# Linear Classifiers and Perceptron

CS678 Advanced Topics in Machine Learning
Thorsten Joachims
Spring 2003

Outline:
- Linear classifiers
- Example: text classification
- Perceptron learning algorithm
- Mistake bound for Perceptron
- Separation margin
- Dual representation

---

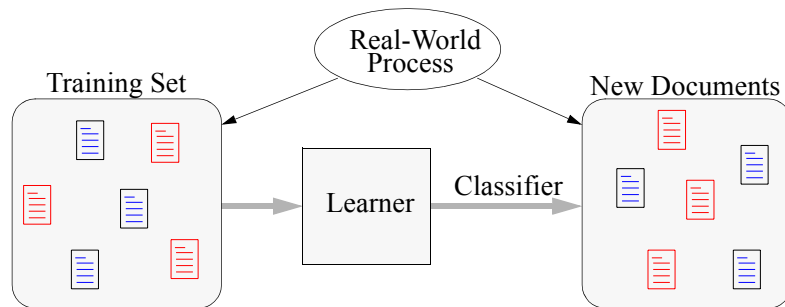# Text Classification

E.D. And F. MAN TO BUY INTO HONG KONG FIRM

The U.K. Based commodity house E.D. And F. Man Ltd and Singapore's Yeo Hiap Seng Ltd jointly announced that Man will buy a substantial stake in Yeo's 71.1 pct held unit, Yeo Hiap Seng Enterprises Ltd. Man will develop the locally listed soft drinks manufacturer into a securities and commodities brokerage arm and will rename the firm Man Pacific (Holdings) Ltd.

About a corportate acquisition?

JA                NEIN

---

# Learning Text Classifiers



Real-World Process

Training Set → Learner → Classifier → New Documents

**Goal:**
- Learner uses training set to find classifier with low prediction error.

---

# Generative vs. Discriminative Training

**Process:**
- Generator: Generates descriptions $\vec{x}$ according to distribution $P(\vec{x})$.
- Teacher: Assigns a value $y$ to each description $\vec{x}$ based on $P(y|\vec{x})$.

=> Training examples $(\vec{x}_1, y_1), ..., (\vec{x}_n, y_n) \sim P(\vec{x}, y)$   $\vec{x}_i \in \Re^N$  $y_i \in \{1, -1\}$
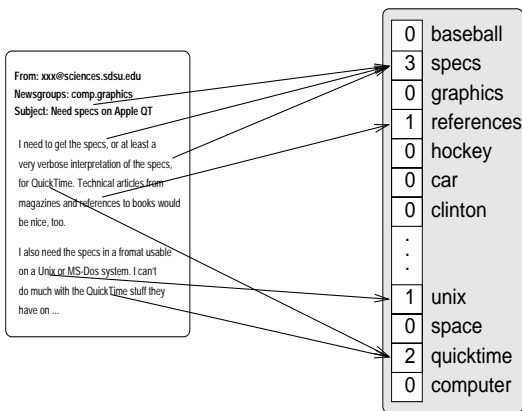
**Generative Training**
- make assumptions about the parametric form of $P(\vec{x}, y)$.
- estimate the parameters of $P(\vec{x}, y)$ from the training data
- derive optimal classifier using Bayes' rule
- example: naive Bayes

**Discriminative Training**
- make assumptions about the set $H$ of classifiers
- estimate error of classifiers in $H$ from the training data
- select classifier with lowest error rate
- example: SVM, decision tree

## Representing Text as Attribute Vectors

| 0 | baseball |
| 3 | specs |
| 0 | graphics |
| 1 | references |
| 0 | hockey |
| 0 | car |
| 0 | clinton |
| . | |
| . | |
| . | |
| 1 | unix |
| 0 | space |
| 2 | quicktime |
| 0 | computer |

**Attributes:** Words (Word-Stems)

**Values:** Occurrence-Frequencies

==> The ordering of words is ignored!

---

## Linear Classifiers (Example)

**Text Classification: Physics (+1) versus Receipes (-1)**

| ID | nuclear $(x_1)$ | atom $(x_2)$ | salt $(x_3)$ | pepper $(x_4)$ | water $(x_5)$ | heat $(x_6)$ | and $(x_7)$ | y |
|---|---|---|---|---|---|---|---|---|
| D1 | 1 | 2 | 0 | 0 | 2 | 0 | 2 | +1 |
| D2 | 0 | 0 | 0 | 3 | 0 | 1 | 1 | -1 |
| D3 | 0 | 2 | 1 | 0 | 0 | 0 | 3 | +1 |
| D4 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | -1 |

| w,b | 2 | 3 | -1 | -3 | -1 | -1 | 0 | b=1 |
|---|---|---|---|---|---|---|---|---|

D1: $\sum_{i=1}^{7} \vec{w_i}\vec{x_i} + b = [2\cdot 1 + 3\cdot 2 + (-1)\cdot 0 + (-3)\cdot 0 + (-1)\cdot 2 + (-1)\cdot 0 + 0\cdot 2] + 1$
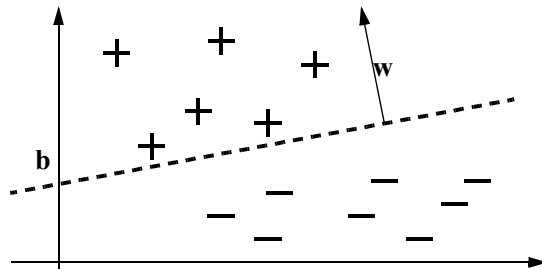
D2: $\sum_{i=1}^{7} \vec{w_i}\vec{x_i} + b = [2\cdot 0 + 3\cdot 0 + (-1)\cdot 0 + (-3)\cdot 3 + (-1)\cdot 0 + (-1)\cdot 1 + 0\cdot 1] + 1$

---

## Linear Classifiers

**Rules of the Form:** weight vector $\vec{w}$, threshold $b$

$$h(\vec{x}) = sign\left[\sum_{i=1}^{N} \vec{w_i}\vec{x_i} + b\right] = \begin{cases} 1 & if \sum_{i=1}^{N} \vec{w_i}\vec{x_i} + b > 0 \\ -1 & else \end{cases}$$

**Geometric Interpretation (Hyperplane):**



---

## Perceptron (Rosenblatt)

Input: $S = \{(\vec{x_1}, y_1), ..., (\vec{x_n}, y_n)\} \quad \vec{x_i} \in \Re^N \quad y_i \in \{1, -1\}$ (linear separable)

- $w_0 \leftarrow 0; b_0 \leftarrow 0; k \leftarrow 0$
- $R = max_i \|\vec{x_i}\|$
- repeat
  - for i=1 to n
    - if $y_i(\vec{w_k} \cdot \vec{x_i} + b_k) \leq 0$
      - $\vec{w_{k+1}} \leftarrow \vec{w_k} + \eta y_i \vec{x_i}$
      - $b_{k+1} \leftarrow b_k + \eta y_i R^2$
      - $k \leftarrow k = 1$
    - endif
  - endfor
- until no mistakes made in the for loop
- return $(\vec{w_k}, b_k)$

## Analysis of Perceptron

**Definition (Margin of an Example):** The margin of an example $(\vec{x}_i, y_i)$ with respect to the hyperplane $(\vec{w}, b)$ is

$$\delta_i = y_i(w \cdot x_i + b)$$

**Definition (Margin of an Example):** The margin of a training set $S = \{(\vec{x}_1, y_1), ..., (\vec{x}_n, y_n)\}$ with respect to the hyperplane $(\vec{w}, b)$ is

$$\delta = min_i \ \ y_i(w \cdot x_i + b)$$

**Theorem (Novikoff):** If for a training set S there exists a weight vector with margin $\delta$, then the perceptron makes at most

$$4\left(\frac{R^2}{\delta^2}\right)$$

mistakes before returning a separating hyperplane.

---

## Dual Perceptron

- For each example $(\vec{x}_i, y_i)$, count with $\alpha_i$ the number of times the perceptron algorithm makes a mistake on it. Then
$$\vec{w} = \sum_{i=1}^{n} \alpha_i y_i \vec{x}_i$$
- $\vec{\alpha} = 0; b_0 \leftarrow 0$ and $R = max_i \|\vec{x}_i\|$
- repeat
  - for i=1 to n
    - if $y_i\left(\sum_{j=1}^{n} \alpha_i y_i (\vec{x}_j \cdot \vec{x}_i) + b_k\right) \leq 0$
      - $\vec{\alpha}_i \leftarrow \alpha_i + 1$
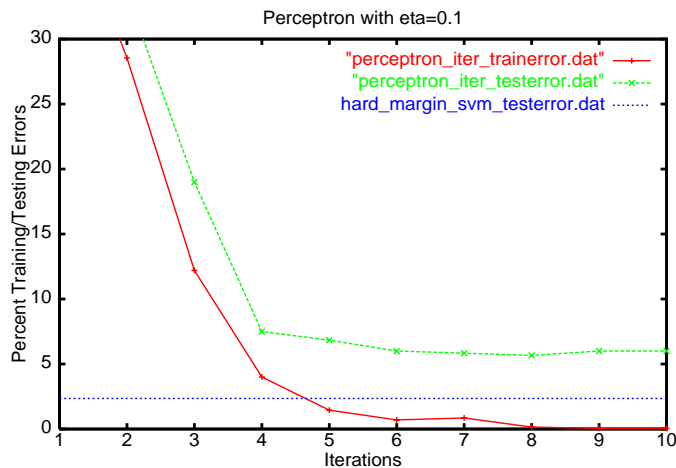      - $b \leftarrow b + y_i R^2$
    - endif
  - endfor
- until no mistakes made in the for loop
- return $(\vec{\alpha}, b)$

---

## Experiment: Perceptron for Text Classification



Perceptron with eta=0.1

"perceptron_iter_trainerror.dat"
"perceptron_iter_testerror.dat"
hard_margin_svm_testerror.dat

Train on 1000 pos / 1000 neg examples for " acq" (Reuters-21578).