $\boldsymbol{y}$ is predicted for example $t$ and caused an update because of misclassification

## Kernel Trick – Perceptron Algorithm

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 $\boldsymbol{y} = \arg\max_{\boldsymbol{y}} \boldsymbol{\omega}^{(i)} \cdot \phi(\boldsymbol{x}_t, \boldsymbol{y})$
5.          if $\boldsymbol{y} \neq \boldsymbol{y}_t$
6.              $\boldsymbol{\omega}^{(i+1)} = \boldsymbol{\omega}^{(i)} + \phi(\boldsymbol{x}_t, \boldsymbol{y}_t) - \phi(\boldsymbol{x}_t, \boldsymbol{y})$
7.              $i = i + 1$
8.      return $\boldsymbol{\omega}^i$

▶ Each feature function $\phi(\boldsymbol{x}_t, \boldsymbol{y}_t)$ is added and $\phi(\boldsymbol{x}_t, \boldsymbol{y})$ is subtracted to $\boldsymbol{\omega}$ say $\alpha_{\boldsymbol{y},t}$ times

  ▶ $\alpha_{\boldsymbol{y},t}$ is proportional to the # of times during learning label $\boldsymbol{y}$ is predicted for example $t$ and caused an update because of misclassification

▶ Thus,

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

## Kernel Trick – Perceptron Algorithm

- We can re-write the argmax function as:

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

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

# Kernel Trick – Perceptron Algorithm

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

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

▶ Testing on unseen instance $\boldsymbol{x}$:

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

# 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
- ▶ Disadvantage: Exact kernel methods depend polynomially on the number of training examples - infeasible for large datasets

# Kernels for Large Training Sets

▶ Alternative to exact kernels: Explicit randomized feature map
[Rahimi and Recht 2007]

  ▶ Shallow neural network by random Fourier/Binning transformation:

    ▶ Random weights from input to hidden units
    ▶ Cosine as transfer function
    ▶ Convex optimization of weights from hidden to output units

## 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

## **Outlook**

- ▶ Multiclass linear models are basic building blocks for further lectures: structured output prediction, graphical models, multilayer perceptron neural networks

# Outlook

▶ Multiclass linear models are basic building blocks for further lectures: structured output prediction, graphical models, multilayer perceptron neural networks

▶ Kernel Machines
  ▶ Kernel Machines introduce nonlinearity by using specific feature maps or kernels
  ▶ Feature map or kernel is not part of optimization problem, thus convex optimization of loss function for linear model possible

▶ Neural Networks
  ▶ Similarities and nonlinear combinations of features are learned: representation learning
  ▶ We lose the advantages of convex optimization since objective functions will be nonconvex

Wrap up and time for questions

# 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. *Machine Learning*, (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. *Computational 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 *Proceedings of the 38th Annual Meeting of the 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. *Journal of Machine Learning Research*, 37:277–296.

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

▶ **SVM:**

▶ Vladimir N. Vapnik. 1998.
*Statistical Learning Theory*. Wiley.

▶ Olivier Chapelle. 2007.

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

▶ Ben Taskar, Carlos Guestrin, and Daphne Koller. 2003.
  Max-margin markov networks. In *Advances in Neural Information Processing 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 *Advances in Neural 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. *Annals of Mathematical Statistics*, 22(3):400–407.

▶ Boris T. Polyak. 1964.

Some methods of speeding up the convergence of iteration methods. *USSR 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. *IEEE Transactons on Information Theory*, 50(9):2050–2057.

▶ Ilya Sutskever, James Martens, George E. Dahl, and Geoffrey E. Hinton. 2013. On the importance of initialization and momentum in deep learning. In *Proceedings of International Conference on Machine Learning (ICML)*, pages 1139–1147.

▶ Leon Bottou, Frank E. Curtis, and Jorge Nocedal. 2016. Optimization methods for large-scale machine learning. *CoRR*, abs/arXiv:1606.04838v1.

## Thanks

- Thanks to Steven Abney, Shay Cohen, Chris Dyer, Philipp Koehn, Ryan McDonald for allowing the use of some their slide materials in this lecture.