

Foundations of Artificial Intelligence

Perceptrons and Optimal Hyperplanes

CS472 – Fall 2007
Thorsten Joachims

Example: Majority-Vote Function

- **Definition: Majority-Vote Function f_{majority}**
 - N binary attributes, i.e. $x \in \{0,1\}^N$
 - If more than N/2 attributes in x are true, then $f_{\text{majority}}(x)=1$, else $f_{\text{majority}}(x)=-1$.
- **How can we represent this function as a decision tree?**
 - Huge and awkward tree!
- **Is there an “easier” representation of f_{majority} ?**

Example: Spam Filtering

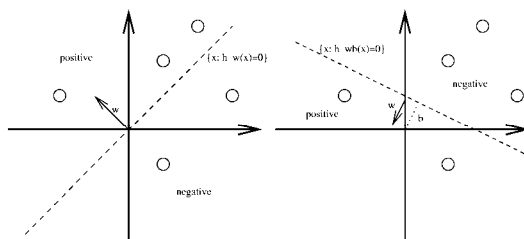
	viagra	learning	the	dating	nigeria	spam?
$\vec{x}_1 =$	1	0	1	0	0	$y_1 = 1$
$\vec{x}_2 =$	0	1	1	0	0	$y_2 = -1$
$\vec{x}_3 =$	0	0	0	0	1	$y_3 = 1$

- **Instance Space X:**
 - Feature vector of word occurrences => binary features
 - N features (N typically > 50000)
- **Target Concept c:**
 - Spam (+1) / Ham (-1)
- **Type of function to learn:**
 - Set of Spam words S, Set of Ham words H
 - Classify as Spam (+1), if more Spam words than Ham words in example.

Linear Classification Rules

- **Hypotheses of the form**
 - unbiased: $h_{\vec{w}}(\vec{x}) = \begin{cases} 1 & w_1x_1 + \dots + w_Nx_N > 0 \\ -1 & \text{else} \end{cases}$
 - biased: $h_{\vec{w},b}(\vec{x}) = \begin{cases} 1 & w_1x_1 + \dots + w_Nx_N + b > 0 \\ -1 & \text{else} \end{cases}$
 - Parameter vector w, scalar b
- **Hypothesis space H**
 - $H_{\text{unbiased}} = \{h_{\vec{w}} : \vec{w} \in \mathbb{R}^N\}$
 - $H_{\text{biased}} = \{h_{\vec{w},b} : \vec{w} \in \mathbb{R}^N, b \in \mathbb{R}\}$
- **Notation**
 - $w_1x_1 + \dots + w_Nx_N = \vec{w} \cdot \vec{x}$ and $\text{sign}(a) = \begin{cases} 1 & a > 0 \\ -1 & \text{else} \end{cases}$
 - $h_{\vec{w}}(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x})$
 - $h_{\vec{w},b}(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b)$

Geometry of Hyperplane Classifiers



- **Linear Classifiers divide instance space as hyperplane**
- **One side positive, other side negative**

(Batch) Perceptron Algorithm

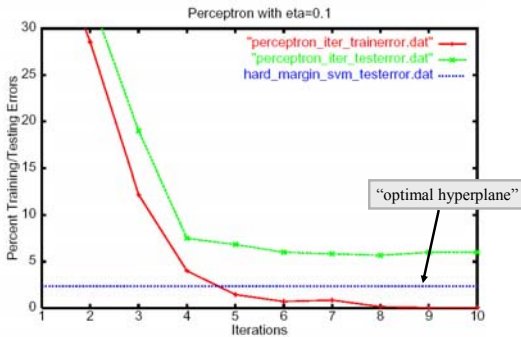
Input: $D = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$, $\vec{x}_i \in \mathbb{R}^N$, $y_i \in \{-1, 1\}$,
 $\eta \in \mathbb{R}$, $I \in [1, 2, \dots]$

Algorithm:

- $\vec{w}_0 = \vec{0}$, $k = 0$
- repeat
 - FOR $i=1$ TO n
 - * IF $y_i(\vec{w}_k \cdot \vec{x}_i) \leq 0$ ### makes mistake
 - $\vec{w}_{k+1} = \vec{w}_k + \eta y_i \vec{x}_i$
 - $k = k + 1$
 - * ENDIF
 - ENDFOR
- until I iterations reached

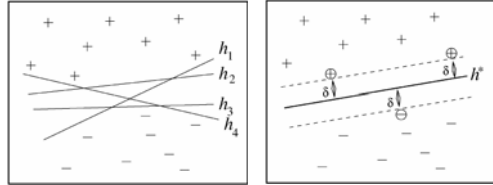
	x_1	x_2	y
$\vec{x}_1 =$	1	2	$y_1 = 1$
$\vec{x}_2 =$	2	1	$y_2 = 1$
$\vec{x}_3 =$	-1	-1	$y_3 = -1$
$\vec{x}_4 =$	-1	1	$y_4 = -1$

Example: Reuters Text Classification



Optimal Hyperplanes

Assumption: Training examples are linearly separable.



Definition: For a linear classifier $h_{\vec{w},b}$, the margin δ of an example (\vec{x}, y) is $\delta = y(\vec{w} \cdot \vec{x} + b)$.

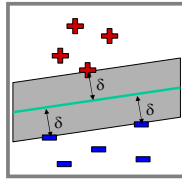
Definition: The margin is called geometric margin, if $\|\vec{w}\| = 1$. Otherwise, functional margin.

Hard-Margin Separation

Goal: Find hyperplane with the largest distance to the closest training examples.

Optimization Problem (Primal):

$$\begin{aligned} \min_{\vec{w}, b} \quad & \frac{1}{2} \vec{w} \cdot \vec{w} \\ \text{s.t.} \quad & y_1(\vec{w} \cdot \vec{x}_1 + b) \geq 1 \\ & \dots \\ & y_n(\vec{w} \cdot \vec{x}_n + b) \geq 1 \end{aligned}$$



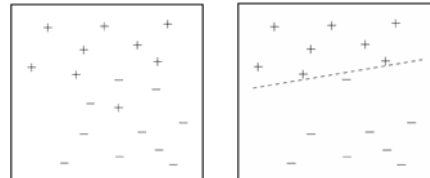
Support Vectors: Examples with minimal distance (i.e. margin).

Definition: The (hard) margin of a linear classifier $h_{\vec{w},b}$ on data D is $\delta = \min_{(\vec{x}, y) \in D} \{y(\vec{w} \cdot \vec{x} + b)\}$.

Non-Separable Training Data

Limitations of hard-margin formulation

- For some training data, there is no separating hyperplane.
- Complete separation (i.e. zero training error) can lead to suboptimal prediction error.



Soft-Margin Separation

Idea: Maximize margin and minimize training error.

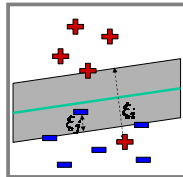
Hard-Margin OP (Primal):

$$\begin{aligned} \min_{\vec{w}, b} \quad & \frac{1}{2} \vec{w} \cdot \vec{w} \\ \text{s.t.} \quad & y_1(\vec{w} \cdot \vec{x}_1 + b) \geq 1 \\ & \dots \\ & y_n(\vec{w} \cdot \vec{x}_n + b) \geq 1 \end{aligned}$$

Soft-Margin OP (Primal):

$$\begin{aligned} \min_{\vec{w}, b, \xi} \quad & \frac{1}{2} \vec{w} \cdot \vec{w} + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_1(\vec{w} \cdot \vec{x}_1 + b) \geq 1 - \xi_1 \wedge \xi_1 \geq 0 \\ & \dots \\ & y_n(\vec{w} \cdot \vec{x}_n + b) \geq 1 - \xi_n \wedge \xi_n \geq 0 \end{aligned}$$

- Slack variable ξ_i measures by how much (x_i, y_i) fails to achieve margin δ
- $\sum \xi_i$ is upper bound on number of training errors
- C is a parameter that controls trade-off between margin and training error.

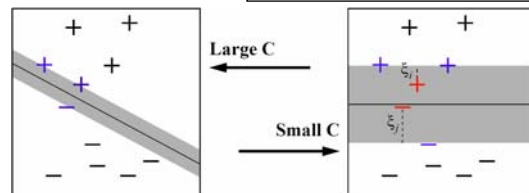


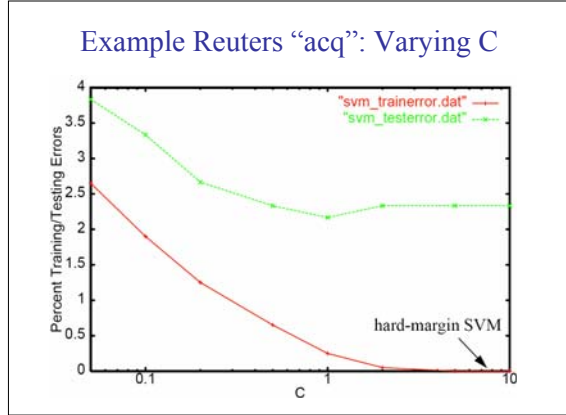
Controlling Soft-Margin Separation

- $\sum \xi_i$ is upper bound on number of training errors
- C is a parameter that controls trade-off between margin and training error.

Soft-Margin OP (Primal):

$$\begin{aligned} \min_{\vec{w}, b, \xi} \quad & \frac{1}{2} \vec{w} \cdot \vec{w} + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_1(\vec{w} \cdot \vec{x}_1 + b) \geq 1 - \xi_1 \wedge \xi_1 \geq 0 \\ & \dots \\ & y_n(\vec{w} \cdot \vec{x}_n + b) \geq 1 - \xi_n \wedge \xi_n \geq 0 \end{aligned}$$





Example: Margin in High-Dimension

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	y
Example 1	1	0	0	1	0	0	0	1
Example 2	1	0	0	0	1	0	0	1
Example 3	0	1	0	0	0	1	0	-1
Example 4	0	1	0	0	0	0	1	-1
	w_1	w_2	w_3	w_4	w_5	w_6	w_7	b
Hyperplane 1	1	1	0	0	0	0	0	2
Hyperplane 2	0	0	0	1	1	-1	-1	0
Hyperplane 3	1	-1	1	0	0	0	0	0
Hyperplane 4	1	-1	0	0	0	0	0	0
Hyperplane 5	0.95	-0.95	0	0.05	0.05	-0.05	-0.05	0