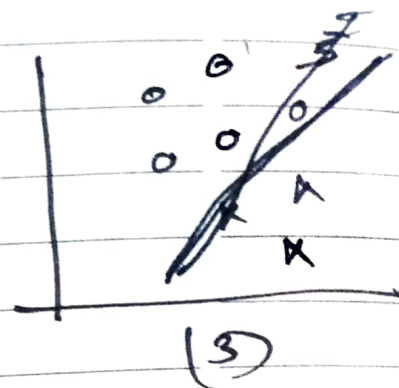
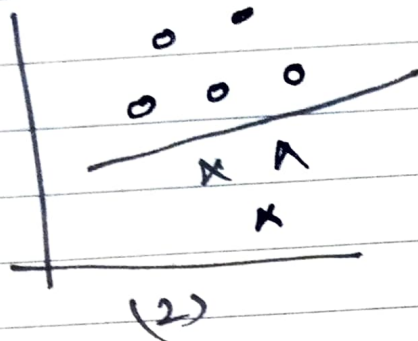
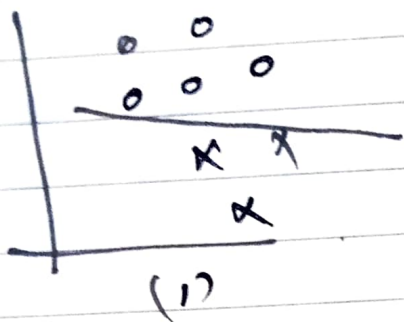


- Perceptron — learns the weights for the linearly separable data. If the data are not linearly separable, Perceptron fails to learn.
- But, if the dimensionality (data representation) is changed (from lower to higher), perceptron is able to learn the weights.
- We did the XOR problem with perceptron by adding extra dimension and moved a point that could not classify properly into the additional dimension so that we could linearly separate the classes.
- kernel function is used to do the same.
- SVM is a ML alg used to provide better performance for ~~reasonably~~ reasonably sized data set.
- don't work well for the large size dataset.
- it is very expensive

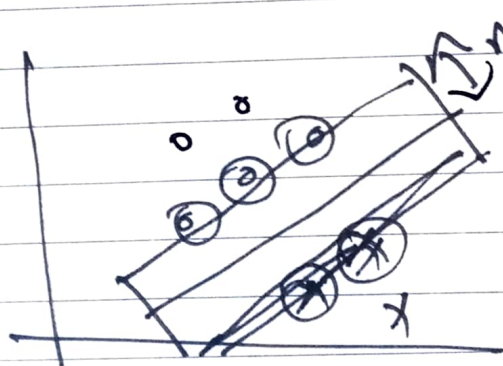
Optimal separation



- If we pick a line from (1) / (2) -

there is a chance that a data point from one class will be on the wrong side of the line.

Margin and support vector.



maximum margin classifier (M).

The data point in class lies closest to the classification line is called support vectors.

- margin should be as large as possible.
- support vectors are the most useful

Appointments

- data points. we can keep only these support vector data points after training and throw away all other data points.
- Support vector data points alone used for classification.
 - we have weight vector w and input vector x .
 - classifier line $y = wx + b$
 - any x value gives +ve value for $y = wx + b$ is above a line \in +ve class.
 - any x value gives -ve value for $y = wx + b$ is below a line \in -ve class.
- $$w \cdot x = \sum_{i=1}^n w_i x_i = w^T x.$$
- any point x where $w^T x + b \geq M$ is +ve
 $\leq -M$ is -ve
 - actual separating hyperplane is given by
 $w^T x + b = 0.$
 - pick a x^+ point on the +ve class boundary line $w^T x^+ = M$. This is support vector

- weight vector w is \perp to the decision line $y = w^T x + b$.
- we make w a unit vector $\|w\| = 1$.
- so margin is $1/\|w\|$.
- we need to find a and b from the training data, find the maximum margin M with min of $w^T x$.
- set $w=0$, max margin M and min

Constrained optimization problem

- target answers for 2 classes be ± 1 instead of 0/1.
- $t_i \times y_i$ will be +ve if both are same else -ve.
- $t_i (w^T x + b) \geq 1$
- min $\frac{1}{2} w^T w$ subject to

t_i	y_i	
+	+	+
+	+	+
-	+	-

$$t_i (w^T x + b) \geq 1 \text{ for all } i=1, 2, \dots$$

Appointments

Karush - Kuhn - Tucker condition.

will be satisfied for the optimal parameters.

$$\lambda_i^* (1 - t_i (w^{*T} x_i + b^*)) = 0.$$

$$1 - t_i (w^{*T} x_i + b^*) \leq 0$$

$$\lambda_i^* \geq 0$$

 λ_i is Lagrange multiplier.if $\lambda_i \neq 0$ then $1 - t_i (w^{*T} x_i + b^*) = 0$

This is true for support vectors.

Lagrangian function

$$L(w, b, \lambda) = \frac{1}{2} w^T w + \sum_{i=1}^n \lambda_i (1 - t_i (w^T x_i + b_i))$$

$$= \frac{1}{2} w^T w + \sum_{i=1}^n \lambda_i - \sum_{i=1}^n t_i w^T x_i \cdot \lambda_i$$

$$+ \sum_{i=1}^n \lambda_i b_i t_i \quad (3)$$

diff w.r.t w .

$$\nabla_w L = w - \sum_{i=1}^n \lambda_i t_i x_i = 0.$$

$$w = \sum_{i=1}^n \lambda_i t_i x_i \quad (1)$$

diff. w.r.t. b

$$\nabla_b L = \sum_{i=1}^n \lambda_i t_i = 0$$

$$\sum_{i=1}^n \lambda_i t_i = 0 \quad (2)$$

Sub (1) and (2) into eq (3)

$$L(w, b, \lambda) = \frac{1}{2} w^t \cdot w + \sum_{i=1}^n \lambda_i$$

$$- \sum_{i=1}^n \lambda_i t_i w^t \cdot x_i - \sum_{i=1}^n \lambda_i t_i b$$

$$= \frac{1}{2} \sum_{i=1}^n \lambda_i t_i x_i \lambda_j t_j x_j - \sum_{i=1}^n \lambda_i t_i x_i \lambda_j t_j x_j$$

$$+ \sum_{i=1}^n \lambda_i$$

$$= -\frac{1}{2} \sum_{i=1}^n \lambda_i \lambda_j \cdot t_i t_j \cdot x_i x_j + \sum_{i=1}^n \lambda_i$$

This eq is dual problem. Constraint is $\lambda_i \geq 0$ for all i and $\sum_{i=1}^n \lambda_i t_i = 0$

— for the support vector $t_i (w^t x_i + b) = 1$

$$L = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \lambda_i \lambda_j t_i t_j x_i x_j$$

22 Sunday max L , always +ve

— if $\lambda_i = 0$, then x_i is not S.V.

— if λ_i is very high for S.V.

— high influence for the position of hyperplane

Appointments

of λ_i is antireadily large, x_i is spurious.

$$b^* = \frac{1}{n_s} \sum_{\text{support vector } j} \left(t_j - \sum_{i=1}^n \lambda_i t_i x_i^T x_j \right).$$

- to make a prediction, for a new point z

$$\underline{w^{*T} \cdot z} + \underline{b^*} = \left(\sum_{i=1}^n \lambda_i \cdot t_i x_i \right)^T z + b^*$$

slack variables for non-linearly separable problems.

- Introduce slack variables $\eta_i \geq 0$ to resolve non-linearly separable prob., so constraint becomes. $t_i (w^T x_i + b) \geq 1 - \eta_i$, for inputs that are correct, we set $\eta_i = 0$

- minimize $w^T \cdot w + c \alpha$. (distance between misclassified pts from the boundary line)

- if c is small, large margin, err-less.

- if c is large, small " , err-high.

- turn this into soft margin classifier

$$L(w, \epsilon) = w^T \cdot w + c \sum_{i=1}^n \epsilon_i$$

$0 \leq \lambda_i \leq c$, supports vectors for those with $\lambda_i > 0$

KKT constraint becomes .

$$\lambda_i^* (1 - t_i (w^* \cdot x_i + b^*) - \eta_i) = 0 \quad (1)$$

$$(1 - \lambda_i^*) \eta_i = 0 \quad (2)$$

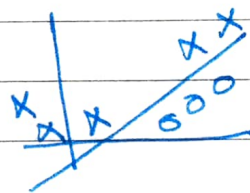
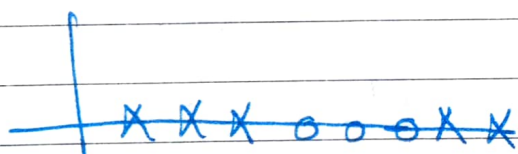
$$\sum_{i=1}^n \lambda_i^* t_i = 0 \quad (3)$$

In (2) if $\lambda_i \neq 0$, then $\eta_i = 0$. which means that they are support vectors.

Kernels

$$w^T x + b = \left(\sum_{i=1}^n \lambda_i t_i \phi(x_i) \right)^T \phi(x) + b$$

↑ the dimensions.



$$x_1^2 + x_1$$

x_1, x_2, \dots, x_d , input $x_1^2, x_2^2, \dots, x_d^2$

$x_1 x_2, x_1 x_3, \dots, x_{d-1} x_d$.

$\phi(x_i)$ - to replace x_i in higher dim.

If $d=3$

$$\phi(x) = (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_3, x_1^2, x_2^2, x_3^2, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \sqrt{2}x_2x_3)$$

$$\phi(x)^T \cdot \phi(y) = 1 + 2 \sum_{i=1}^d x_i y_i + \sum_{i=1}^d x_i^2 y_i^2 + 2 \sum_{i,j=1, i < j}^d x_i x_j y_i y_j$$

↓ reduced to

$$(1 + x^T y)^2$$

kernels.

1) polynomials up some degree s , in element x_k
 $k(x, y) = (1 + x^T y)^s$

for $s=1$, linear kernel.

2) sigmoid function of x_k 's. with parameters $k \in \mathbb{R}$ and σ
 $k(x, y) = \tanh(k x^T y - \sigma)$

3) rbf kernel $k(x, y) = \exp\left(-\frac{(x-y)^2}{2\sigma^2}\right)$
 kernels in 2-d.

26

116-249

Thursday

April 2018

Wk - 17

Appointments

Extension of SVM

$C_1 \quad C_2 \quad C_3$

$C_1 / C_2, C_3$
SVM1

$C_2 / C_1, C_3$
SVM2

$C_3 / C_1, C_2$
SVM3

y largest (farther away) from the boundary in two side.