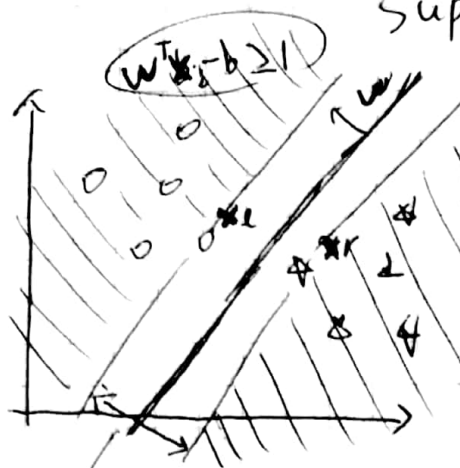


Support Vector Machine (SVM)



- Finding hyperplane which separates two classes of data and maximizes the distance to the closest point from either class.

- Support vectors: Samples on the margin (edge of the boundary line); most difficult to classify.

Consider a normal vector w to the dividing hyperplane

$$\left. \begin{array}{l} \text{LHS data satisfy } w^T x_i - b \geq 1 \\ \text{RHS } w^T x_i - b \leq -1 \end{array} \right\} \begin{array}{l} y_i (w^T x_i - b) \geq 1, \\ y_i = \begin{cases} 1, \text{ LHS} \\ -1, \text{ RHS} \end{cases} \end{array} \quad \text{constraint}$$

Difference between two hyperplanes

$$w^T x - b = 1, w^T x - b = -1 \text{ is,}$$

$$(x_l - x_r) \cdot \frac{w}{\|w\|} = \frac{(x_l w - x_r w)}{\|w\|} = \frac{(1+b) - (-1+b)}{\|w\|} = \frac{2}{\|w\|} \quad \text{maximize}$$

→ minimize $\|w\|$, where

$$y_i (w^T x_i - b) \geq 1 \text{ for } i=1, \dots, N \text{ (\# of data samples)}$$

Lagrange method

$$\text{minimize } \sum_i c_i - \frac{1}{2} \sum_i \sum_j y_i c_i (x_i^T x_j) y_j c_j, \text{ where}$$

$$\sum_i c_i y_i = 0, 0 \leq c_i \leq \frac{1}{2n\lambda}$$

- Kernel trick: Defines inner product in the transformed space; can learn non-linear classification by applying linear classification on transformed data points.

$$\rightarrow \sum_i c_i - \frac{1}{2} \sum_i \sum_j y_i c_i k(x_i, x_j) y_j c_j \quad \left| \quad k(x_i, x_j) = \phi(x_i) \phi(x_j) \right|$$

Kernel matrix $K = [k_{ij}]$ where $k_{ij} = k(x_i, x_j)$ is positive semidefinite

$$\text{b.c. } \forall x, x^T K x = \sum_{i,j} x_i x_j k_{ij} = \sum_{i,j} x_i x_j \phi(x_i) \phi(x_j) = \left\| \sum_i x_i \phi(x_i) \right\|_2^2 \geq 0$$

and kernel function is symmetric, $k(x_i, x_j) = k(x_j, x_i)$.

= Mercer's theorem: for any symmetric function k that is positive semidefinite, \exists transform $\phi_i, \lambda_i \geq 0$ such that $k(x, x') = \sum_i \lambda_i \phi_i(x) \phi_i(x')$.

By the Mercer's theorem, we can see that the kernel function must be symmetric, p.s.d.

probabilistic LDA (PLDA)

- In Gaussian mixture model

$P(x|y) = N(x|y, \Phi_w)$ where $x \rightarrow$ example, $y \rightarrow$ latent var. (mean of mixture)
instead of having discrete prior prob. of y ,
model $\boxed{P(y) = N(y|m, \Phi_b)}$ (Gaussian prior).

When $\Phi_w \rightarrow$ p.d and $\Phi_b \rightarrow$ p.s.d then $\exists V$ s.t.

$$V^T \Phi_b V = \Psi \text{ and } V^T \Phi_w V = I \rightarrow \Phi_w = A A^T, \Phi_b = A \Psi A^T, A = V^{-T}$$

Then $\boxed{P(y) = N(y|m, \Phi_b) = \frac{N(y|m, A \Psi A^T)}{m + A * N(v|0, \Psi)}} \quad \textcircled{1}$

$\boxed{P(x|y) = N(x|y, \Phi_w) = \frac{N(x|m + A v, A A^T)}{m + A * N(u|v, I)}} \quad \textcircled{2}$

~~where~~ $\rightarrow \boxed{\begin{matrix} y = m + A v & \textcircled{1} & v \sim N(0, \Psi) \\ x = m + A u & \textcircled{2} & u \sim N(v, I) \end{matrix}}$

Params of PLDA: m , covariance Ψ , loading matrix A
(mean) (Φ_b, Φ_w)

- while LDA is efficient on classifying known data, by making prior of the class centers continuous, unknown data can be dealt more nicely.