LINEAR CLASSIFIERS

0

***** The Problem: Consider a two class task with ω_1 , ω_2

$$g(\underline{x}) = \underline{w}^T \underline{x} + w_0 = 0 =$$

$$w_1 x_1 + w_2 x_2 + \dots + w_l x_l + w$$

Assume $\underline{x}_1, \underline{x}_2$ on the decision hyperplane: $0 = \underline{w}^T \underline{x}_1 + w_0 = \underline{w}^T \underline{x}_2 + w_0 \Longrightarrow$ $\underline{w}^T (\underline{x}_1 - \underline{x}_2) = 0 \quad \forall \underline{x}_1, \underline{x}_2$

LINEAR CLASSIFIERS



The Perceptron Algorithm

➤ Assume linearly separable classes, i.e.,

$$\exists \underline{w} * \begin{cases} \underline{w} *^{T} \underline{x} > 0 & \forall \underline{x} \in \omega_{1} \\ \underline{w} *^{T} \underline{x} < 0 & \forall \underline{x} \in \omega_{2} \end{cases}$$

The case
$$\underline{w}^{*T} \underline{x} + w_0^*$$

falls under the above formulation, since

•
$$\underline{w}' \equiv \begin{bmatrix} \underline{w}^* \\ w_0^* \end{bmatrix}$$
, $\underline{x}' = \begin{bmatrix} \underline{x} \\ 1 \end{bmatrix}$

•
$$\underline{w}^{*T} \underline{x} + w_0^* = \underline{w'}^T \underline{x'} = 0$$

Our goal: Compute a solution, i.e., a hyperplane <u>w</u>, so that

$$\underline{w}^{T} \underline{x} \gtrless 0 \implies \underline{x} \in \begin{cases} \omega_{1} \\ \omega_{2} \end{cases}$$

• The steps

- 1. Define a cost function to be minimized
- 2. Choose an algorithm to minimize the cost function
- 3. The minimum corresponds to a solution

The Cost Function

$$J(\underline{w}) = \sum_{\underline{x} \in Y} (\delta_{\underline{x}} \underline{w}^T \underline{x})$$

• Where *Y* is the subset of the vectors wrongly classified by <u>w</u>. When *Y*=(empty set) a solution is achieved and

$$J(\underline{w}) = 0$$

• Otherwise:

$$\delta_x = -1 \text{ if } \underline{x} \in Y \text{ and } \underline{x} \in \omega_1$$

 $\delta_x = +1 \text{ if } \underline{x} \in Y \text{ and } \underline{x} \in \omega_2$
 $J(\underline{w}) \ge 0$

LINEAR CLASSIFIERS **>** The Perceptron Algorithm

• $J(\underline{w})$ is piecewise linear (WHY?)



≻ The Algorithm

• The philosophy of the gradient descent is adopted.

$$\underline{w}(\text{new}) = \underline{w}(\text{old}) + \Delta \underline{w}$$

$$\Delta \underline{w} = -\mu \frac{\partial J(\underline{w})}{\partial \underline{w}} | \underline{w} = \underline{w}(\text{old})$$
• Wherever valid

• This is the celebrated Perceptron Algorithm

$$\frac{\partial J(\underline{w})}{\partial \underline{w}} = \frac{\partial}{\partial \underline{w}} \left(\sum_{\underline{x} \in Y} \delta_{\underline{x}} \underline{w}^T \underline{x} \right) = \sum_{\underline{x} \in Y} \delta_{\underline{x}} \underline{x}$$
$$\underline{w}(t+1) = \underline{w}(t) - \rho_t \sum_{\underline{x} \in Y} \delta_{\underline{x}} \underline{x}$$



The perceptron algorithm converges in a finite number of iteration steps to a solution if

$$\lim_{t \to \infty} \sum_{k=0}^{t} \rho_k \to \infty, \qquad \qquad \lim_{t \to \infty} \sum_{k=0}^{t} \rho_k^2$$

e.g.,: $\rho_t = \frac{c}{t}$

 $< +\infty$

✤ A useful variant of the perceptron algorithm

$$\begin{split} \underline{w}(t+1) &= \underline{w}(t) + \rho \underline{x}_{(t)} , \begin{cases} \underline{w}^{T}(t) \underline{x}_{(t)} \leq 0 \\ \underline{x}_{(t)} \in \omega_{1} \end{cases} \\ \underline{w}(t+1) &= \underline{w}(t) - \rho \underline{x}_{(t)} , \begin{cases} \underline{w}^{T}(t) \underline{x}_{(t)} \geq 0 \\ \underline{x}_{(t)} \in \omega_{2} \end{cases} \\ \underline{w}(t+1) &= \underline{w}(t) \end{cases} , \text{ otherwise} \end{split}$$

≻ It is a reward and punishment type of algorithm

Example 3.2

Figure 3.4 shows four points in the two-dimensional space. Points (-1, 0), (0, 1) belong to class ω_1 , and points (0, -1), (1, 0) belong to class ω_2 . The goal of this example is to design a linear classifier using the perceptron algorithm in its reward and punishment form. The parameter ρ is set equal to one, and the initial weight vector is chosen as $w(0) = [0, 0, 0]^T$ in the extended three-dimensional space. According to (3.21)–(3.23), the following computations are in order:



LINEAR CLASSIFIERS **>** The Perceptron Algorithm

$$\boldsymbol{w}^{T}(0) \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} = 0, \quad \boldsymbol{w}(1) = \boldsymbol{w}(0) + \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$$

$$\boldsymbol{w}^{T}(1) \begin{bmatrix} 0\\1\\1 \end{bmatrix} = 1 > 0, \quad \boldsymbol{w}(2) = \boldsymbol{w}(1)$$

$$\boldsymbol{w}^{T}(2) \begin{bmatrix} 0\\-1\\1 \end{bmatrix} = 1 > 0, \quad \boldsymbol{w}(3) = \boldsymbol{w}(2) - \begin{bmatrix} 0\\-1\\1 \end{bmatrix} = \begin{bmatrix} -1\\1\\0 \end{bmatrix}$$

$$\boldsymbol{w}^{T}(3) \begin{bmatrix} 1\\0\\1 \end{bmatrix} = -1 < 0, \quad \boldsymbol{w}(4) = \boldsymbol{w}(3)$$

11

LINEAR CLASSIFIERS The Perceptron Algorithm

Step 5.

$$\boldsymbol{w}^{T}(4) \begin{bmatrix} -1\\0\\1 \end{bmatrix} = 1 > 0, \quad \boldsymbol{w}(5) = \boldsymbol{w}(4)$$

Step 6.

$$\boldsymbol{w}^{T}(5) \begin{bmatrix} 0\\1\\1 \end{bmatrix} = 1 > 0, \quad \boldsymbol{w}(6) = \boldsymbol{w}(5)$$

Step 7.

$$\boldsymbol{w}^{T}(6) \begin{bmatrix} 0\\-1\\1 \end{bmatrix} = -1 < 0, \quad \boldsymbol{w}(7) = \boldsymbol{w}(6)$$

Since for four consecutive steps no correction is needed, all points are correctly classified and the algorithm terminates. The solution is $w = [-1, 1, 0]^T$. That is, the resulting linear classifier is $-x_1 + x_2 = 0$, and it is the line passing through the origin shown in Figure 3.4.





- > The network is called **perceptron or neuron**
- It is a learning machine that learns from the training vectors via the perceptron algorithm

The Perceptron Algorithm

- Choose w(0) randomly
- Choose ρ_0
- $\bullet t = 0$
- Repeat
 - $Y = \emptyset$
 - For i = 1 to N• If $\delta_{x_i} \boldsymbol{w}(t)^T \boldsymbol{x}_i \ge 0$ then $Y = Y \cup \{\boldsymbol{x}_i\}$
 - End {For}
 - $\boldsymbol{w}(t+1) = \boldsymbol{w}(t) \rho_t \sum_{\boldsymbol{x} \in Y} \delta_{\boldsymbol{x}} \boldsymbol{x}$
 - Adjust ρ_t
 - t = t + 1
- Until $Y = \emptyset$

LINEAR CLASSIFIERS The Perceptron Algorithm

Example: At some stage *t* the perceptron algorithm results in



Least Squares Methods

- If classes are linearly separable,
 the perceptron output results in ±1
- ➤ If classes are <u>NOT</u> linearly separable, we shall compute the weights W₁, W₂,..., W₀

so that the difference between

- The actual output of the classifier, $\underline{w}^T \underline{x}$, and
- The desired outputs, e.g. $\begin{cases} +1 \text{ if } \underline{x} \in \omega_1 \\ -1 \text{ if } \underline{x} \in \omega_2 \end{cases}$

- > SMALL, in the mean square error sense, means to choose \underline{W} so that the cost function
 - $J(\underline{w}) \equiv E[(\underline{y} \underline{w}^T \underline{x})^2]$ is minimum
 - $\underline{\hat{w}} = \arg\min_{\underline{w}} J(\underline{w})$
 - *y* the corresponding desired responses (targets)

LINEAR CLASSIFIERS Least Squares Methods Mean Square Error Estimation

> Minimizing $J(\underline{w})$ w.r. to \underline{w} results in :

$$\frac{\partial J(\underline{w})}{\partial \underline{w}} = \frac{\partial}{\partial \underline{w}} E[(y - \underline{w}^T x)^2] = 0$$
$$= 2E[\underline{x}(y - \underline{x}^T \underline{w})] \Rightarrow$$
$$E[\underline{x} \underline{x}^T] \underline{w} = E[\underline{x} y] \Rightarrow$$
$$\underline{\hat{w}} = R_x^{-1} E[\underline{x} y]$$

where R_x is the autocorrelation matrix

$$R_{x} \equiv E[\underline{x}\underline{x}^{T}] = \begin{bmatrix} E[x_{1}x_{1}] & E[x_{1}x_{2}]... & E[x_{1}x_{l}] \\ & \\ E[x_{1}x_{1}] & E[x_{1}x_{2}]... & E[x_{l}x_{l}] \end{bmatrix}$$

and
$$E[\underline{x}y] = \begin{bmatrix} E[x_{1}y] \\ ... \\ E[x_{l}y] \end{bmatrix}$$
 is the crosscorrelation vector

LINEAR CLASSIFIERS
Least Squares Methods
Mean Square Error Estimation

Multi-class generalization

• The goal is to compute *M* linear discriminant functions:

$$g_i(\underline{x}) = \underline{w}_i^T \underline{x}$$

according to the MSE.

• Adopt as desired responses y_i :

 $y_i = 1$ if $\underline{x} \in \omega_i$ $y_i = 0$ otherwise

• Let

$$\underline{y} = \begin{bmatrix} y_1, y_2, \dots, y_M \end{bmatrix}^T$$

• and the matrix

$$W = \left[\underline{w}_1, \underline{w}_2, \dots, \underline{w}_M\right]$$

LINEAR CLASSIFIERS
Least Squares Methods
Mean Square Error Estimation

• The goal is to compute *W*:

$$\hat{W} = \arg\min_{W} E\left[\left\|\underline{y} - W^T \underline{x}\right\|^2\right] = \arg\min_{W} E\left[\sum_{i=1}^{M} \left(y_i - \underline{w}_i^T \cdot \underline{x}\right)^2\right]$$

• The above is equivalent to a number *M* of MSE minimization problems. That is:

Design each \underline{w}_i so that its desired output is 1 for $\underline{x} \in \omega_i$ and 0 for any other class.

- Remark: The MSE criterion belongs to a more general class of cost function with the following important property:
 - The value of $g_i(\underline{x})$ is an estimate, in the MSE sense, of the a-posteriori probability $P(\omega_i | \underline{x})$, provided that the desired responses used during training are $y_i = \overline{1}, \underline{x} \in \omega_i$ and 0 otherwise.

Sum of Error Squares Estimation

SMALL in the sum of error squares sense means

$$> \frac{\partial J(\underline{w})}{\partial \underline{w}} = \frac{\partial}{\partial \underline{w}} \sum_{i=1}^{N} (y_i - \underline{w}^T \underline{x}_i)^2 = 0 \Longrightarrow$$

$$(\sum_{i=1}^{N} \underline{x}_{i} \underline{x}_{i}^{T}) \underline{w} = \sum_{i=1}^{N} \underline{x}_{i} y_{i}$$

✤ Pseudoinverse Matrix➢ Define

$$X = \begin{bmatrix} \underline{x}_{1}^{T} \\ \underline{x}_{2}^{T} \\ \dots \\ \underline{x}_{N}^{T} \end{bmatrix} \quad (an \ N \times l \ matrix) \qquad \underline{y} = \begin{bmatrix} y_{1} \\ \dots \\ y_{N} \end{bmatrix} \quad corresponding \\ desired \ responses$$

>
$$X^T = [\underline{x}_1, \underline{x}_2, ..., \underline{x}_N]$$
 (an $l \times N$ matrix)

$$X^T X = \sum_{i=1}^N \underline{x}_i \underline{x}_i^T$$
$$X^T \underline{y} = \sum_{i=1}^N \underline{x}_i y_i$$

LINEAR CLASSIFIERS Least Squares Methods Sum of Error Squares Estimation

Thus

$$\left(\sum_{i=1}^{N} \underline{x}_{i}^{T} \underline{x}_{i}\right) \underline{\hat{w}} = \left(\sum_{i=1}^{N} \underline{x}_{i} y_{i}\right)$$
$$\left(X^{T} X\right) \underline{\hat{w}} = X^{T} \underline{y} \Longrightarrow$$
$$\underline{\hat{w}} = (X^{T} X)^{-1} X^{T} \underline{y}$$
$$= X^{\neq} \underline{y}$$

$$X^{\#} \equiv (X^T X)^{-1} X^T \qquad \mathbf{I}$$

Pseudoinverse of X

Assume $N = l \implies X$ square and invertible. Then

$$(X^T X)^{-1} X^T = X^{-1} X^{-T} X^T = X^{-1} \Longrightarrow$$

$$X^{\#} = X^{-1}$$

LINEAR CLASSIFIERS
Least Squares Methods
Sum of Error Squares Estimation

Assume N > l. Then, in general, there is no solution to satisfy all equations simultaneously:

$$X \underline{w} = \underline{y} : \begin{cases} \underline{x}_{1}^{T} \underline{w} = y_{1} \\ \underline{x}_{2}^{T} \underline{w} = y_{2} \\ \dots \\ \underline{x}_{N}^{T} \underline{w} = y_{N} \end{cases} N$$

N equations > l unknowns

The "solution" $\underline{w} = X^{\#} \underline{y}$ corresponds to the minimum sum of squares solution

LINEAR CLASSIFIERS
Least Squares Methods
Sum of Error Squares Estimation



26

$$w_1 x_1 + w_2 x_2 + w_0 = 0$$

$$X^{T}X = \begin{bmatrix} 2.8 & 2.24 & 4.8 \\ 2.24 & 2.41 & 4.7 \\ 4.8 & 4.7 & 10 \end{bmatrix}, X^{T}\underline{y} = \begin{bmatrix} -1.6 \\ 0.1 \\ 0.0 \end{bmatrix}$$

$$\underline{w} = (X^T X)^{-1} X^T \underline{y} = \begin{bmatrix} -3.13 \\ 0.24 \\ 1.34 \end{bmatrix}$$

27

LINEAR CLASSIFIERS
Least Squares Methods
Mean Square Error Regression

- ➤ Mean square error regression: Let $y \in \Re^M$, $\underline{x} \in \Re^\ell$ be jointly distributed random vectors with a joint pdf $p(\underline{x}, y)$
 - The goal: Given the value of \underline{x} estimate the value of \underline{y} . In the pattern recognition framework, given \underline{x} one wants to estimate the respective label $y = \pm 1$.
 - The MSE estimate \hat{y} of \underline{y} given \underline{x} is defined as:

$$\underline{\hat{y}} = \arg\min_{\tilde{y}} E\left[\left\|y - \tilde{y}\right\|^{2}\right]$$

• It turns out that:

$$\underline{\hat{y}} = E[\underline{y} \mid \underline{x}] \equiv \int_{-\infty}^{+\infty} \underline{y} p(\underline{y} \mid \underline{x}) d\underline{y}$$

The above is known as the regression of \underline{y} given \underline{x} and it is, in general, a non-linear function of \underline{x} . If $p(\underline{x}, \underline{y})$ is Gaussian the MSE regressor is linear.

The Bias – Variance Dilemma

A classifier $g(\underline{x})$ is a learning machine that tries to predict the class label y of \underline{x} . In practice, a finite data set D is used for its training. Let us write $g(\underline{x}; D)$. Observe that:

- For some training sets, $D = \{(y_i, \underline{x}_i)\}_{i=1}^N$, the training may result to good estimates, for some others the result may be worse.
- The average performance of the classifier can be tested against the MSE optimal value, in the mean squares sense, that is:

$$E_{D}\left[\left(g\left(\underline{x};D\right)-E\left[y\mid\underline{x}\right]\right)^{2}\right]$$

where E_D is the mean over all possible data sets D.

• The above is written as:

$$E_{D}\left[\left(g\left(\underline{x};D\right)-E\left[y\mid\underline{x}\right]\right)^{2}\right]=\left(E_{D}\left[g\left(\underline{x};D\right)\right]-E\left[y\mid\underline{x}\right]\right)^{2}+E_{D}\left[\left(g\left(\underline{x};D\right)-E_{D}\left[g\left(\underline{x};D\right)\right]\right)^{2}\right]$$

- In the above, the first term is the contribution of the bias and the second term is the contribution of the variance.
- For a finite *D*, there is a trade-off between the two terms. Increasing bias it reduces variance and vice verse. This is known as the bias-variance dilemma.
- Using a complex model results in low-bias but a high variance, as one changes from one training set to another. Using a simple model results in high bias but low variance.



FIGURE 3.8

The data points are spread around the f(x) curve. The line g(x) = 0 exhibits zero variance but high bias. The high degree polynomial curve, $g_1(x) = 0$, always passes through the training points and leads to low bias (zero bias at the training points) but to high variance.

31

***** LOGISTIC DISCRIMINATION

≻ Let an *M*-class task, $\omega_1, \omega_2, ..., \omega_M$. In logistic discrimination, the logarithm of the likelihood ratios are modeled via linear functions, i.e.,

$$\ln\left(\frac{P\left(\omega_{i} \mid \underline{x}\right)}{P\left(\omega_{M} \mid \underline{x}\right)}\right) = w_{i,0} + \underline{w}_{i}^{T} \underline{x}, \quad i = 1, 2, ..., M - 1$$

➤ Taking into account that

$$\sum_{i=1}^{M} P(\omega_i \mid \underline{x}) = 1$$

it can be easily shown that the above is equivalent with modeling posterior probabilities as:

LINEAR CLASSIFIERS
Logistic Discrimination

$$P(\omega_{M} \mid \underline{x}) = \frac{1}{1 + \sum_{i=1}^{M-1} \exp\left(w_{i,0} + \underline{w}_{i}^{T} \underline{x}\right)}$$
$$P(\omega_{i} \mid \underline{x}) = \frac{\exp\left(w_{i,0} + \underline{w}_{i}^{T} \underline{x}\right)}{1 + \sum_{i=1}^{M-1} \exp\left(w_{i,0} + \underline{w}_{i}^{T} \underline{x}\right)}, i = 1, 2, \dots, M - 1$$

 \succ For the two-class case it turns out that

$$P(\omega_2 \mid \underline{x}) = \frac{1}{1 + \exp(w_0 + \underline{w}^T \underline{x})}$$
$$P(\omega_1 \mid \underline{x}) = \frac{\exp(w_0 + \underline{w}^T \underline{x})}{1 + \exp(w_0 + \underline{w}^T \underline{x})}$$

LINEAR CLASSIFIERS Logistic Discrimination

The unknown parameters $\underline{w}_i, w_{i,0}, i = 1, 2, ..., M-1$ are usually estimated by maximum likelihood arguments.

> Logistic discrimination is a useful tool, since

- it allows linear modeling, and
- at the same time ensures posterior probabilities to add to one.

LINEAR CLASSIFIERS > Support Vector Machines

Support Vector Machines (SVM)

The goal: Given two linearly separable classes, design the classifier $g(x) = w^T x + w_0 = 0$

that leaves the maximum margin from both classes



> Margin: Each hyperplane is characterized by

- Its direction in space, i.e., \underline{W}
- Its position in space, i.e., W_0



• For EACH direction, <u>w</u>, choose the hyperplane that leaves the SAME distance from the nearest points from each class. The margin is twice this distance.

- > The distance of a point $\underline{\hat{x}}$ from a hyperplane $z_{\hat{x}} = \frac{g(\underline{\hat{x}})}{\|\underline{w}\|}$
- Scale, $\underline{w}, \underline{w}_0$, so that at the nearest points from each class the discriminant function is ± 1 :

$$|g(\underline{x})| = 1 \{g(\underline{x}) = +1 \text{ for } \omega_1 \text{ and } g(\underline{x}) = -1 \text{ for } \omega_2 \}$$

 \succ Thus the margin is given by

$$\frac{1}{\|\underline{w}\|} + \frac{1}{\|\underline{w}\|} = \frac{2}{\|w\|}$$

> Also, the following is valid
$$\begin{cases} \underline{w}^T \underline{x} + w_0 \ge 1 & \forall \underline{x} \in \omega_1 \\ \underline{w}^T \underline{x} + w_0 \le -1 & \forall \underline{x} \in \omega_2 \end{cases}$$

SVM (linear) classifier

$$g(\underline{x}) = \underline{w}^T \underline{x} + w_0$$

➤ Minimize

$$J(\underline{w}) = \frac{1}{2} \left\| \underline{w} \right\|^2$$

Subject to

$$y_i(\underline{w}^T \underline{x}_i + w_0) \ge 1, \ i = 1, 2, ..., N$$

$$y_i = 1 \quad \text{, for } \underline{x}_i \in \omega_1,$$

$$y_i = -1, \text{ for } \underline{x}_i \in \omega_2$$

> The above is justified since by minimizing $\|\underline{w}\|$

the margin
$$\frac{2}{\|w\|}$$
 is maximized

LINEAR CLASSIFIERS > Support Vector Machines

The above is a quadratic optimization task, subject to a set of linear inequality constraints. The Karush-Kuhn-Tucker conditions state that the minimizer satisfies:

• (1)
$$\frac{\partial}{\partial \underline{w}} L(\underline{w}, w_0, \underline{\lambda}) = \underline{0}$$

• (2)
$$\frac{\partial}{\partial w_0} L(\underline{w}, w_0, \underline{\lambda}) = 0$$

• (3)
$$\lambda_i \ge 0, i = 1, 2, ..., N$$

• (4)
$$\lambda_i \left[y_i (\underline{w}^T \underline{x}_i + w_0) - 1 \right] = 0, i = 1, 2, ..., N$$

• Where
$$L(\bullet, \bullet, \bullet)$$
 is the Lagrangian
 $L(\underline{w}, w_0, \underline{\lambda}) \equiv \frac{1}{2} \underline{w}^T \underline{w} - \sum_{i=1}^N \lambda_i [y_i(\underline{w}^T \underline{x}_i + w_0)]$

39

The solution: from the above, it turns out that

•
$$\underline{W} = \sum_{i=1}^{N} \lambda_i y_i \underline{x}_i$$

•
$$\sum_{i=1}^N \lambda_i y_i = 0$$

> Remarks:

• The Lagrange multipliers can be either zero or positive. Thus,

$$- \underline{w} = \sum_{i=1}^{N_s} \lambda_i y_i \underline{x}_i$$

where $N_s \leq N$, corresponding to positive Lagrange multipliers

- From constraint (4) above, i.e., $\lambda_i [y_i (w^T x_i + w_0) - 1] = 0, \quad i = 1, 2, ..., N$

the vectors contributing to \underline{W} satisfy

$$\underline{w}^T \underline{x}_i + w_0 = \pm 1$$

- These vectors are known as SUPPORT VECTORS and are the closest vectors, from each class, to the classifier.
- Once \underline{w} is computed, w_0 is determined from conditions (4).
- The optimal hyperplane classifier of a support vector machine is UNIQUE.
- Although the solution is unique, the resulting Lagrange multipliers are not unique.

LINEAR CLASSIFIERS > Support Vector Machines

Dual Problem Formulation

- The SVM formulation is a convex programming problem, with
 - Convex cost function
 - Convex region of feasible solutions
- Thus, its solution can be achieved by its dual problem, i.e.,

- maximize
$$L(\underline{w}, w_0, \underline{\lambda})$$

- subject to $\underline{w} = \sum_{i=1}^{N} \lambda_i y_i \underline{x}_i$
 $\sum_{i=1}^{N} \lambda_i y_i = 0$
 $\underline{\lambda} \ge \underline{0}$

• Combine the above to obtain

- maximize
$$(\sum_{i=1}^{N} \lambda_i - \frac{1}{2} \sum_{ij} \lambda_i \lambda_j y_i y_j \underline{x}_i^T \underline{x}_j)$$

– subject to

$$\sum_{i=1}^{N} \lambda_{i} y_{i} = 0$$
$$\underline{\lambda} \ge \underline{0}$$

> Remarks:

• Support vectors enter via inner products

➢ Non-Separable classes



In this case, there is no hyperplane such that

$$\underline{w}^T \underline{x} + w_0 \ge 1, \quad \forall \underline{x}$$



• Recall that the margin is defined as twice the distance between the following two hyperplanes

$$\underline{w}^{T} \underline{x} + w_{0} = 1$$

and
$$\underline{w}^{T} \underline{x} + w_{0} = -1$$

The training vectors belong to <u>one</u> of <u>three</u> possible categories

1) Vectors outside the band which are correctly classified, i.e.,

 $y_i(\underline{w}^T \underline{x} + w_0) > 1$

2) Vectors inside the band, and correctly classified, i.e.,

$$0 \le y_i(\underline{w}^T \underline{x} + w_0) < 1$$

3) Vectors misclassified, i.e., $y_i(\underline{w}^T \underline{x} + w_0) < 0$

LINEAR CLASSIFIERS > Support Vector Machines

 \succ All three cases above can be represented as

$$y_i(\underline{w}^T \underline{x} + w_0) \ge 1 - \xi_i$$

1)
$$\rightarrow \xi_i = 0$$

- 2) $\rightarrow 0 < \xi_i \le 1$ 3) $\rightarrow 1 < \xi_i$

ξ_i are known as slack variables

LINEAR CLASSIFIERS > Support Vector Machines

 \succ The goal of the optimization is now two-fold

- Maximize margin
- Minimize the number of patterns with $\xi_i > 0$, One way to achieve this goal is via the cost

$$J(\underline{w}, w_0, \underline{\xi}) = \frac{1}{2} \left\| \underline{w} \right\|^2 + C \sum_{i=1}^N I(\xi_i)$$

where C is a constant and

$$I(\xi_i) = \begin{cases} 1 & \xi_i > 0 \\ 0 & \xi_i = 0 \end{cases}$$

• *I(.)* is not differentiable. In practice, we use an approximation

$$J(\underline{w}, w_0, \underline{\xi}) = \frac{1}{2} \left\| \underline{w} \right\|^2 + C \sum_{i=1}^N \xi_i$$

• Following a similar procedure as before we obtain

➢ KKT conditions

(1)
$$\underline{w} = \sum_{i=1}^{N} \lambda_i y_i \underline{x}_i$$

(2) $\sum_{i=1}^{N} \lambda_i y_i = 0$
(3) $C - \mu_i - \lambda_i = 0, i = 1, 2, ..., N$
(4) $\lambda_i [y_i (\underline{w}^T \underline{x}_i + w_0) - 1 + \xi_i] = 0, \quad i = 1, 2, ..., N$
(5) $\mu_i \xi_i = 0, \quad i = 1, 2, ..., N$
(6) $\mu_i, \lambda_i \ge 0, \quad i = 1, 2, ..., N$

> The associated dual problem

Maximize
$$\frac{\lambda}{2} \left(\sum_{i=1}^{N} \lambda_{i} - \frac{1}{2} \sum_{i,j} \lambda_{i} \lambda_{j} y_{i} y_{j} \underline{x}_{i}^{T} \underline{x}_{j} \right)$$
subject to
$$0 \le \lambda_{i} \le C, \ i = 1, 2, ..., N$$
$$\sum_{i=1}^{N} \lambda_{i} y_{i} = 0$$

<u>Remarks</u>: The only difference with the separable class case is the existence of *C* in the constraints

> Training the SVM

A major problem is the high computational cost. To this end, decomposition techniques are used. The rationale behind them consists of the following:

- Start with an arbitrary data subset (working set) that can fit in the memory. Perform optimization, via a general purpose optimizer.
- Resulting support vectors remain in the working set, while others are replaced by new ones (outside the set) that violate severely the KKT conditions.
- Repeat the procedure.
- The above procedure guarantees that the cost function decreases.
- Platt's SMO algorithm chooses a working set of two samples, thus analytic optimization solution can be obtained.



FIGURE 3.13

An example of two nonseparable classes and the resulting SVM linear classifier (full line) with the associated margin (dotted lines) for the values (a) C = 0.2 and (b) C = 1000. In the latter case, the location and direction of the classifier as well as the width of the margin have changed in order to include a smaller number of points inside the margin.

> Multi-class generalization

Although theoretical generalizations exist, the most popular in practice is to look at the problem as M two-class problems (one against all).

