• Assumption: we *know* the proper *forms* for the *discriminant functions*, and use the samples to estimate the values of parameters of the classifier

- The problem of minimizing a criterion function
  - Sample risk, or training error the average loss incurred in classifying the set of training samples

#### **Linear Discriminant Functions**

• A classifier that uses linear discriminant functions is called "*a linear machine*"

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

• The decision surfaces for a linear machine are pieces of hyperplanes defined by:

$$g_i(\mathbf{x}) = g_j(\mathbf{x})$$

#### **Linear Discriminant Functions**

The two-category case

$$g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x}) = \mathbf{w}_1^T \mathbf{x} + w_{10} - (\mathbf{w}_2^T \mathbf{x} + w_{20})$$
$$= (\mathbf{w}_1 - \mathbf{w}_2)^T \mathbf{x} + (w_{10} - w_{20}) = \mathbf{w}^T \mathbf{x} + w_0$$

# W is the weight vector $w_0$ is the bias, or threshold weight

**Decision rule:** decide  $\omega_1$  if g(x)>0,  $\omega_2$  otherwise



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 output unit sums all these products and emits a +1 if  $\mathbf{w}^t \mathbf{x} + w_0 > 0$  or a -1 otherwise.

- g(x) = 0 defines a hyperplane separating the feature space into decision regions R<sub>1</sub> and R<sub>2</sub>
- g(x) gives an algebraic measure of the distance from x to the hyperplane H



$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

Since  $g(\mathbf{x}_p) = 0$ :  $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = r \|\mathbf{w}\|$ 

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = r \|\mathbf{w}\|$$
$$\mathbb{I}$$
$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}$$

• The orientation of H is determined by the normal vector w, and the location of H is determined by the bias  $w_0$ 



#### The Multicategory Case



One versus one



Figure 5.3: Linear decision boundaries for a four-class problem. The top figure shows  $\omega_i/\text{not }\omega_i$  dichotomies while the bottom figure shows  $\omega_i/\omega_j$  dichotomies. The pink regions have ambiguous category assignments.

.

#### The Multicategory Case

- To avoid ambiguous regions:
  - Define c linear discriminant functions

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0} \qquad i = 1, ..., c$$

 Assign x to class ω<sub>i</sub> if g<sub>i</sub>(x) > g<sub>j</sub>(x) for all j different from i

#### The Multicategory Case



Figure 5.4: Decision boundaries produced by a linear machine for a three-class problem and a five-class problem.

10

.

#### **Generalized Linear Discriminant Functions**

• Linear discriminant function:

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = w_0 + \sum_{i=1}^d w_i x_i$$

• Polynomial discriminant functions:

$$g(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j + \sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d w_{ijk} x_i x_j x_k + \dots$$

#### **Generalized Linear Discriminant Functions**

• Generalized linear discriminant function:

$$g(\mathbf{x}) = \sum_{i=1}^{\hat{d}} a_i y_i(\mathbf{x}) = \mathbf{a}^T \mathbf{y}$$



#### **Generalized Linear Discriminant Functions**

• Generalized linear discriminant function:

$$g(\mathbf{x}) = \sum_{i=1}^{\hat{d}} a_i y_i(\mathbf{x}) = \mathbf{a}^T \mathbf{y}$$

- The resulting discriminant function is not linear in x, but it is linear in y
- The functions  $y_i(x)$  map points in ddimensional x-space to points in  $\hat{d}$  dimensional y-space



Figure 5.5: The mapping  $\mathbf{y} = (1, x, x^2)^t$  takes a line and transforms it to a parabola in three dimensions. A plane splits the resulting  $\mathbf{y}$  space into regions corresponding to two categories, and this in turn gives a non-simply connected decision region in the one-dimensional x space.



Figure 5.6: The two-dimensional input space  $\mathbf{x}$  is mapped through a polynomial function f to  $\mathbf{y}$ . Here the mapping is  $y_1 = x_1$ ,  $y_2 = x_2$  and  $y_3 \propto x_1x_2$ . A linear discriminant in this transformed space is a hyperplane, which cuts the surface. Points to the positive side of the hyperplane  $\hat{H}$  correspond to category  $\omega_1$ , and those beneath it  $\omega_2$ . Here, in terms of the  $\mathbf{x}$  space,  $\mathcal{R}_1$  is a not simply connected.

# http://www.youtube.com/watch?v=3liCbRZPrZ

.

#### **Generalized Linear Discriminant Functions**



- We have a set of n samples  $y_1, ..., y_n$  some labelled  $\omega_1$  and some labelled  $\omega_2$
- We want to use these samples to determine a in  $g(\mathbf{x}) = \mathbf{a}^T \mathbf{y}$

• If we can find such a that all of the samples are correctly classified, the samples are *linearly separable* 

- $g(\mathbf{x}) = \mathbf{a}^T \mathbf{y}$
- $\mathbf{y}_i$  is classified correctly if  $\mathbf{a}^T \mathbf{y}_i > 0$  and  $\mathbf{y}_i$  is labelled  $\boldsymbol{\omega}_1$  or  $\mathbf{a}^T \mathbf{y}_i < 0$  and  $\mathbf{y}_i$  is labelled  $\boldsymbol{\omega}_2$
- Normalization: replace all samples labelled  $\omega_{\text{2}}$  by their negatives
- We'll look for such solution vector a that a<sup>T</sup>y<sub>i</sub> > 0 for all of the samples

- a can be thought as a point in weight space
- Each y<sub>i</sub> places a constraint on the possible location of a
- The solution vector must be on the positive side of every hyperplane





Figure 5.8: Four training samples (black for  $\omega_1$ , red for  $\omega_2$ ) and the solution region in feature space. The figure on the left shows the raw data; the solution vectors leads to a plane that separates the patterns from the two categories. In the figure on the right, the red points have been "normalized" — i.e., changed in sign. Now the solution vector leads to a plane that places all "normalized" points on the same side.



Figure 5.9: The effect of the margin on the solution region. At the left, the case of no margin (b = 0) equivalent to a case such as shown at the left in Fig. 5.8. At the right is the case b > 0, shrinking the solution region by margins  $b/||\mathbf{y}_i||$ .

- Gradient descent procedures: define a criterion function J(a) that is minimized if a is a solution vector
- Basic gradient descent:
  - Start with arbitrary weight vector a(1)
  - Compute  $\nabla J(\mathbf{a}(1))$
  - Repeat iterations (moving weight vector in the direction of steepest descent):

$$\mathbf{a}(k+1) = \mathbf{a}(k) - \eta(k)\nabla J(\mathbf{a}(k))$$

 $\eta(k)$  is learning rate

Algorithm 1 (Basic gradient descent)

- 1 <u>begin</u> <u>initialize</u> a, criterion  $\theta, \eta(\cdot), k = 0$
- $2 \quad \underline{\mathbf{do}} \ k \leftarrow k+1$
- 3  $\mathbf{a} \leftarrow \mathbf{a} \eta(k) \nabla J(\mathbf{a})$
- $4 \quad \underline{\text{until}} \ \eta(k) \boldsymbol{\nabla} J(\mathbf{a}) < \theta$
- 5 <u>return</u> a
- $6 \, \underline{\mathrm{end}}$

#### Choice of the learning rate is important

- If it is too small, convergence is too slow
- If it is too large, the correction process will overshoot

• Setting the learning rate:

$$\eta(k) = \frac{\|\nabla J\|^2}{\nabla J^T \mathbf{H} \nabla J}$$

H is the Hessian matrix of second partial derivatives  $\partial^2 J/\partial a_i \partial a_j$ 

Algorithm 2 (Newton descent)

$$\begin{array}{cccc} 1 & \underline{\text{begin initialize}} & \mathbf{a}, \text{criterion } \theta \\ 2 & \underline{\text{do}} \\ 3 & \mathbf{a} \leftarrow \mathbf{a} - \mathbf{H}^{-1} \nabla J(\mathbf{a}) \\ 4 & \underline{\text{until}} & \mathbf{H}^{-1} \nabla J(\mathbf{a}) < \theta \\ 5 & \underline{\text{return }} \mathbf{a} \\ 6 & \text{end} \end{array}$$



Figure 5.10: The sequence of weight vectors given by a simple gradient descent method (red) and by Newton's (second order) algorithm (black). Newton's method typically leads to greater improvement per step, even when using optimal learning rates for both methods. However the added computational burden of inverting the Hessian matrix used in Newton's method is not always justified, and simple descent may suffice.

- Problem: to construct a criterion function for solving  $\mathbf{a}^T \mathbf{y}_i > 0$
- Possibility: Let J(a; y<sub>1</sub>, ..., y<sub>n</sub>) be the number of samples misclassified by a.
- Disadvantage: this function is piecewise constant



• Perceptron criterion function:

$$J_p(\mathbf{a}) = \sum_{y \in \mathcal{Y}} (-\mathbf{a}^T \mathbf{y})$$

where  $\mathcal{Y}(a)$  is the set of samples misclassified by a 5

•  $J_p(a)$  is proportional to the sum of the distances from the misclassified samples to the decision boundary



• 
$$\nabla J_p = \sum_{y \in \mathcal{Y}} (-\mathbf{y})$$

# (since the jth component of the gradient of $\mathit{J}_p$ is $\partial \mathit{J}_p/\partial a_j$ )

Algorithm 3 (Batch Perceptron)

$$1 \underline{\text{begin initialize}}_{\mathbf{2}} \mathbf{a}, \eta(\cdot), \text{criterion } \theta, k = 0$$

$$2 \underline{\text{do}} k \leftarrow k + 1$$

$$3 \mathbf{a} \leftarrow \mathbf{a} + \eta(k) \sum_{\mathbf{y} \in \mathcal{Y}_k} \mathbf{y}$$

$$4 \underline{\text{until}} \eta(k) \sum_{\mathbf{y} \in \mathcal{Y}_k} \mathbf{y} < \theta$$

$$5 \underline{\text{return a}}_{\mathbf{6} \underline{\text{end}}}$$



Figure 5.12: The Perceptron criterion,  $J_p$  is plotted as a function of the weights  $a_1$ and  $a_2$  for a three-pattern problem. The weight vector begins at **0**, and the algorithm sequentially adds to it vectors equal to the "normalized" misclassified patterns themselves. In the example shown, this sequence is  $\mathbf{y}_2, \mathbf{y}_3, \mathbf{y}_1, \mathbf{y}_3$ , at which time the vector lies in the solution region and iteration terminates. Note that the second update (by  $\mathbf{y}_3$ ) takes the candidate vector *farther* from the solution region than after the first update (cf. Theorem 5.1. (In an alternate, batch method, *all* the misclassified points are added at each iteration step leading to a smoother trajectory in weight space.)

• Batch *Perceptron algorithm* : the next weight vector is obtained by adding some multiple of the sum of the misclassified samples to the present weight vector

Algorithm 5 (Variable increment Perceptron with margin)

 $1 \underline{\text{begin initialize}}_{2} \mathbf{a}, \text{criterion } \theta, \text{margin } b, \eta(\cdot), k = 0$   $2 \underline{do} \ k \leftarrow k + 1$   $3 \underline{if} \ \mathbf{a}^{t} \mathbf{y}_{k} + b < 0 \ \underline{then} \ \mathbf{a} \leftarrow \mathbf{a} + \eta(k) \mathbf{y}_{k}$   $4 \underline{until} \ \mathbf{a}^{t} \mathbf{y}_{k} + b > 0 \text{ for all } k$   $5 \underline{return} \mathbf{a}$   $6 \underline{end}$ 

#### **Descent algorithm:**

• Another possible criterion function is:

$$J_q(\mathbf{a}) = \sum_{y \in \mathcal{Y}} (\mathbf{a}^T \mathbf{y})^2$$

# where $\mathcal{Y}(a)$ is the set of samples misclassified by a

- Advantage: gradient is continuous
- Disadvantages: 1) the function is very smooth near the boundary of the solution region; 2) value is dominated by the longest sample vectors



• To avoid these problems, we introduce the criterion function as:

$$J_r(\mathbf{a}) = \frac{1}{2} \sum_{y \in \mathcal{Y}} \frac{(\mathbf{a}^T \mathbf{y} - b)^2}{\|\mathbf{y}\|^2}$$

where  $\mathcal{Y}(\mathbf{a})$  is the set of samples for which  $\mathbf{a}^T \mathbf{y} \leq b$ 

$$\nabla J_r = \sum_{y \in \mathcal{Y}} \frac{\mathbf{a}^T \mathbf{y} - b}{\|\mathbf{y}\|^2} \mathbf{y}$$



#### • The relaxation algorithm becomes:

Algorithm 8 (Batch relaxation with margin) 1 <u>begin</u> initialize  $\mathbf{a}, \eta(\cdot), k = 0$ do  $k \leftarrow k + 1 \mod n$  $\mathcal{D}$  $\mathcal{Y}_k = \{\}$ 3 i = 04 <u>do</u>  $j \leftarrow j + 1$ 5 $\underline{\mathbf{if}}$   $\mathbf{a}^T \mathbf{y}_j \leq b$   $\underline{\mathbf{then}}$  Append  $\mathbf{y}_j$  to  $\mathcal{Y}_k$ 6 until j = n7  $\mathbf{a} \leftarrow \mathbf{a} + \eta(k) \sum_{\mathbf{y} \in \mathcal{V}} \frac{b - \mathbf{a}^t \mathbf{y}}{\|\mathbf{y}\|^2} \mathbf{y}$ 8 until  $\mathcal{Y}_k = \{\}$ g10 return a 11 end

#### • Single-sample relaxation rule with margin:

Algorithm 9 (Single-sample relaxation with margin)

• In general,  $0 < \eta < 2$ 



Figure 5.15: At the left, underrelaxation  $(\eta < 1)$  leads to needlessly slow descent, or even failure to converge. Overrelaxation  $(1 < \eta < 2)$ , shown in the middle) describes overshooting; nevertheless convergence will ultimately be achieved.

#### Nonseparable behavior

- Perceptron and relaxation procedures are called error-correcting procedures
- In practice they are used if one can believe that the error rate for the optimal LDF is low
- Which technique should one apply if sample set is unseparable?

#### **Minimum Squared Error Procedures**

- We will consider now a criterion function that involves *all* of the samples
- We shall try to resolve  $a^T y_i = b_i$ where  $b_i$  are some arbitrarily specified positive constants

$$\begin{pmatrix} Y_{10} & Y_{11} & \dots & Y_{1d} \\ Y_{20} & Y_{21} & \dots & Y_{2d} \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ Y_{n0} & Y_{n1} & \dots & Y_{nd} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ \vdots \\ \vdots \\ b_n \end{pmatrix} \equiv \mathbf{Ya} = \mathbf{b}$$

#### **Minimum Squared Error Procedures**

• **Y**a = b

• If Y were nonsingular, we could obtain a solution as  $a = Y^{-1}b$ 

- Usually Y has more rows than columns → no exact solution for a exists
- We'll search to minimize the squared length of the error vector e = Ya b:

$$J_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{i=1}^n (\mathbf{a}^T \mathbf{y}_i - b_i)^2$$

#### **Minimum Squared Error Procedures**

• We minimize the sum of squared error by a gradient search procedure:

$$\nabla J_s = \sum_{i=1}^n 2(\mathbf{a}^T \mathbf{y}_i - b_i) \mathbf{y}_i = 2\mathbf{Y}^T (\mathbf{Y}\mathbf{a} - \mathbf{b})$$
$$\nabla J_s = 2\mathbf{Y}^T (\mathbf{Y}\mathbf{a} - \mathbf{b}) = 0$$
$$\mathbf{Y}^T \mathbf{Y}\mathbf{a} = \mathbf{Y}^T \mathbf{b}$$

• If  $\mathbf{Y}^T \mathbf{Y}$  is nonsingular, we can solve for a uniquely as:  $\mathbf{a} = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{b} = \mathbf{Y}^{\dagger} \mathbf{b}$ 

where  $\mathbf{Y}^{\dagger} \equiv (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T$  is called the *pseudoinverse* of **Y** 

Example 1: Constructing a linear classifier by matrix pseudoinverse

Suppose we have the following two-dimensional points for two categories:  $\omega_1$ :  $(1,2)^t$  and  $(2,0)^t$ , and  $\omega_2$ :  $(3,1)^t$  and  $(2,3)^t$ , as shown in black and red, respectively, in the figure.

Our matrix  ${\bf Y}$  is therefore

$$\mathbf{Y} = \begin{pmatrix} 1 & 1 & 2\\ 1 & 2 & 0\\ -1 & -3 & -1\\ -1 & -2 & -3 \end{pmatrix}$$

and after a few simple calculations we find that its pseudoinverse is

$$\mathbf{Y}^{\dagger} \equiv \lim_{\epsilon \to 0} (\mathbf{Y}^{t} \mathbf{Y} + \epsilon \mathbf{I})^{-1} \mathbf{Y}^{t} = \begin{pmatrix} 5/4 & 13/12 & 3/4 & 7/12 \\ -1/2 & -1/6 & -1/2 & -1/6 \\ 0 & -1/3 & 0 & -1/3 \end{pmatrix}$$



We arbitrarily let all the margins be equal, i.e.,  $\mathbf{b} = (1, 1, 1, 1)^t$ . Our solution is  $\mathbf{a} = \mathbf{Y}^{\dagger}\mathbf{b} = (11/3, -4/3, -2/3)^t$ , and leads to the decision boundary shown in the figure. Other choices for **b** would typically lead to different decision boundaries, of course.

#### Demo of MSE procedure:

#### **Support Vector Machines**

- Tutorial
  - <u>http://videolectures.net/mlss06tw\_lin\_svm/</u>
- •Library
  - http://www.csie.ntu.edu.tw/~cjlin/libsvm/