

2009年10月13日
下午 10:50

Linear Models for Classification

C_k : K discrete class $k = 1, \dots, K$

usually C_k are disjoint

decision boundaries (D-1)-dim hyperplanes
within D-dim input space

input space is divided into decision regions

linearly separable

↳ number of classes
↳ dimensions

binary-class $t \in \{0, 1\}$
multi-class $t = (0, 1, 0, 0, 0)^T$ / 1-of-k scheme

different approaches

- ① discriminant functions
- ② generative $p(C_k|x) = \frac{p(x|C_k) p(C_k)}{p(x)}$
- ③ directly model the conditional $p(C_k|x)$

2009年10月20日
下午 03:22

Discriminant Functions

linear discriminant

$$y(x) = W^T x + w_0 \leftarrow \text{bias} \quad \left| \begin{array}{l} \text{negative bias} \\ \equiv \text{threshold} \end{array} \right.$$

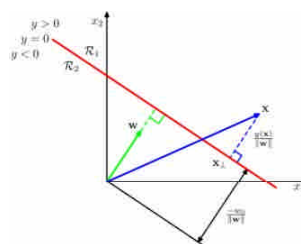
$$\begin{cases} C_1, & \text{if } y(x) > 0, \\ C_2, & \text{otherwise.} \end{cases}$$

decision plane : $y(x) = 0$

W is orthogonal to every vector lying within the decision surface

$$y(x) = W^T x + w_0 = 0$$

$$\rightarrow \frac{W^T x}{\|W\|} = -\frac{w_0}{\|W\|}$$



$$x = x_{\perp} + r \frac{W}{\|W\|}$$

$$\begin{cases} y(x) = W^T x + w_0 \\ y(x_{\perp}) = W^T x_{\perp} + w_0 = 0 \end{cases}$$

$$y(x) = W^T x + w_0 = W^T x_{\perp} + w_0 + r \frac{W^T W}{\|W\|}$$

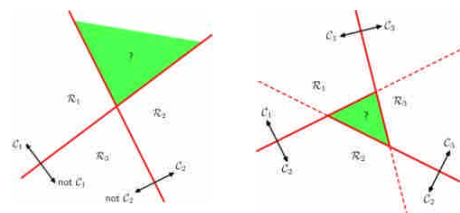
$$y(x) = r \|W\|$$

$$r = \frac{y(x)}{\|W\|}$$

Multiple Classes

one-versus-the-rest classifier
one-versus-one classifier

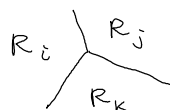
ambiguous region



alternative way

$$y_k(x) = w_k^T x + w_{k0}$$

$x \in C_k$ if $y_k(x) > y_j(x)$
for all $j \neq k$



convex region

$$\hat{x} = \lambda x_A + (1-\lambda)x_B \quad 0 \leq \lambda \leq 1$$

$$\Rightarrow y_k(\hat{x}) = \lambda y_k(x_A) + (1-\lambda)y_k(x_B)$$

$$y_k(x_A) > y_j(x_A), y_k(x_B) > y_j(x_B)$$

for all $j \neq k$

$$\Rightarrow y_k(\hat{x}) > y_j(\hat{x})$$

Least Squares for Classification

K binary coding

$$y_k(x) = w_k^T x + w_{k0}$$

$$y(x) = \tilde{w}^T \tilde{x}$$

$$y_k = \tilde{w}_k^T \tilde{x}$$

$$\tilde{w} = [\tilde{w}_1, \tilde{w}_2, \dots, \tilde{w}_K]$$

dummy input

$$x_0 = 1$$

$$\tilde{w} = (w_0, w)$$

$$\tilde{x} = (x_0, x)$$

$\{x_n, t_n\}$ training data

$$\tilde{X} = \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_N^T \end{bmatrix} \quad T = \begin{bmatrix} t_1^T \\ \vdots \\ t_N^T \end{bmatrix}$$

sum-of-squares error function

$$E_D(\tilde{w}) = \frac{1}{2} \text{Tr} \{ (\tilde{X} \tilde{w} - T)^T (\tilde{X} \tilde{w} - T) \}$$

set the derivatives w.r.t. \tilde{w} to zero

$$\tilde{w}^* = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T T = \tilde{X}^+ T$$

$$y(x) = \tilde{w}^T \tilde{x} = T^T (\tilde{X}^+)^T \tilde{x}$$

closed-form
solution

Problems: ① outliers Fig 4.4

② assume a Gaussian conditional
binary target vector t distribution
not Gaussian

2009年10月20日
下午 08:11

Fisher's linear discriminant or (LDA)

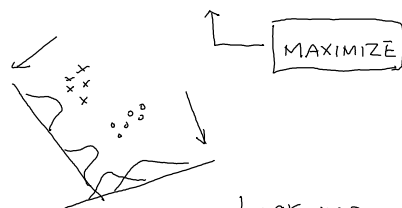
N_1 points in C_1 N_2 points in C_2

$$m_1 = \frac{1}{N_1} \sum_{n \in C_1} x_n, \quad m_2 = \frac{1}{N_2} \sum_{n \in C_2} x_n$$

projected onto w

$$m_2 - m_1 = w^T (m_2 - m_1)$$

$$m_k = w^T m_k$$



$$\sum_i w_i^2 = 1$$

Lagrange multiplier

$$\begin{aligned} &\text{maximize } w^T (m_2 - m_1) \\ &\text{subject to } w^T w = 1 \end{aligned}$$

$$w \propto (m_2 - m_1)$$

$$\begin{aligned} &(m_2 - m_1) + 2\lambda w \\ &= 0 \end{aligned}$$

NOT GOOD ENOUGH

MINIMIZE VARIANCE

$$s_k^2 = \sum_{n \in C_k} (y_n - m_k)^2 \quad | \quad y_n = w^T x_n$$

total within-class variance $s_1^2 + s_2^2$

2009年10月20日
下午 08:52

$$J(w) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} \quad \text{Fisher criterion}$$

$$J(w) = \frac{w^T S_B w}{w^T S_W w}$$

between-class scatter matrix: $S_B = (m_2 - m_1)(m_2 - m_1)^T$

(total) within-class scatter matrix:

$$S_W = \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T + \sum_{n \in C_2} (x_n - m_2)(x_n - m_2)^T$$

$$\frac{\partial}{\partial w} J(w) = 0$$

$$\frac{2 S_B w}{w^T S_W w} + \frac{w^T S_B w}{(w^T S_W w)^2} (-2 S_W w) = 0$$

$$\underbrace{(w^T S_B w)}_{\text{scalar}} S_W w = \underbrace{(w^T S_W w)}_{\text{scalar}} S_B w$$

$$\text{since } S_B w = (m_2 - m_1)(m_2 - m_1)^T w$$

$S_B w$ is in the direction of $(m_2 - m_1)$

we have

$$w \propto S_W^{-1} (m_2 - m_1)$$

RELATION TO LEAST SQUARES

2009年10月20日

下午 09:06

least-squares approach: making the model predictions as close as possible to a set of target values

Fisher criterion: maximum class separation in the output space

Fisher criterion as a special case of least squares

Let the target values be $\frac{N}{N_1}$ and $-\frac{N}{N_2}$
 \uparrow \uparrow
 $|C_1|$ $|C_2|$

sum-of-squares error function

$$E = \frac{1}{2} \sum_{n=1}^N (w^T x_n + w_0 - t_n)^2$$

$$\frac{\partial E}{\partial w_0} = 0 \quad \sum_{n=1}^N (w^T x_n + w_0 - t_n) = 0$$

$$w_0 = -w^T m \quad \leftarrow \begin{cases} \sum t_n \\ = N_1 \frac{N}{N_1} - N_2 \frac{N}{N_2} \\ = 0 \end{cases}$$

$$\frac{\partial E}{\partial w} = 0 \quad \sum_{n=1}^N (w^T x_n + w_0 - t_n) x_n = 0$$

$$\sum_{n=1}^N w^T x_n x_n + N m w_0 - \sum_{n=1}^N t_n x_n = 0$$

after one-page derivation

$$\left(S_W + \frac{N_1 N_2}{N} S_B \right) w = N (m_1 - m_2) \quad \left| \begin{aligned} m &= \frac{1}{N} \sum_{n=1}^N x_n \\ &= \frac{1}{N} (N_1 m_1 + N_2 m_2) \end{aligned} \right.$$

$$S_B w \propto (m_2 - m_1) \Rightarrow w \propto S_W^{-1} (m_2 - m_1) \quad \text{same as Fisher}$$

FISHER'S DISCRIMINANT FOR MULTIPLE CLASSES

2009年10月20日

下午 09:24

$K > 2$

assume that the dimensionality D of the input space is greater than the number K of classes

we want to get D' linear features

$$y_k = w_k^T x \quad k = 1, \dots, D'$$

weight vectors $\{w_k\}$

$$y = W^T x$$

within-class $S_W = \sum_{k=1}^K S_k$

$$S_k = \sum_{n \in C_k} (x_n - m_k)(x_n - m_k)^T$$

$$m_k = \frac{1}{N_k} \sum_{n \in C_k} x_n$$

total between-class $S_T = \sum_{n=1}^N (x_n - m)(x_n - m)^T$

$$m = \frac{1}{N} \sum_{k=1}^K N_k m_k \quad N = \sum_{k=1}^K N_k$$

$$S_T = S_W + S_B$$

$$S_B = \sum_{k=1}^K N_k (m_k - m)(m_k - m)^T$$

2009年10月20日
下午 09:32

in the projected D' -dimensional y -space

$$S_W = \sum_{k=1}^K \sum_{n \in G_k} (y_n - \mu_k)(y_n - \mu_k)^T$$

$$S_B = \sum_{k=1}^K N_k (\mu_k - \mu)(\mu_k - \mu)^T$$

$$\mu_k = \frac{1}{N_k} \sum_{n \in G_k} y_n, \quad \mu = \frac{1}{N} \sum_{k=1}^K N_k \mu_k$$

criterion?

$$J(W) = \text{Tr} \{ S_W^{-1} S_B \}$$

\Downarrow

$$J(W) = \text{Tr} \{ (W S_W W^T)^{-1} (W S_B W^T) \}$$

other criteria

$$J_2 = \ln |S_W^{-1} S_B| = \ln |S_B| - \ln |S_W|$$

$$J_3 = \text{Tr} \{ S_B \} - \mu (\text{Tr} \{ S_W \} - C)$$

$$J_4 = \frac{\text{Tr} \{ S_B \}}{\text{Tr} \{ S_W \}}$$

How to find W that maximizes $J(W)$?

2009年10月20日
下午 09:44

$$J(W) = \text{Tr} \{ S_W^{-1} S_B \} = \text{Tr} \{ (W S_W W^T)^{-1} (W S_B W^T) \}$$

$$\frac{\partial J}{\partial W} = 0 \Rightarrow -2 S_W W S_W^{-1} S_B S_W^{-1} + 2 S_B W S_W^{-1} = 0$$

$$(S_W^{-1} S_B) W = W (S_W^{-1} S_B)$$

S_W, S_B simultaneously diagonalized by B

$$B^T S_B B = \Lambda, \quad B^T S_W B = I$$

If we apply the transformation B to y
the value of the criterion will not change

$$\text{Tr} \{ \tilde{S}_W^{-1} \tilde{S}_B \} = \text{Tr} \{ (B^T S_W B)^{-1} (B^T S_B B) \}$$

$$= \text{Tr} \{ B^{-1} S_W^{-1} S_B B \}$$

$$= \text{Tr} \{ S_W^{-1} S_B B B^{-1} \}$$

$$= \text{Tr} \{ S_W^{-1} S_B \}$$

$$S_B = B^{-T} \Lambda B^{-1}, \quad S_W = B^{-T} B^{-1}$$

$$(S_W^{-1} S_B) W = W (B^{-T} B^{-1})^{-1} (B^{-T} \Lambda B^{-1}) = W B \Lambda B^{-1}$$

$$(S_W^{-1} S_B) (W B) = (W B) \Lambda$$

\uparrow eigenvalues of $(S_W^{-1} S_B)$

solution $\left[\begin{array}{l} \text{Find the eigenvectors corresponding to} \\ \text{the } D' \text{ largest eigenvalues of } S_W^{-1} S_B \end{array} \right.$

2009年10月20日
下午 10:08

THE PERCEPTRON ALGORITHM

$$y(x) = f(w^T \phi(x))$$

$f(\cdot)$ nonlinear activation function

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases}$$

$$\phi_0(x) = 1 \quad \text{bias component}$$

$$t = +1 \quad \text{for class } C_1$$

$$t = -1 \quad \text{for class } C_2$$

Goal: find w such that $w^T \phi(x_n) t_n > 0$
for $n = 1, \dots, N$

perceptron criterion

$$E_p(w) = - \sum_{n \in N} w^T \phi(x_n) t_n$$

minimize the criterion using stochastic gradient descent

one input at a time

$$\text{if } w^{(c)T} \phi(x_n) t_n < 0$$

x_n is misclassified by $w^{(c)}$

$$\text{update } w^{(c+1)} = w^{(c)} - \eta \nabla E_p(w) = w^{(c)} + \eta \phi(x_n) t_n$$

η is the learning rate parameter

2009年10月20日
下午 10:17

$$\begin{aligned} E_p(w^{(c+1)}) &= -w^{(c+1)T} \phi_n t_n \\ &= -w^{(c)T} \phi_n t_n - (\phi_n t_n)^T (\phi_n t_n) \\ &< -w^{(c)T} \phi_n t_n = E_p(w^{(c)}) \end{aligned}$$

$E_p(w)$ decreases for (x_n, t_n)

(The algorithm converges in a finite number of steps)
if the data are linearly separable

PERCEPTRON CONVERGENCE THEOREM

Let $R = \max \|\phi(x_n)\|$

Suppose that there exists a vector w^* such that
 $\|w^*\| = 1$ and $t_n (w^{*T} \phi(x_n)) \geq \gamma$ for $n = 1, \dots, N$

then the number of mistakes made by the perceptron algorithm on the training data is at most $(\frac{R}{\gamma})^2$

Proof.

$$w^{(c+1)T} w^* = w^{(c)T} w^* + \eta (\phi(x_n) t_n)^T w^* \geq w^{(c)T} w^* + \eta \gamma$$

$$\text{apply the inequality } \tau \text{ times} \Rightarrow w^{(\tau+1)T} w^* \geq \tau \eta \gamma \quad (1)$$

$$\|w^{(z+1)}\|^2 = \|w^{(z)}\|^2 + \underbrace{2\eta t_n w^T \phi(x_n) + \eta^2 t_n^2 \|\phi(x_n)\|^2}_{< 0}$$

$$\begin{aligned} \|w^{(z+1)}\|^2 &\leq \|w^{(z)}\|^2 + \eta^2 \|\phi(x_n)\|^2 \quad | \quad t_n^2 = 1 \\ &\leq \|w^{(z)}\|^2 + \eta^2 R^2 \end{aligned}$$

apply the inequality τ times $\Rightarrow \|w^{(z+1)}\|^2 \leq \tau \eta^2 R^2$ (2)

Combine (1) and (2)

$$\|w^*\| \sqrt{\tau} \eta R \geq \|w^*\| \|w^{(z+1)}\| \geq \underbrace{w^{(z+1)T} w^*}_{\text{Cauchy-Schwarz inequality}} \geq \tau \eta \gamma$$

$$\|w^*\|^2 \tau \eta^2 R^2 \geq \tau^2 \eta^2 \gamma^2$$

$$\tau \leq \left(\frac{R}{\gamma}\right)^2$$

Problems of the perceptron algorithm

- ① cannot distinguish between a non separable problem and one that is simply slow to converge
- ② solution is not unique
- ③ not linearly separable \rightarrow never converges
- ④ no probabilistic output
- ⑤ binary classification
- ⑥ linear combination of fixed basis functions