

Announcements

- HW1 due today
- Most of last lecture was on the blackboard.
- Gather hand data set today.

CSF 402 - Fall 05

Linear Discriminant Functions (Sections 5.1-5.2)

•

Linear discriminant functions and decisions surfaces

3

• Definition

It is a function that is a linear combination of the components of x

$$g(x) = w^T x + w_0 \quad (1)$$

where w is the weight vector and w_0 the bias

- A two-category classifier with a discriminant function of the form (1) uses the following rule:
Decide ω_1 if $g(x) > 0$ and ω_2 if $g(x) < 0$
 \Leftrightarrow Decide ω_1 if $w^T x > -w_0$ and ω_2 otherwise
If $g(x) = 0 \Rightarrow x$ is assigned to either class

Pattern Classification, Ch4 (Part 1)

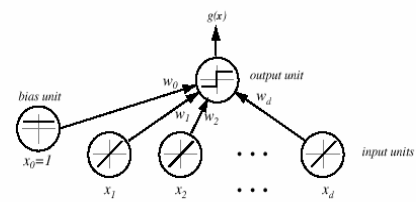


FIGURE 5.1. A simple linear classifier having d input units, each corresponding to the values of the components of an input vector. Each input feature value x_i is multiplied by its corresponding weight w_i ; the effective input at the output unit is the sum all these products, $\sum w_i x_i$. We show in each unit its effective input-output function. Thus each of the d input units is linear, emitting exactly the value of its corresponding feature value. The single bias unit unit always emits the constant value 1.0. The single output unit emits a +1 if $w^T x + w_0 > 0$ or a -1 otherwise. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Pattern Classification, Ch4 (Part 1)

- The equation $g(x) = 0$ defines the **decision surface** that separates points assigned to the category ω_1 from points assigned to the category ω_2
- When $g(x)$ is linear, the decision surface is a hyperplane
- Algebraic measure of the distance from x to the hyperplane (interesting result!)

Pattern Classification, Ch4 (Part 1)

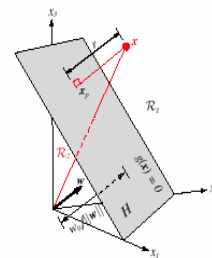


FIGURE 5.2. The linear decision boundary H , where $g(x) = w^T x + w_0 = 0$, separates the feature space into two half-spaces \mathcal{R}_1 (where $g(x) > 0$) and \mathcal{R}_2 (where $g(x) < 0$). From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Pattern Classification, Ch4 (Part 1)

7

$x = x_p + \frac{r \cdot w}{\|w\|}$ (since w is colinear with $x - x_p$ and $\frac{w}{\|w\|} = 1$)
 since $g(x) = 0$ and $w' \cdot w = \|w\|^2$
 therefore $r = \frac{g(x)}{\|w\|}$
 in particular $d(0, H) = \frac{w_0}{\|w\|}$

- In conclusion, a linear discriminant function divides the feature space by a hyperplane decision surface
- The orientation of the surface is determined by the normal vector w and the location of the surface is determined by the bias

Pattern Classification, Ch4 (Part 1)

8

The multi-category case

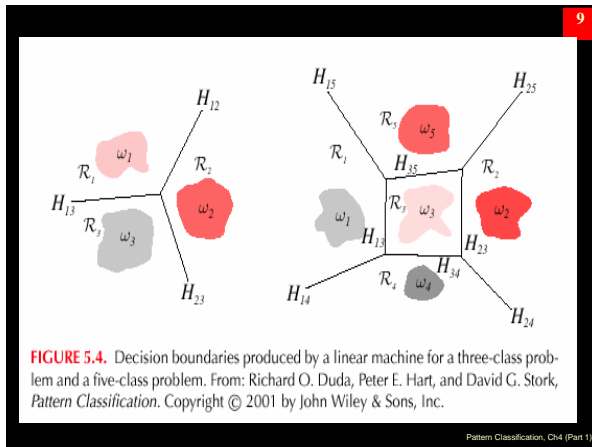
- We define c linear discriminant functions

$$g_i(x) = w_i'x + w_{i0} \quad i = 1, \dots, c$$
 and assign x to ω_i if $g_i(x) > g_j(x) \forall j \neq i$, in case of ties, the classification is undefined
- In this case, the classifier is a "linear machine"
- A linear machine divides the feature space into c decision regions, with $g_i(x)$ being the largest discriminant if x is in the region R_i
- For a two contiguous regions R_i and R_j , the boundary that separates them is a portion of hyperplane H_{ij} defined by:

$$g_i(x) = g_j(x) \iff (w_i - w_j)'x + (w_{i0} - w_{j0}) = 0$$
- $w_i - w_j$ is normal to H_{ij} and

$$d(x, H_{ij}) = \frac{g_i - g_j}{\|w_i - w_j\|}$$

Pattern Classification, Ch4 (Part 1)



10

- It is easy to show that the decision regions for a linear machine are convex, this restriction limits the flexibility and accuracy of the classifier

Pattern Classification, Ch4 (Part 1)

Perceptron

output

$$o = \begin{cases} 1 & \text{if } net > \theta \\ 0 & \text{otherwise} \end{cases}$$

$$net = \sum_{i=1}^n w_i x_i$$

input

Linear, threshold units

x_i : inputs

w_i : weights

θ : threshold

$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_1 x_1 + \dots + w_n x_n > \theta \\ -1 & \text{otherwise.} \end{cases}$$

CS348, Fall 2001 © David Kriegman, 2001

Perceptron

output

$$o = \begin{cases} 1 & \text{if } net > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$net = \sum_{i=1}^n w_i x_i$$

input

The threshold can be easily forced to 0 by introducing an additional weight input $w_0 = \theta$

$x_0 = 1$

w_0

$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise.} \end{cases}$$

CS348, Fall 2001 © David Kriegman, 2001

How powerful is a perceptron? Threshold = 0

Inverter

input x_1	output
0	1
1	0

Boolean AND

input x_1	input x_2	output
0	0	0
0	1	0
1	0	0
1	1	1

Boolean OR

input x_1	input x_2	output
0	0	0
0	1	1
1	0	1
1	1	1

Boolean XOR

input x_1	input x_2	output
0	0	0
0	1	1
1	0	1
1	1	0

CS348, Fall 2001 © David Kriegman, 2001

Concept Space & Linear Separability

Linear Separability

Linear Separability

Linear Separability

CS348, Fall 2001 © David Kriegman, 2001

Training Perceptron

Perceptron Training Rule

$$w_i \leftarrow w_i + \Delta w_i$$

$\Delta w_i = \eta (t - o) x_i$

new weight ← old weight + increment
 increment ← step size × (target - perceptron output) × input

Converges, if...

- ... training data linearly separable
- ... step size η sufficiently small
- ... no "hidden" units

Where:

- t is target value
- o is perceptron output
- η is small constant (e.g., .1) called *learning rate*

CS348, Fall 2001 © David Kriegman, 2001

Gradient Descent

- Learn w_i 's that minimize squared error

$$E[\vec{w}] = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$D = \text{training data}$

16

CS348, Fall 2001 © David Kriegman, 2001

Gradient Descent

Gradient: $\nabla E[\vec{w}] = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$

Training rule: $\Delta \vec{w} = -\eta \nabla E[\vec{w}]$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

CS348, Fall 2001 © David Kriegman, 2001

Gradient Descent

- To find the best direction in the feature space we compute the gradient of E with respect to each of the components of \vec{w}

$$\nabla E(\vec{w}) \equiv \left[\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_n} \right]$$

- This vector specifies the direction that produces the steepest increase in E ;
- We want to modify \vec{w} in the direction of $-\nabla E(\vec{w})$

$$\vec{w} = \vec{w} + \Delta \vec{w}$$

- Where:

$$\Delta \vec{w} = -\mathbf{R} \nabla E(\vec{w})$$

CS348, Fall 2001 © David Kriegman, 2001

Batch Learning

- Initialize each w_i to small random value
- Repeat until termination:
 - $\Delta w_i = 0$
 - For each training example d do
 - $o_d \leftarrow \sigma(\sum_i w_i x_{i,d})$
 - $\Delta w_i \leftarrow \Delta w_i + \eta (t_d - o_d) o_d (1 - o_d) x_{i,d}$
 - $w_i \leftarrow w_i + \Delta w_i$

CS348, Fall 2001

© David Kriegman, 2001

Incremental (Online) Learning

- Initialize each w_i to small random value
- Repeat until termination:
 - For each training example d do
 - $\Delta w_i = 0$
 - $o_d \leftarrow \sum_i w_i x_{i,d}$
 - $\Delta w_i \leftarrow \Delta w_i + \eta (t_d - o_d) o_d (1 - o_d) x_{i,d}$
 - $w_i \leftarrow w_i + \Delta w_i$

20

CS348, Fall 2001

© David Kriegman, 2001

Summary: Single Layer Networks

- Variety of update rules
 - Multiplicative $\Delta w_i = \mathbf{R}(t_d - o_d) x_{id}$
 - Additive
- Batch and incremental algorithms
- Various convergence and efficiency conditions
- There are other ways to learn linear functions
 - Linear Programming (general purpose)
 - Probabilistic Classifiers (some assumption)
- Although simple and restrictive -- linear predictors perform very well on many realistic problems
- However, the representational restriction is limiting in many applications

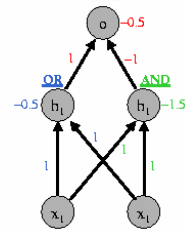
CS348, Fall 2001

© David Kriegman, 2001

Increasing Expressiveness: Multi-Layer Neural Networks

Boolean XOR

input x_1	input x_2	output
0	0	0
0	1	1
1	0	1
1	1	0



2-layer Neural Net

CS348, Fall 2001

© David Kriegman, 2001