

(1)

Alpaydin.

Note: Chp 10.

Linear Discrimination

4/27/2008

Recall discriminant functions $g_j(x), j=1, \dots, K$.

choose C_i of $g_i(x) = \max_{j=1}^K g_j(x)$

In previous chapters, the $g_j(x)$ were obtained from Bayesian decision theory by estimating probabilities $p(x|C_i), p(C_i)$ from the data.

Now discriminant-based classification where we assume a model for the discriminant

Model $g_i(x|\Phi_i)$ Φ_i set of parameters
 \rightarrow now need to learn the Φ_i from data.

In this lecture (Chp 10), we assume model
 $g_i(x|\underline{w}_i, w_{i0}) = \underline{w}_i^T x + w_{i0}$
 linear discriminant.

(recall linear discriminant arises from Bayes decision if the distribution $p(x|C_i)$ are Gaussian with identical covariance).

Simple generalizations

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

for other generalizations - see kernel trick.

(2)

Geometry of Linear Discrimination

Two classes : $g(\underline{x}) = g_1(\underline{x}) - g_2(\underline{x})$

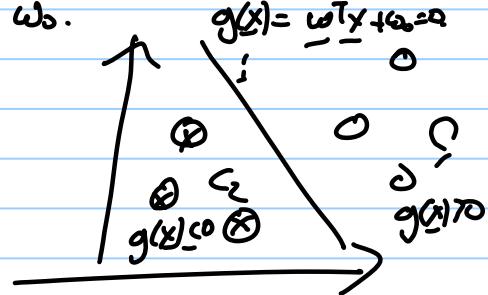
$$= (\underline{\omega}_1 - \underline{\omega}_2)^T \underline{x} + (\omega_{10} - \omega_{20})$$

$$= \underline{\omega}^T \underline{x} + \omega_0$$

choose C_1 if $g(\underline{x}) > 0$
 C_2 otherwise

$g(\underline{x}) = 0$ defines a hyperplane
it divides the space into

two regions $g(\underline{x}) > 0$ C_1
& $g(\underline{x}) < 0$ C_2 .



The normal to the hyperplane is \underline{w} .

closest distance to the origin is $\frac{|\omega_0|}{\|\omega\|}$

If $K > 2$ classes :

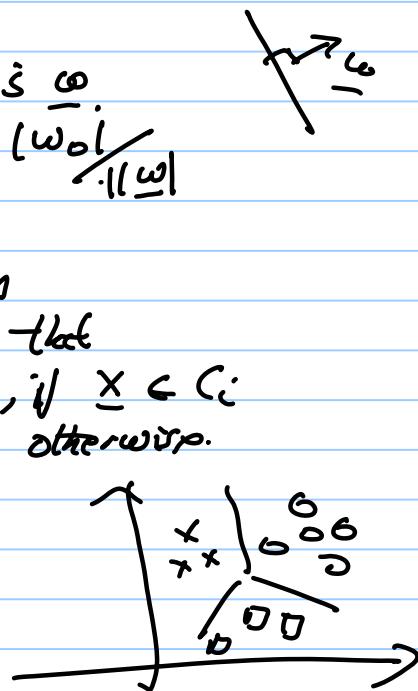
K discriminant functions

Adjust the parameters $\underline{\omega}_i, \omega_{i0}$, so that

$$g_i(\underline{x} | \underline{\omega}_i, \omega_{i0}) > 0, \text{ if } \underline{x} \in C_i$$

$$\leq 0, \text{ otherwise.}$$

In general, it will
not be possible to linearly
separate the data
(see support vector machine)



(3)

Support Vector Machines

Two class

$$X = \{\underline{x}^t, r^t\}, \quad r^t = +1, \text{ if } \underline{x}^t \in C_1, \quad r^t = -1, \text{ if } \underline{x}^t \in C_2.$$

Want to find $\underline{\omega}$ & w_0 s.t.

$$\begin{aligned}\underline{\omega}^T \underline{x}^t + w_0 &\geq 1 \quad \text{if } r^t = +1 \\ \underline{\omega}^T \underline{x}^t + w_0 &\leq -1 \quad \text{if } r^t = -1.\end{aligned}$$

re-express as $r^t (\underline{\omega}^T \underline{x}^t + w_0) \geq 1$

why "1"
require the
data to be
past the margin,

$$\min \frac{1}{2} |\underline{\omega}|^2 \quad \text{subject to} \quad r^t (\underline{\omega}^T \underline{x}^t + w_0) \geq 1, \forall t$$

Make the margin $\frac{1}{|\underline{\omega}|}$ as big as possible. (why? more likely to generalize)

Distance of \underline{x}^t to the plane $\underline{\omega}^T \underline{x} + w_0 = 0$

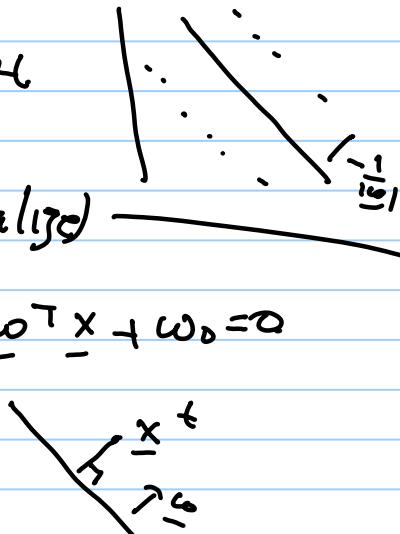
$$\text{line } \underline{x}(t) = \underline{x}^t - \hat{\lambda} \underline{\omega}$$

this hits the plane when.

$$\underline{\omega}^T \{ \underline{x}^t - \hat{\lambda} \underline{\omega} \} + w_0 = 0 \Rightarrow \hat{\lambda} = \frac{\underline{\omega}^T \underline{x}^t + w_0}{|\underline{\omega}|^2}$$

distance to plane is $|\hat{\lambda} \underline{\omega}| = \frac{|\underline{\omega}^T \underline{x}^t + w_0|}{|\underline{\omega}|}$.

sign-distance is $\frac{\underline{\omega}^T \underline{x}^t + w_0}{|\underline{\omega}|}$



4)

Formalize as constrained optimization problem
 → find separating plane with biggest margin.

$$L_p = \frac{1}{2} \|\underline{\omega}\|^2 - \sum_{t=1}^N \alpha^t [r^t (\underline{\omega}^T \underline{x}_t + \omega_0) - 1]$$

$\{\alpha^t\}$ Lagrange parameters
constrained so that $\alpha^t \geq 0$

minimized const. $\underline{\omega}$

maximized const. α^t .

note if $r^t (\underline{\omega}^T \underline{x}_t + \omega_0) - 1 > 0$

then we maximize α^t by setting $\alpha^t = 0$

convex optimization problem.

$$\frac{\partial L_p}{\partial \underline{\omega}} = 0 \Rightarrow \underline{\omega} = \sum_t \alpha^t r^t \underline{x}^t$$

$$\frac{\partial L_p}{\partial \alpha^t} = 0 \Rightarrow \sum_t \alpha^t r^t = 0.$$

Substituting back into L_p gives the dual

$$L_d = -\frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s (\underline{x}^t)^T \underline{x}^s + \sum_t \alpha^t.$$

goal maximize const. $\{\alpha^t\}$ subject
to $\sum_t \alpha^t r^t = 0$ & $\alpha^t \geq 0$, $\forall t$

can be solved by quadratic
optimization method. time complexity $O(N^3)$
space complexity $O(N^2)$.

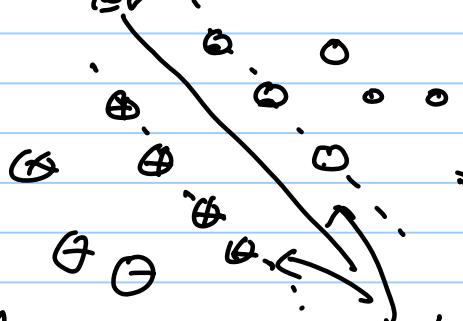
(5)

Support Vectors :

the solution is of form $\underline{\omega} = \sum_t \alpha^t r^t \underline{x}^t$.

But $\alpha^t = 0$ unless $r^t (\underline{\omega}^\top \underline{x}^t - \omega_0) = 0$
 which happens only if \underline{x}^t is on the margin
 - i.e. distance of \underline{x}^t to the hypersphere is $\frac{1}{\|\underline{\omega}\|}$

So only the data points on
 the margin determine the
 separating hyperplane.



To determine ω_0 , we only
 need to find one support vector \underline{x}^s
 - any data vector with $\alpha^s > 0$ - .

Then $\omega_0 = 1 - \underline{\omega}^\top \underline{x}^s$
 where $\underline{\hat{\omega}} = \sum_t \alpha^t r^t \underline{x}^t$

But, for stability, average over all the support
 vectors

$$\omega_0 = 1 - \frac{1}{|S|} \sum_{s \in S} \underline{\hat{\omega}}^\top \underline{x}^s$$

where $S = \{t : \alpha^t > 0\}$

(6)

Nonsparable Case

Define $\overbrace{\text{slack variables}}^{\xi^t \geq 0} \xi^t$ to allow a data point to move.

$$r^t (\underline{\omega}^T \underline{x}^t + \omega_0) \geq 1 - \xi^t$$

rewrite as $r^t (\underline{\omega}^T (\underline{x}^t + r^t \underline{\omega} \xi^t) + \omega_0) \geq 1$.

think of this as moving the data point to put it on the correct side of the margin.

penalize slackness by $\sum_t \xi^t$

$$\underset{\text{minimize}}{\text{minimize}} \quad L_p = \frac{1}{2} \|\underline{\omega}\|^2 + C \sum_t \xi^t - \sum_t \alpha^t \left\{ r^t (\underline{\omega}^T \underline{x}^t + \omega_0) - 1 - \xi^t \right\}$$

$$\text{subject w.r.t. } (\alpha^t), (\mu^t), \quad \text{if } \xi^t > 0, \text{ set } \mu^t = 0.$$

To obtain the dual $\frac{\partial L_p}{\partial \underline{\omega}} = 0 \quad \frac{\partial L_p}{\partial \xi^t} = 0$

$$L_d = \sum_t \alpha^t - \frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s (\underline{x}^t)^T \underline{x}^s$$

$$\text{subject to} \quad \sum_t \alpha^t r^t = 0 \quad 0 \leq \alpha^t \leq C, \forall t$$

Solve the dual - quadratic optimizer - to obtain the $\hat{\alpha}^t$, solve $\hat{\omega} = \sum_t \hat{\alpha}^t r^t \underline{x}^t$ to get $\hat{\omega}$, estimate ω_0 from support vectors as before.

(7)

The Kernel Trick.

Example Go from X -space to Z -space

$$z_1 = x_1, z_2 = x_2, z_3 = x_1^2, z_4 = x_2^2, z_5 = x_1 x_2$$

take $\underline{z} = (z_1 \dots z_5)$ as input.

linear function in z -space is non-linear in X -space

$$\underline{z} = \underline{\Phi}(x) \quad , \quad z_j = \phi_j(x), j=1 \dots k.$$

$$g(\underline{z}) = \underline{w}^\top \underline{\Phi}(x)$$

The solution of SVM is $\underline{w} = \sum_t \alpha^t r^t \underline{z}^t = \sum_t \alpha^t r^t \underline{\Phi}(x^t)$

The discriminant is $g(x) = \underline{w}^\top \underline{\Phi}(x) = \sum_t \alpha^t r^t \underline{\Phi}(x^t)^\top \underline{\Phi}(x)$

kernel trick - all that matters is

$$K(x^t, x) = \underline{\Phi}(x^t)^\top \underline{\Phi}(x)$$

we don't care what $\underline{\Phi}(x)$ is.

$$g(x) = \sum_t \alpha^t r^t K(x^t, x)$$

Kernels \rightarrow three main types.

$$\circ \text{polys} \quad K(x^t, x) = (x^t \cdot x + 1)^d$$

$$\circ \text{radial basis function} \quad K(x^t, x) = \exp\left\{-\frac{\|x^t - x\|}{\sigma}\right\}$$

$$\circ \text{sigmoid function} \quad K(x^t, x) = \tanh(2x^t \cdot x + 1)$$