

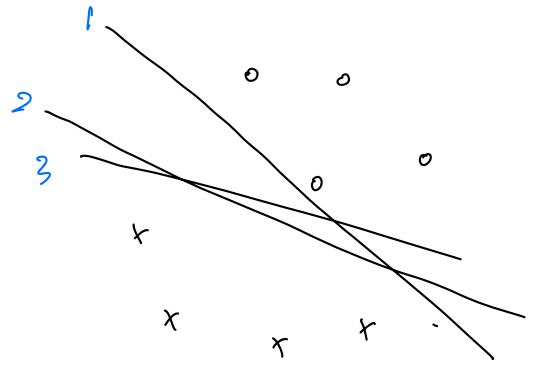
SVM

perceptron

maximize correctness

$$h_{\theta}(x) = \begin{cases} 1 & \theta^T x \geq 0 \\ 0 & \theta^T x < 0 \end{cases}$$

$$\theta_j = \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \quad \checkmark$$

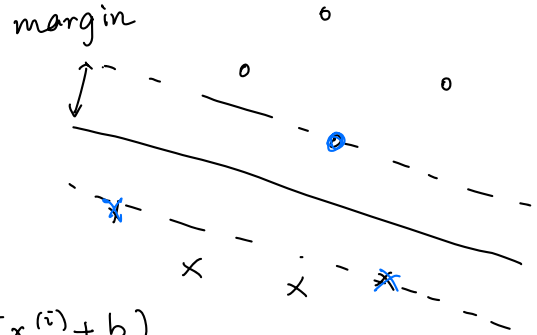


⇓

SVM

maximize the margin

$$h_{w,b}(x) = \begin{cases} 1 & w^T x + b \geq 0 \\ -1 & w^T x + b \leq 0 \end{cases}$$



functional margin

$$\hat{\gamma} = \min_i \gamma^{(i)} = \min_i y^{(i)} (w^T x^{(i)} + b)$$

geometric margin

$$\begin{aligned} \gamma &= \min_i \gamma^{(i)} = \min_i y^{(i)} \left(\frac{w^T}{\|w\|} x^{(i)} + \frac{b}{\|w\|} \right) \\ &= \frac{1}{\|w\|} \hat{\gamma} \end{aligned}$$

max γ

* prime problem

$$\min \frac{1}{2} \|w\|^2$$

$$\text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq 1 \quad i = 1, \dots, m$$

dual problem

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \langle \alpha_i y^{(i)} x^{(i)}, \alpha_j y^{(j)} x^{(j)} \rangle$$

$$\text{s.t. } \alpha_i \geq 0 \quad i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

solution

$$w^* = \sum_{i=1}^m \alpha_i^* y^{(i)} x^{(i)} \quad \checkmark$$

$$\text{(after solving } \alpha^*) \quad b^* = -\frac{1}{2} \left(\max_{i: y^{(i)} = -1} w^{*T} x^{(i)} + \min_{i: y^{(i)} = 1} w^{*T} x^{(i)} \right) \quad \checkmark$$

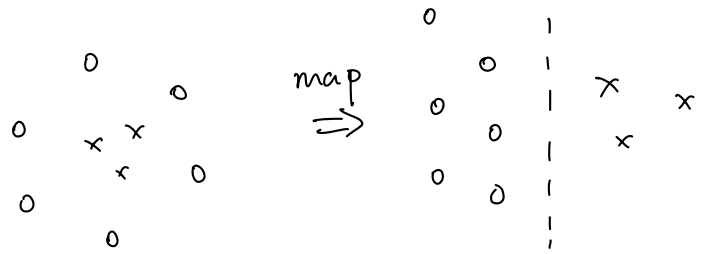
prediction $\text{sign} [w^{*T} z + b]$ for new sample z

kernel trick

$$x \in \mathbb{R}^d \Rightarrow \phi(x) \in \mathbb{R}^D \quad D \gg d$$

$$\underline{K(x, x')} \triangleq \phi^T(x) \phi(x')$$

$\in \mathbb{R}$



$$\text{dual: } \max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j K(x^{(i)}, x^{(j)})$$

s.t. $\alpha_i \geq 0 \quad i=1, \dots, m \quad \sum_{i=1}^m \alpha_i y^{(i)} = 0$

$$\text{prediction: } b^* = y^{(j)} - \sum_{i=1}^m \alpha_i^* y^{(i)} K(x^{(i)}, x^{(j)})$$

$$\text{sign} \left[\sum_{i=1}^m \alpha_i^* y^{(i)} K(x^{(i)}, z) + b^* \right]$$

Neural Network.

back propagation.

Bias & Variance

$$\hat{h}(x) \quad h(x)$$

$$\text{Bias}(\hat{h}) = E_D [\hat{h}(x) - h(x)]$$

$$\text{Var}(\hat{h}) = E_D [\hat{h}(x)^2] - E_D [\hat{h}(x)]^2$$

$$\sigma^2 = E_D [(h(x) - y)^2]$$

D: all choices of training data sampled from $P_{X,Y}$

$$\text{MSE Decomposition: } \text{MSE} = E_D [(\hat{h}(x) - y)^2]$$

$$= \text{Bias}(\hat{h}) + \text{Var}(\hat{h}) + \sigma^2$$

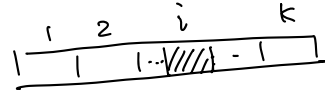
Model Selection :

models + hyperparameters.

→ cross validation :

Hold-out : train, val, test split.

K-fold :



take average.

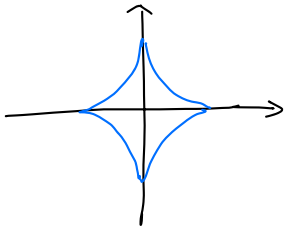
Leave-One-Out :

$k = n$

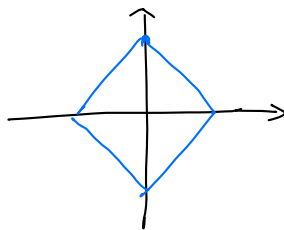
Regularization :

→ norm $L_q = \lambda \cdot \frac{1}{2} \|\theta\|_q^2 + L(X, Y; \theta)$. ✓

$q = 0.5$



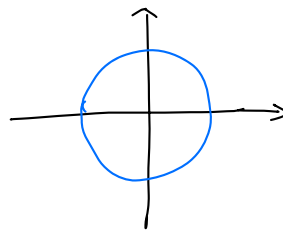
$q = 1$



(LASSO)

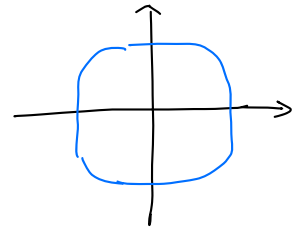
sparse solution
 $q \leq 1$

$q = 2$



(ridge)

$q = 4$



→ Data Augmentation

→ Parameter sharing

→ Drop out