# Classification#1
## Discriminant functions

Machine learning, 2021

Mansoor Rezghi

Department of Computer science, TMU

Ref: CB, AW, DU

# References

- Bishop: chap4

- R. O. Duda, P.E. Hart, D. G. Stork, Pattern Classification, Second Edition, Wiley, 2001.(DU)-chap5

* Linear regression: norm 1& norm 2

* Logistic Regression as discriminant function classifier

* Perceptron classifer

* Support vector machine classifier

supervised | Regres $\longrightarrow$ linn   min $\sum \text{loss}(\hat{y}_i, y_i)$

$\quad\quad\quad\quad \hookrightarrow$ opht $<$   | GD , SGD   Min-|   $x_{u+1} = x_u \neq \alpha \boxed{\nabla f_u}$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \alpha \boxed{g} \rightarrow$ subg —

Classifier   $\{(x_i, y_i)\}$   $(x, y)$

$f_u \longleftarrow$   $\downarrow$ Labe   $? y \in \{c_1-, c_u\}$

# Classification

Given Data set $x_i$ with known labels $t_i$

Predict the label of test data $x$

$\{C_i\} : C_u\}$

$C_i \quad g_i(u)$

$x \in C_\ell \qquad g_\ell(n) > g_i(n) \quad i \neq \ell$



$C_1 \rightarrow g_1$

$i$

$C_k \rightarrow g_k$

$g_j(n) = g_j(x)$

$\{x, g\}$ $\quad \bar{x}$

$P(x, C)$

$=\frac{q}{i}$

$P(\bar{x}, C_1)$

$\lor$

$P(\bar{x}, C_2)$

**Discriminant function**

For each class find discriminat function $g_i(x)$

Assign $x$ to class $C_l$ if $l = \arg\max_{1 i \leq K} g_i(x)$

$P(x,C)$ or $P(C|x)$

Generative Models ✓

Discriminative Models

$P(x, C)$

$P(x, C)$

$\bar{x} \quad P(\bar{x}, C_1)$

$\lor$

$P(\bar{x}, C_2)$

$$P(x, C_1) > P(x, C_2)$$

Assign $X \longrightarrow C_1$

$$P(A, B) = P(A|B) P(B)$$

$$P(C_1|x) P(x) > P(C_2|x) P(x)$$

$$P(C_1|x) = \frac{P(x|C_1) P(C_1)}{P(x)}$$

$$\boxed{P(C_1|x)} > P(C_2|x)$$

$$\frac{P(x|C_1) P(C_1)}{P(x)} > \frac{P(x|C_2) P(C_2)}{P(x)}$$

$$\boxed{P(x|C_1)} P(C_1) > P(x|C_2) P(C_2)$$

Generative.

$$P(x|C_1) k$$

$$P(x=1..|C_1) =$$

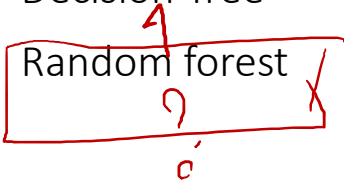$$\boxed{P(x)} = \frac{P(x|C_1) P(C_1)}{} + \frac{P(x|C_2) P(C_2)}{}$$

# Classification
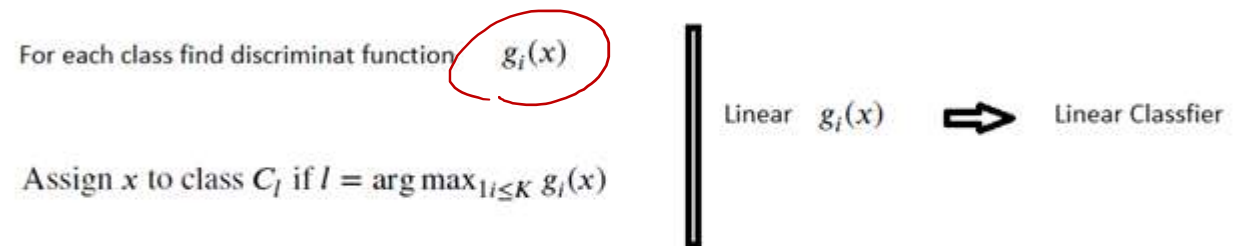
Generative models

Discriminative functions

- □ Least squares, Perceptron

- □ Logistic regression classifier

- □ Optimization based methods

- □ Support Vector machine

- □ ........

Discriminative models: Logistic Regression

- KNN

- Nave Bayes

- Decision Tree

- Random forest

# Discriminant function

For each class find discriminat function $g_i(x)$

Assign $x$ to class $C_l$ if $l = \arg\max_{1 \leq i \leq K} g_i(x)$
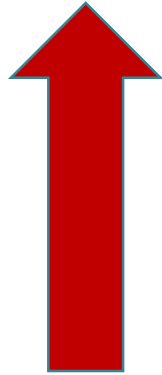
Linear $g_i(x)$ $\Rightarrow$ Linear Classfier

# Discriminant Function

Linear separable classes : classes can be separated via linear function(Linear decision surface)

$$g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0,$$

By increasing the dimension

Kevul

$\leftarrow \phi(x)$

$x$

Nonlinear separable classes : classes can not be separated via linear function(NonLinear decision surface)

?

# Nonlinear to Linear

Quadratic discriminant function

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

$\mathbf{x} \subset \mathbb{R}^d$

d=2

$$\bar{x} = \begin{pmatrix} x_0 \\ x_1 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \end{pmatrix} \qquad \bar{W} = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_{12} + w_{21} \\ w_{11} \\ w_{22} \end{pmatrix}$$

$$g(\mathbf{x}) = g(\bar{x}) = \bar{w}_0 + \sum_{i=1}^{5} \bar{w}_i \bar{x}_i$$

# Nonlinear to Linear separable(R. O. Duda)

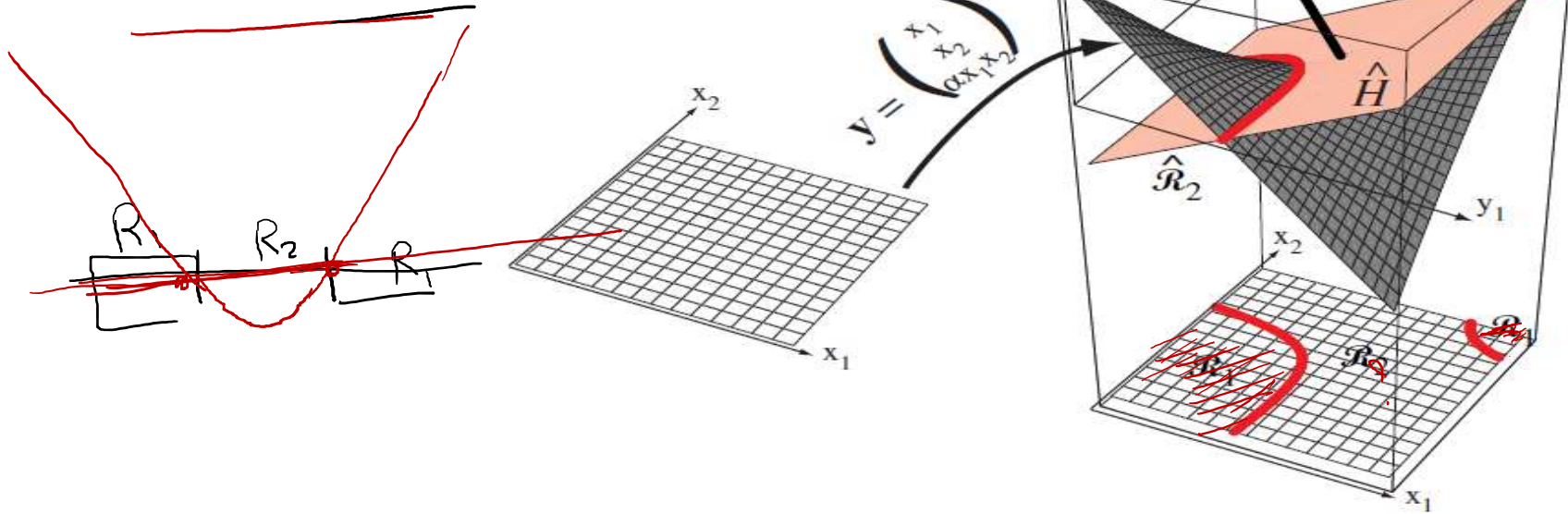$y = x^2$

$x \rightarrow \begin{pmatrix} x \\ x^2 \end{pmatrix}$

kernl.



Figure 5.6: The two-dimensional input space x is mapped through a polynomial function $f$ to y. Here the mapping is $y_1 = x_1$, $y_2 = x_2$ and $y_3 \propto x_1 x_2$. 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 x space, $\mathcal{R}_1$ is a not simply connected.

$$\boxed{x \in C_i \qquad g_i(x) \geqslant g_j(x)} \to a$$

$$\begin{cases} & x \in C_j \to (x_i, t_i) \qquad t_i = \begin{cases} 1 & x_i \in C_1 \\ -1 & x_i \in C_2 \end{cases} \qquad (x_i, t_i) \qquad t_i = \begin{cases} 1 & x \in C_1 \\ 0 & x \in C_2 \end{cases} \\ \\ & (x, y_i) \qquad y_i \in \{C_1 \dots C_k\} \to t_i = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \end{pmatrix} \to \ell^{th} \text{ position} \\ & \qquad\qquad\qquad\qquad x_i \in C_\ell \nearrow \end{cases}$$

/     $g(x) \mapsto t_i$

$x_i \to t_i$

$g(x_i) \to t_i$

$\begin{array}{l} C_1 \quad g_1 \qquad g_1(x) > g_2(x) \\ \\ C_2 \quad g_2 \quad h(x) = g_1(x) - g_2(x) \\ \\ \qquad\qquad h(x) > 0 \end{array}$

# Linear Regression Based classifier(Bishop-chap4)

Two- class

$\{(x_i, d_i)\}$

$d_i \in C_1 \Rightarrow t_i = \begin{cases} 1 & x_i \in C_1 \\ -1 & x_i \in C_2 \end{cases}$

$\in C_2$

316

$g_i(x)$

if $\dfrac{g(x_i) \cap t_i}{w^T x_i \cap t_i}$

$g(x) = w^T x$

$\dfrac{1}{N} \sum loss(w^T x_i, t_i) = \sum_i \dfrac{(w^T x_i - t_i)^r}{|w^T x_i - t_i|} \begin{cases} +\lambda \, \|w\|_2 \\ +\lambda \, \|w\|_1 \end{cases}$

$h(x) > 0 \to x \in C_1$

$< 0 \quad x \in C_2$

$(x_i, t_i)$

$\mathbb{R}^d$

$\mathbb{R}^{N+1}$

$\circ \circ \circ \circ \circ$

$x \times \times \times \times$

$x$

$(x(t))$ $\boxed{\text{Line Separable}}$ $\Gamma(h(x)) = \begin{cases} 1 \to \\ -1 \end{cases}$

$\circ\circ\circ\circ \quad \times\times\times\times\times$ $\mathbb{R}$

$\Leftarrow$

$C_2 \qquad C_1 \quad x \in \mathbb{R}^d$

$x_i$

$y = \dfrac{w^T x + w_0}{w^T x} = g h(x)$

$x \in C_1 \quad \boxed{h(x_i) \sim 1}$

$\forall x \; \boxed{h(x) > 0} \to x \in C_1$

$< 0 \to x \in C_2$

$x \in C_2 \quad h(x_i) \sim -1$

$\boxed{\Gamma(t) = \begin{cases} 1 & t > 0 \\ -1 & t < 0 \end{cases}}$

$\boxed{\Gamma(h(x))}$

$(x_i, \boxed{y_i})$  $y_i = \begin{cases} 1 & x \in C_1 \\ -1 & x \in C_2 \end{cases}$

$l_i$

$y_i = \begin{cases} 1 & x \in C_1 \\ 0 & x \in C_2 \end{cases}$

$\omega^T x + \omega_i \rangle$

$x > a$

$\omega^T x_i + \omega \cup t_i$

$f(x) = \omega^T x_i + \omega_.$

$\omega^T x + \omega_. \rangle 0 \Rightarrow$

$ax + b \cup t_i$

$C_1 \Leftarrow ax + b \rangle 0$

$\boxed{x > -\dfrac{b}{a}}$

$a_1 x_1 + a_2 x_2 + b \geqslant 0$

$\boxed{a_1 x_1 + a_1 x_2 = -b}$

# Regression Based Classifier:

Input $\{X_i\}$

$X_i \in \{C_1 \cdots C_k\}$

$X_i \longrightarrow t_i = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$

if $X_i \in C_\ell$  $t_i = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{pmatrix}$ ⟶ $\ell$ th position.

$$c_1, c_2 \dots, c_k$$

$$\downarrow \qquad \downarrow \qquad \downarrow$$

$$g_1 \qquad g_2 \qquad g_k$$

$$x \in A \quad \forall x \in C_\ell \quad g_\ell(x) > g_i(x)$$

$$i \neq \ell$$

$$g_1(x) = w_1^T x$$

$$g_2(x) = w_2^T x$$

$$\vdots$$

$$g_k(x) = w_k^T x$$

$$g(x) = \begin{pmatrix} g_1(x) \\ \vdots \\ g_k(x) \end{pmatrix} = \begin{pmatrix} w_1^T x \\ w_2^T x \\ w_k^T x \end{pmatrix} = W^T x$$

$$W \in R^{d \times k}$$

$$g(x) =$$

$$x \in C_2$$

$$R_x \qquad t = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$g(x_i) \cong t_i$$

$$\begin{pmatrix} w_1^T x_i \\ w_2^T x_i \\ w_k^T x_i \end{pmatrix} \sim \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$x_i \in \mathbb{R}^d \rightarrow t_i \in \mathbb{R}^k \qquad g(x) = W^T x \qquad \lambda \mathcal{R}(W)$$

$$W \in \mathbb{R}^{d \times k}$$

$$\sum_i loss(g(x_i), t_i) = \sum \| W^T x_i - t_i \|_2^2 \quad \bigg| + \begin{cases} \lambda \| W \|_F^2 \\ \lambda \| W \|_1 \longrightarrow \text{sparse} \\ \lambda \| W \|_{2,1} \end{cases}$$

$$= \sum \| W^T x_i - t_i \|_1$$

Rogressia



ω₂  not ω₂  
ω₁  
ω₁  
not ω₁  
ambiguous region  
ω₂  
ω₃  
not ω₄  
ω₄  
ω₄  
not ω₃  ω₃

ω₁  
ω₂  
ω₁  
ω₃  
ω₁  
ω₁  
ambiguous region  
ω₄  
ω₂  
ω₃  
ω₂  ω₃  
ω₄  
ω₂  ω₄  
ω₃  
ω₄

# Drawbacks of LS method

$l \; \sigma(t) = 0$

Regresion

Sensitive to

☐ Number of train data of each class

☐ Outlier

$\vec{w}^T \vec{x}_i \sim y_i = \{ \; 1$

$$\sum loss(\vec{w}^T \vec{x}_i, y_i) \stackrel{?}{=}$$

$$= \sum loss(\sigma(\vec{w}^T \vec{x}_i), y_i)$$

$l. \sigma(t) =$

$l. \sigma(t) = 1$

$t \to \infty$

Least squres

Bia

$\sigma(\vec{w}^T x_i) = \{ \begin{array}{l} 1 \\ -1 \end{array}$  $\vec{w}^T x_i$

$$\sigma(t) = \frac{1}{1 + \bar{e}^t}$$

$$\min \sum \text{loss}\left(\sigma(w^T x_i), t_i\right)$$

$$\frac{}{\sum \left(\sigma(w^T x_i) - t_i\right)^2}$$

$(x_i, t_i)$

$t_i = \begin{cases} 1 & x \in C_1 \\ 0 & x \in C_2 \end{cases}$

# Logistic Regression

$$\sum \left|\sigma(w^T x_i) - t_i\right|$$

$$\sigma(t) = \frac{1}{1 + \bar{e}^t}$$

- Deterministic viewpoint

- Statistical viewpoint Discriminative model

$$\sum \text{loss}\left(\sigma(w^T x_i), t_i\right)$$

$$\sigma(0) = \frac{1}{2}$$

$\sigma(w^T x_i) = .9 \quad , t_i = 1$

$l. \; \sigma(t) = 1$

$t \to \infty$

$\sigma(w^T x_i) = .1 \quad t_i = 1$

$l. \; \sigma(t) = 0$

$t \to -\infty$

x x x x

o o o o o

# Logistic Regression classifier-Two class

$$y(\mathbf{x}) = f\left(\mathbf{w}^{\mathrm{T}}\phi(\mathbf{x})\right) \quad f(a) = \sigma(a) = \frac{1}{1+\exp(-a)}$$

$$y(\mathbf{x}) = f\left(\mathbf{w}^{\mathrm{T}}\phi(\mathbf{x})\right) \quad f(a) = \sigma(a) = \frac{1}{1+\exp(-a)}$$

$$f_{\mathbf{w}}(\phi(\mathbf{x})) = f\left(\mathbf{w}^{\mathrm{T}}\phi(\mathbf{x})\right)$$

Non-convex objective function

$$J(\mathbf{w}) = \sum_{i=1}^{N}\left(y^{(i)} - f_{\mathbf{w}}(\phi(\mathbf{x}^{(i)}))\right)^2$$

$$J(\mathbf{w}) = \sum_{i=1}^{N} \text{Cost}\left(y^{(i)}, f_{\mathbf{w}}(\phi(\mathbf{x}^{(i)}))\right)$$

$y^{(i)}$

$\sigma(W^{\top}a_i) \simeq 1 \rightarrow \sigma$

$f_W(\phi(x^i)) \simeq 1$

$-\phi^{(i)}$

$\text{hand log} \quad s = \log\left(\pm\right)$

$t = e^s$

$$\text{Cost}\left(y^{(i)}, f_W(\phi(x^{(i)}))\right) =$$

$$\simeq -\left[y^{(i)} \ln f_W(\phi(x^{(i)})) + (1 - y^{(i)}) \ln(1 - f_W(\phi(x^{(i)})))\right]$$

**For** $y^{(i)} = 0$

$$f_W(\phi(x^{(i)})) \simeq 0 \implies \text{Cost}\left(y^{(i)}, f_W(\phi(x^{(i)}))\right) \simeq 0$$

$$f_W(\phi(x^{(i)})) \simeq 1 \implies \text{Cost}\left(y^{(i)}, f_W(\phi(x^{(i)}))\right) \simeq \text{Inf}$$

**For** $y^{(i)} = 1$

$$f_W(\phi(x^{(i)})) \simeq 1 \implies \text{Cost}\left(y^{(i)}, f_W(\phi(x^{(i)}))\right) \simeq 0$$

$$f_W(\phi(x^{(i)})) \simeq 0 \implies \text{Cost}\left(y^{(i)}, f_W(\phi(x^{(i)}))\right) \simeq \text{Inf}$$

-log(x)

-log(1-x)

$y^{(i)} = 0 \implies$

$y^{(i)} = 1$

**Convex** $\quad J(\mathbf{w}) = \sum_{i=1}^{N} \text{Cost}\left(y^{(i)}, f_W(\phi(x^{(i)}))\right) \rightarrow$

$$z^i \, ,$$

$$\text{if } y_i = 0 \xrightarrow{\text{اذا}} \sigma(\vec{w}^T x_i) \simeq 0$$

$$\text{if } y_i = 0 \qquad \boxed{\ln \frac{1}{1-\sigma(\vec{w}^T x_i)}}$$

$$\ln \frac{1}{\sigma(\vec{w}^T x_i)}$$

$$\sigma(\vec{w}^T x_i) \simeq 0 \qquad \ln \frac{1}{1-\sigma(\vec{w}^T x_i)} \simeq 0$$

$$\sigma(\vec{w}^T x_i) \simeq 1 \qquad \ln \frac{1}{1-\sigma(\vec{w}^T x_i)} \simeq \infty$$

$$\text{if } y_i = 1 \qquad \ln \frac{1}{\sigma(\vec{w}^T x_i)}$$

$$\sigma(\vec{w}^T x_i) \simeq 1 \rightarrow \ln \frac{1}{\sigma(\vec{w}^T x_i)} = 0$$

$$\sigma(\vec{w}^T x_i) \simeq 0 \qquad \ln \frac{1}{\sigma(\vec{w}^T x_i)} = \infty$$

$$\text{cost}\left(\sigma(w^T x_i), y_i\right) = \begin{cases} \ln \dfrac{1}{\sigma(w^T x_i)} & y_i = 1 \\[4mm] \ln \dfrac{1}{1 - \sigma(w^T x_i)} & y_i = 0 \end{cases}$$
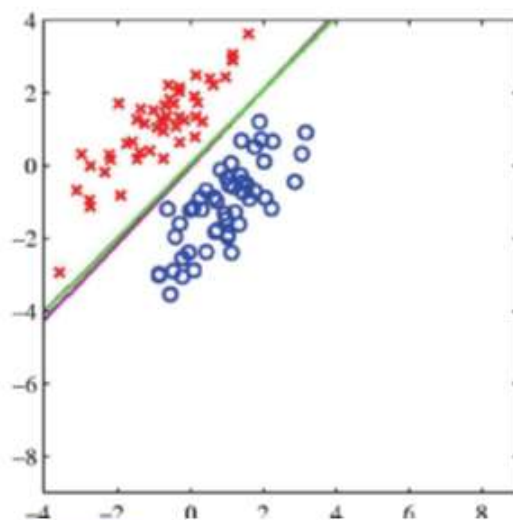
$$y_i \ln \frac{1}{\sigma(w^T x_i)} + (1 - y_i) \ln \frac{1}{1 - \sigma(w^T x_i)} = \text{loss}\left(\sigma(w^T x_i), y_i\right)$$

$$\min \sum_{i=1}^{N} \text{loss}\left(\sigma(w^T x_i), y_i\right) = \sum \overset{\text{Convex}}{y_i \ln \frac{1}{\sigma(w^T x_i)}} + (1 - y_i) \ln\left(\frac{1}{1 - \sigma(w^T x_i)}\right)$$

$$\max \sum y_i \ln \sigma(w^T x_i) + \underset{\text{Concave}}{(1 - y_i) \ln\left(1 - \sigma(w^T x_i)\right)}$$

logistic
regression

least squares
regression

Zemel

$\delta'(\quad)$

Logistic Regression-Multiclass objective function: Other Viewpoint

$$\max \; Loss(y, \hat{y}) = \Pi_j \hat{y}_j^{y_j},$$

$$\Downarrow$$

$$\max \; Loss(y, \hat{y}) = \ln \Pi_j \hat{y}_j^{y_j} = \sum_j y_j \ln \hat{y}_j \quad \Rightarrow \quad \max \; \sum_i Loss(y_i, \hat{y}_i) = \sum_i \sum_j y_{ji} \ln \hat{y}_{ji}$$

# Perceptron

Input data

$$\{(x_i, t_i)\}, \qquad t_i = \begin{cases} 1 & x_i \in C_1 \\ 0 & x_i \in C_2 \end{cases}$$

Loss functions:
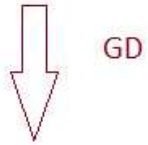
- $\text{Loss(w)} = |M|$,

  $M = \{\text{miss classified train data by } w^T x \text{ hyperplane}\}$

Model:

Find hyperplane $w^t x$ s.t

$$\forall x_i, \quad t_i(w^T x_i) \geq 0$$

- $\text{Loss(w)} = -\sum_{x_i \in M} t_i(w^T x_i)$ Perceptron

$\text{Loss(w)}=-\sum_{x_i \in M} t_i(w^T x_i)$ Perceptron

$\Downarrow$ GD

$w_{k+1} = w_k + \eta \sum_{x_n \in M} t_n x_n,$

$\Downarrow$ SGD

$w_{k+1} = w_k + \eta t_n x_n \qquad w_{k+1} = w_k + t_n x_n$

$\Rightarrow$

**Perceptron**

For i=1 ... maxiter

For train data like $(\phi_i, t_i), i = 1, .., N$ train data

$y_i = \mathbf{w}^T \phi_i$

If $t_i y_i < 0$ Then

$\quad w = w + t_i \phi_i$

However, the *perceptron convergence theorem* states that if there exists an exact solution (in other words, if the training data set is linearly separable), then the perceptron learning algorithm is guaranteed to find an exact solution in a finite number of steps. Proofs of this theorem can be found for example in Rosenblatt (1962),

Frank Rosenblatt. Principles of Neurodynamics: Perceptron and the Theory of Brain Mechanisms. Spartan Books, Washington, D.C., 1962

Albert B. J. Novikoff. On convergence proofs for perceptrons. In Proceedings of the Symposium on Mathematical Theory of Automata, volume 12, Brooklyn, New York, 1962.

http://ciml.info/dl/v0_99/ciml-v0_99-ch04.pdf

$$margin(\mathbf{D}, \boldsymbol{w}, b) = \begin{cases} \min_{(x,y) \in \mathbf{D}} y(\boldsymbol{w} \cdot \boldsymbol{x} + b) & \text{if } \boldsymbol{w} \text{ separates } \mathbf{D} \\ -\infty & \text{otherwise} \end{cases}$$

$$margin(\mathbf{D}) = \sup_{w,b} margin(\mathbf{D}, \boldsymbol{w}, b)$$

**Theorem 2** (Perceptron Convergence Theorem). *Suppose the perceptron algorithm is run on a linearly separable data set* $\mathbf{D}$ *with margin* $\gamma > 0$. *Assume that* $||\boldsymbol{x}|| \leq 1$ *for all* $\boldsymbol{x} \in \mathbf{D}$. *Then the algorithm will converge after at most* $\frac{1}{\gamma^2}$ *updates.*

- $\text{Loss(w)} = -\sum_{x_i \in M}(w^T x_i)^2$

- $\text{Loss(w)} = -\sum_{x_i \in M}\dfrac{(w^T x_i - b)^2}{\|w\|^2}$

# Multicategory Generalizations

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

$g_i(\mathbf{x}) = \mathbf{a}_i^t \mathbf{y} \quad i = 1, ..., c,$ where again $\mathbf{x}$ is assigned to $\omega_i$ if $g_i(\mathbf{x}) > g_j(\mathbf{x})$ for all $j \neq i$.

$\mathbf{y}_1, \mathbf{y}_2, ..., \mathbf{y}_n,$ with $n_i$ in the subset $\mathcal{Y}_i$ labelled $\omega_i$.

Duda: 5.12.1

For the data from linear separable multi-class, there exist a set of vectors $W_i, \quad i = 1, ..., k$ such that if $\phi_k \in \mathcal{C}_i$, then

$$\hat{w}_i^T \phi_k > \hat{w}_j^T \phi_k \text{ For } i \neq j$$

$$\hat{\alpha} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_c \end{bmatrix}$$

$$\eta_{ij} = \begin{bmatrix} 0 \\ \vdots \\ \mathbf{y} \\ 0 \\ \vdots \\ -\mathbf{y} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} \\ \\ \Rightarrow i \\ \\ \\ \Rightarrow j \\ \\ \\ \\ \end{matrix}$$

$\Rightarrow$

$$\hat{\alpha}^t \eta_{ij} > 0$$

$$\mathbf{a}_i^t \mathbf{y}_k - \mathbf{a}_j^t \mathbf{y}_k > 0$$
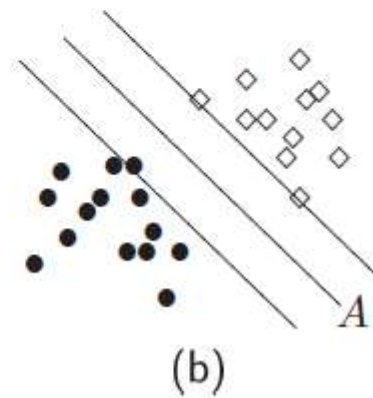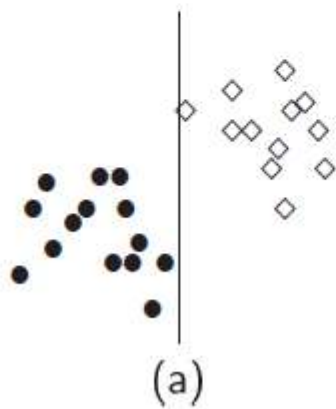
# Support Vector Machine

SVM

# Support vector machine(Maximum marginal classification) Webb (P.249)

$$g(x) = w^T x + w_0$$

$$w^T x + w_0 \begin{cases} > & 0 \\ < & 0 \end{cases} \Rightarrow x \in \begin{cases} \omega_1 \text{ with corresponding numeric value, } y_i = +1 \\ \omega_2 \text{ with corresponding numeric value, } y_i = -1 \end{cases}$$

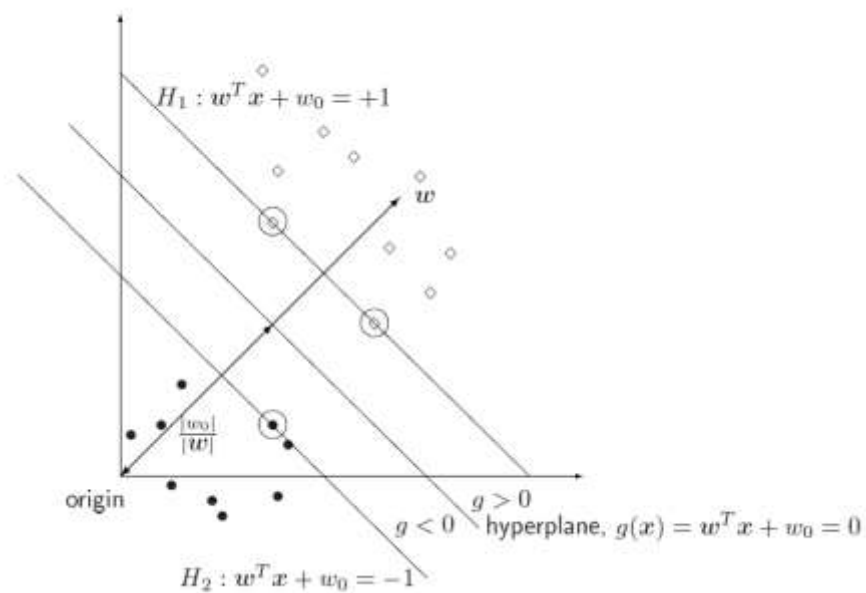$$y_i(w^T x_i + w_0) > 0 \text{ for all } i$$



(a)                    (b)

$$y_i(\boldsymbol{w}^T\boldsymbol{x}_i + w_0) \geq b$$

$\boldsymbol{w}^T\boldsymbol{x}_i + w_0 \geq +1 \quad$ for $y_i = +1$

$\boldsymbol{w}^T\boldsymbol{x}_i + w_0 \leq -1 \quad$ for $y_i = -1$

$H_1 : \boldsymbol{w}^T\boldsymbol{x} + w_0 = +1$

$\boldsymbol{w}$

min $\dfrac{1}{2}\|\mathbf{W}\|^2$

$y_i(\boldsymbol{w}^T\boldsymbol{x}_i + w_0) \geq 1 \quad i = 1, \ldots, n$

$\dfrac{|w_0|}{|w|}$

origin

$g > 0$

$g < 0$ hyperplane, $g(x) = \boldsymbol{w}^T\boldsymbol{x} + w_0 = 0$

$H_2 : \boldsymbol{w}^T\boldsymbol{x} + w_0 = -1$

# Convex Optimization

Optimality Conditions

$$\frac{\partial L_p}{\partial w} = w - \sum_{i=1}^{n} \alpha_i y_i x_i = 0$$

$$\frac{\partial L_p}{\partial w_0} = -\sum_{i=1}^{n} \alpha_i y_i = 0$$

$$y_i(x_i^T w + w_0) - 1 \geq 0$$

$$\alpha_i \geq 0$$

$$\alpha_i(y_i(x_i^T w + w_0) - 1) = 0$$

# Dual Problem

$$L_p = \frac{1}{2} w^T w - \sum_{i=1}^{n} \alpha_i (y_i (w^T x_i + w_0) - 1) \qquad \sum_{i=1}^{n} \alpha_i y_i = 0$$

$$w = \sum_{i=1}^{n} \alpha_i y_i x_i$$

**Dual Form**

$$\textbf{max} \quad L_D = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

$$\textbf{s.t}$$

$$\alpha_i \geq 0 \quad \sum_{i=1}^{n} \alpha_i y_i = 0$$

# Discrimination

$$w = \sum_{i=1}^{n} \alpha_i y_i x_i$$

W?

W0 : By Slackness

$$\alpha_i (y_i(x_i^T w + w_0) - 1) = 0$$

$$\alpha_i = 0$$

$$\alpha_i \neq 0 \implies (y_i(x_i^T w + w_0) - 1) = 0$$

**Support Vector**
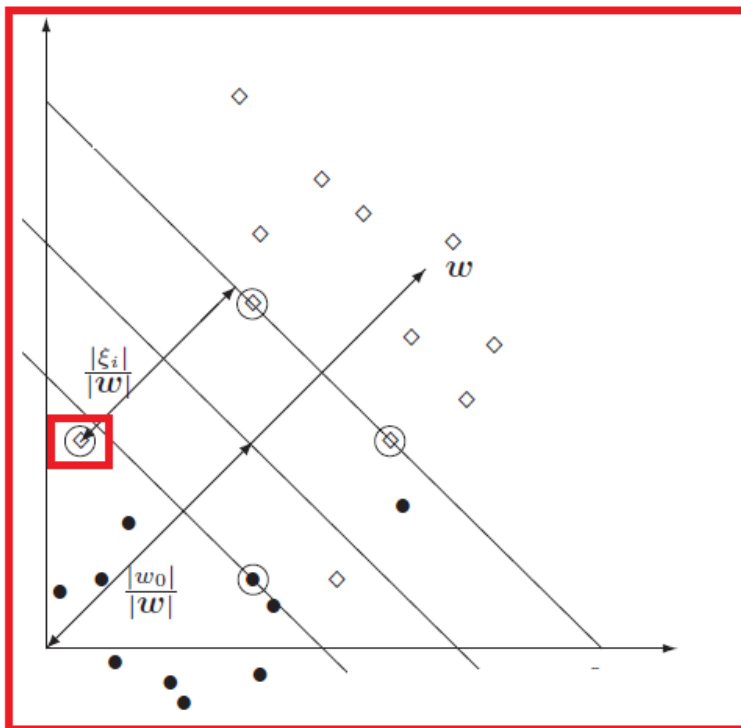
$$n_{\mathcal{SV}} w_0 + w^T \sum_{i \in \mathcal{SV}} x_i = \sum_{i \in \mathcal{SV}} y_i$$

$$w^T x + w_0 = \sum_{i \in \mathcal{SV}} \alpha_i y_i x_i^T x - \frac{1}{n_{\mathcal{SV}}} \sum_{i \in \mathcal{SV}} \sum_{j \in \mathcal{SV}} \alpha_i y_i x_i^T x_j + \frac{1}{n_{\mathcal{SV}}} \sum_{i \in \mathcal{SV}} y_i$$

$$\sum_{i \in \mathcal{SV}} \alpha_i y_i x_i^T x - \frac{1}{n_{\mathcal{SV}}} \sum_{i \in \mathcal{SV}} \sum_{j \in \mathcal{SV}} \alpha_i y_i x_i^T x_j + \frac{1}{n_{\mathcal{SV}}} \sum_{i \in \mathcal{SV}} y_i > 0 \quad \Longrightarrow \quad \text{assign } x \text{ to } \omega_1$$

# SVM for Linear non-separable



**Constraints:**

$$w^T x_i + w_0 \geq +1 - \xi_i \quad \text{for } y_i = +1$$
$$w^T x_i + w_0 \leq -1 + \xi_i \quad \text{for } y_i = -1$$
$$\xi_i \geq 0 \quad\quad\quad\quad i = 1, \ldots, n$$

$$y_i(w^T x_i + w_0) \geq 1 - \xi_i \quad i = 1, \ldots, n$$
$$\xi_i \geq 0 \quad\quad\quad\quad i = 1, \ldots, n$$

**Objective function**

$$\frac{1}{2} w^T w + C \sum_i \xi_i$$

**Dual Problem**

$$\text{Max} \quad L_D = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

$$\text{s.t} \quad \sum_{i=1}^{n} \alpha_i y_i = 0$$

$$0 \leq \alpha_i \leq C$$

# Complementarity conditions

$$\alpha_i (y_i (x_i^T w + w_0) - 1 + \xi_i) = 0$$

$$r_i \xi_i = (C - \alpha_i) \xi_i = 0$$

Patterns for which $\alpha_i > 0$ are termed the support vectors

$$y_i (x_i^T w + w_0) - 1 + \xi_i = 0$$

$0 < \alpha_i < C \implies \xi_i = 0$

$\xi_i \neq 0 \implies \alpha_i = C$

$x_i$ are misclassified if $\xi_i > 1$.

If $\xi_i < 1$, they are classified correctly, but

lie closer to the separating hyperplane than $1/|w|$

$\mathcal{SV}$ is the set of support vectors with associated values of $\alpha_i$ satisfying $0 < \alpha_i \leq C$

$\widetilde{\mathcal{SV}}$ is the set of $n_{\widetilde{\mathcal{SV}}}$ support vectors satisfying $0 < \alpha_i < C$

$$\sum_{i \in \mathcal{SV}} \alpha_i y_i x_i^T x + \frac{1}{n_{\widetilde{\mathcal{SV}}}} \left\{ \sum_{j \in \widetilde{\mathcal{SV}}} y_j - \sum_{i \in \mathcal{SV}, j \in \widetilde{\mathcal{SV}}} \alpha_i y_i x_i^T x_j \right\} > 0$$