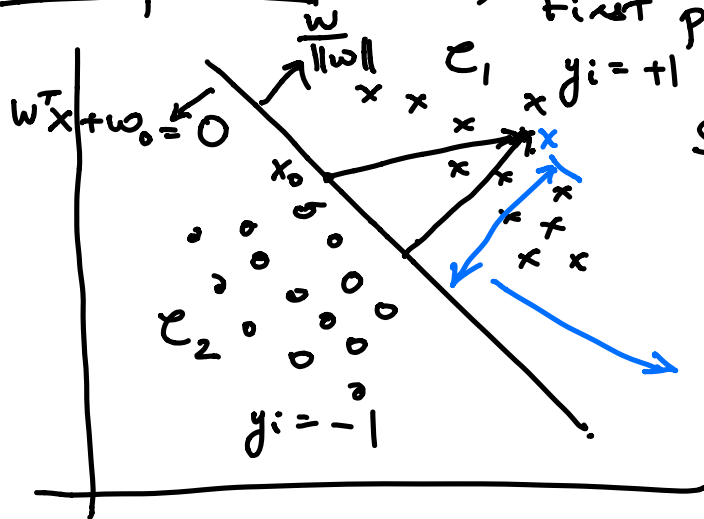


Mar 6, 2020

# Perceptron

First proposed in 1962



Signed distance of  $x$  to hyperplane :

$$\begin{aligned} & \frac{w^T}{\|w\|} (x - x_0) \\ &= \frac{w^T x - w^T x_0}{\|w\|} \\ &= \frac{w^T x + w_0}{\|w\|} \end{aligned}$$

(because  $w^T x_0 + w_0 = 0$ )

For class  $C_1$ ,  $w^T x_i + w_0 > 0$  for all points  $x_i$  in  $C_1$  that are correctly classified

For class  $C_2$ ,  $w^T x_i + w_0 < 0$  for all points  $x_i$  in  $C_2$  that are correctly classified

For correctly classified points,  
 $y_i (w^T x_i + w_0) > 0,$

For misclassified points :

$$y_i (w^T x_i + w_0) < 0.$$

Perceptron Criterion

$$\min_{D(w, w_0)} D(w, w_0) = \left[ \sum_{i \in M} y_i (w^T x_i + w_0) \right]$$

$\rightarrow M$  indexes the misclassified points

$$\nabla_w D(w, w_0) = - \sum_{i \in M} y_i x_i = - \left( \sum_{i \in C_1 \cap M} x_i - \sum_{i \in C_2 \cap M} x_i \right)$$

$$\nabla_{w_0} D(w, w_0) = - \sum_{i \in M} y_i = - (N_1 - N_2)$$

where  $N_i$  is the number of misclassified points in  $C_i$

Perceptron Update Rule

$$\begin{bmatrix} w \\ w_0 \end{bmatrix} \leftarrow \begin{bmatrix} w \\ w_0 \end{bmatrix} + \eta \begin{bmatrix} y_i x_i \\ y_i \end{bmatrix}$$

Stochastic Gradient Descent : look at  $x_i$ , and if  $x_i$  is misclassified go in direction of negative gradient (but only contribution to gradient by  $x_i$ ).

## Perceptron pseudocode

Start with some  $w, w_0$

Repeat

```
for  $i = 1$  to  $n$  do
  if  $y_i(w^T x_i + w_0) < 0$  then
     $w \leftarrow w + \eta y_i x_i$ 
     $w_0 \leftarrow w_0 + \eta y_i$ 
  end if
end for
```

until there are no mistakes (misclassifications)  
within the for loop

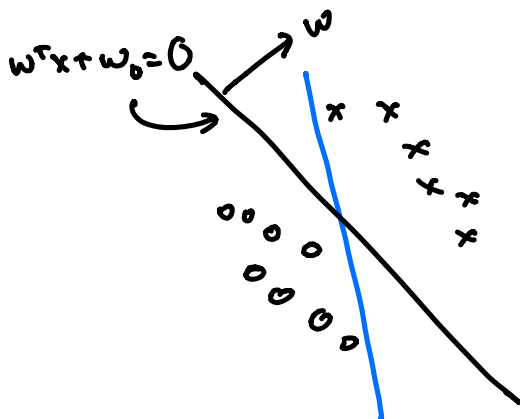
Perceptron algorithm is guaranteed to find a separating hyperplane if the data is linearly separable.

## Drawbacks of Perceptron:

1. If the data is linearly separable, the hyperplane output depends on the order in which points (data) is presented to the algorithm
2. Number of iterations might be large
3. If classes are not linearly separable, then the algorithm will not converge - cycles can

develop that are not easy to detect.

## Support Vector Machines



Signed distance of  $x$   
to hyperplane:

$$\frac{w^T x - w^T x_0}{\|w\|} = \frac{w^T x + w_0}{\|w\|}$$

$y_i \frac{(w^T x_i + w_0)}{\|w\|}$  - Distance of each training point  $x_i$  to the hyperplane

Suppose we put requirement that each of these distances is greater than  $C$

$$y_i \frac{(w^T x_i + w_0)}{\|w\|} \geq C$$

$$\Rightarrow y_i (w^T x_i + w_0) \geq C \cdot \|w\|$$

In linear support vector machines, goal is to maximize  $C$ .

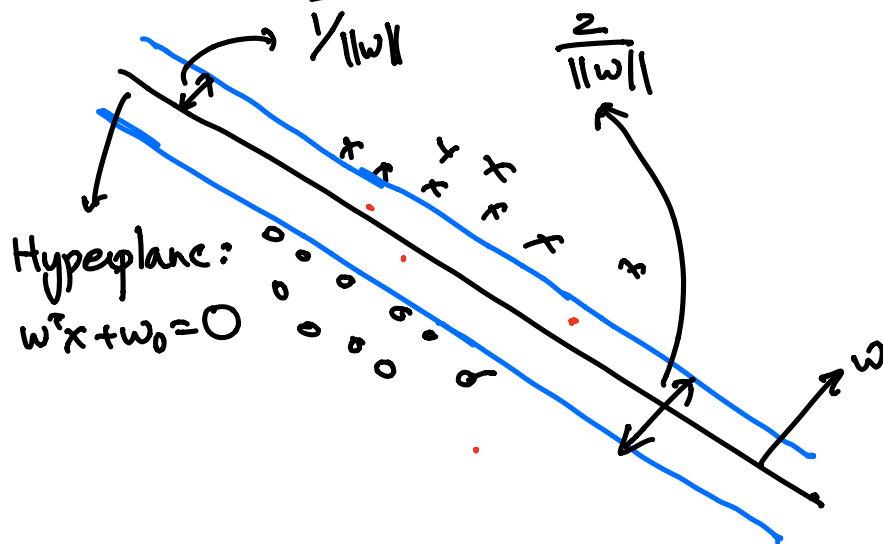
Maximize  $C$   
 $w, w_0$   
such that  $y_i (w^T x_i + w_0) \geq C \|w\|$

Fix  $C\|w\| = 1$  (since I can arbitrarily scale  $w$  &  $w_0$ )

$$\Rightarrow C = \frac{1}{\|w\|}$$

Primal  
SVM  
Problem:

$$\begin{aligned} &\text{Minimize } \|w\|^2 \\ &\text{such that } y_i(w^T x_i + w_0) \geq 1, \quad i=1, 2, \dots, N \end{aligned}$$



General constrained optimization problem

$$\min_x f_0(x)$$

$$\text{st } f_i(x) \leq 0, \quad i=1, 2, \dots, m$$

$$h_i(x) = 0, \quad i=1, 2, \dots, p$$

Here  $x \in \mathbb{R}^n$  (i.e.,  $f_i: \mathbb{R}^n \rightarrow \mathbb{R}, h_i: \mathbb{R}^n \rightarrow \mathbb{R}$ )

① } Primal

Lagrangian  $L: \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)$$

where  $\lambda_i$  is Lagrange multiplier for  $i$ -th inequality constraint

&  $\nu_i$  " " " " " equality "

$\lambda$  &  $\nu$  are also dual variables

Lagrange dual function:

$$g(\lambda, \nu) = \inf_x L(x, \lambda, \nu) = \inf_x \left( f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x) \right)$$

$\lambda$  &  $\nu$  are dual feasible if  $\lambda \geq 0$  &  
 $g(\lambda, \nu) > -\infty$

Fact 1 .  $g(\lambda, \nu) \leq p^*$  for any dual feasible  $\lambda, \nu$   
(where  $p^*$  is optimal value of (P))  
 $p^* = f_0(x^*)$

Fact 2 . If  $\exists$  dual feasible  $\lambda^*, \nu^*$  ( $\lambda^* \geq 0$ )  
& primal feasible  $x^*$  st  $g(\lambda^*, \nu^*) = p^* = f_0(x^*)$

then strong duality is said to hold

Dual Problem :

$$\boxed{\begin{array}{l} \max_{\lambda, \nu} g(\lambda, \nu) \\ \text{st } \lambda \geq 0 \end{array}}$$

Suppose strong duality holds for  $\lambda^*, v^*, x^*$

$$\begin{aligned}
 \text{Then } f_0(x^*) &= g(\lambda^*, v^*) = \inf_x L(x, \lambda^*, v^*) \\
 &= \inf_x \left( f_0(x) + \sum_{i=1}^m \lambda_i^* f_i(x) + \sum_{i=1}^p v_i^* h_i(x) \right) \\
 &\leq f_0(x^*) + \underbrace{\sum_{i=1}^m \lambda_i^* f_i(x^*)}_{\leq 0} + \underbrace{\sum_{i=1}^p v_i^* h_i(x^*)}_{=0} \\
 &\leq f_0(x^*)
 \end{aligned}$$

So, the above two inequalities must hold with equality

①  $x^*$  must be minimizer of  $L(x, \lambda^*, v^*)$

②  $\sum_{i=1}^m \lambda_i^* f_i(x^*) = 0 \Rightarrow \lambda_i^* f_i(x^*) = 0, i=1, 2, \dots, m$   
 if  $\lambda_i^* > 0$ , then  $f_i(x^*) = 0$

Karush Kuhn Tucker if  $f_i(x^*) < 0$ , then  $\lambda_i^* = 0$

↑  
 KKT Conditions : Provide certificate of optimality

when problem ① is convex

KKT conditions for  $x^*, \lambda^*, v^*$  :

$$\text{Primal Feasible} \begin{cases} f_i(x^*) \leq 0, & i=1,2,\dots,m \\ h_i(x^*) = 0, & i=1,2,\dots,p \end{cases}$$

$$\text{Dual Feasible: } \lambda_i^* \geq 0, \quad i=1,2,\dots,m$$

$$\text{Complementary Slackness} \rightarrow \lambda_i^* f_i(x^*) = 0, \quad i=1,2,\dots,m$$

$$x^* = \underset{x}{\text{argmin}} L(x, \lambda^*, \nu^*)$$

$$\rightarrow \nabla f_0(x^*) + \sum \lambda_i^* \nabla f_i(x^*) + \sum \nu_i^* \nabla h_i(x^*) = 0$$

SVM Problem:

$$\text{Primal: } \min_{w, w_0} \frac{1}{2} \|w\|^2$$

$$\text{st } y_i(w^T x_i + w_0) \geq 1, \quad \text{for } i=1,2,\dots,N$$

$$\downarrow$$

$$1 - y_i(w^T x_i + w_0) \leq 0 \quad \text{for } i=1,2,\dots,N$$

Lagrangian

$$L(w, w_0, \alpha) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^N \alpha_i (1 - y_i(w^T x_i + w_0))$$

$$\nabla_w L(w, w_0, \alpha) = 0 \Rightarrow w + \sum_{i=1}^N -\alpha_i y_i x_i = 0 \Rightarrow w = \sum_{i=1}^N \alpha_i y_i x_i$$

$$\nabla_{w_0} L(w, w_0, \alpha) = 0 \Rightarrow -\sum_{i=1}^N \alpha_i y_i = 0 \Rightarrow \sum_{i=1}^N \alpha_i y_i = 0$$



Dual function:  $g(\alpha) = \inf_{w, w_0} L(w, w_0, \alpha)$

$$g(\alpha) = \frac{1}{2} \left\| \sum_{i=1}^N \alpha_i y_i x_i \right\|^2 + \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \alpha_i y_i x_i^T \left( \frac{\sum_{j=1}^N \alpha_j y_j x_j}{\sum_{i=1}^N \alpha_i y_i w_0} \right)$$

$$g(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j$$

SVM Dual:  $\max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j$

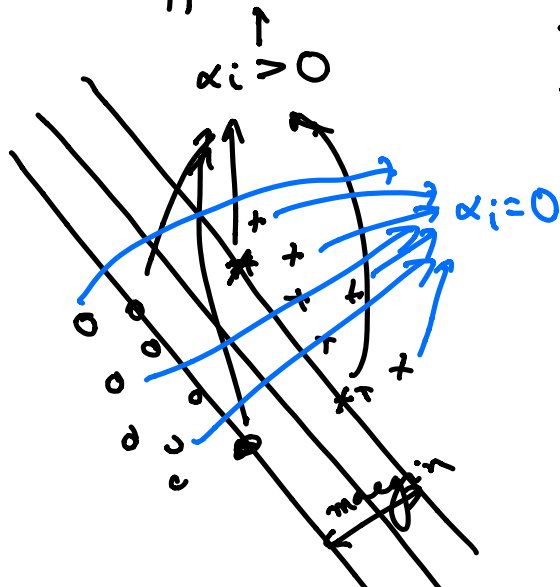
such that  $\alpha_i \geq 0, i=1, 2, \dots, N$   
 &  $\sum_{i=1}^N \alpha_i y_i = 0$

Complementary Slackness

At optimality,  $\alpha_i, w, w_0$  must satisfy:

$$\alpha_i (1 - y_i (w^T x_i + w_0)) = 0 \text{ for all } i$$

"Support" vectors



If  $\alpha_i > 0$ , then  $y_i (w^T x_i + w_0) = 1$

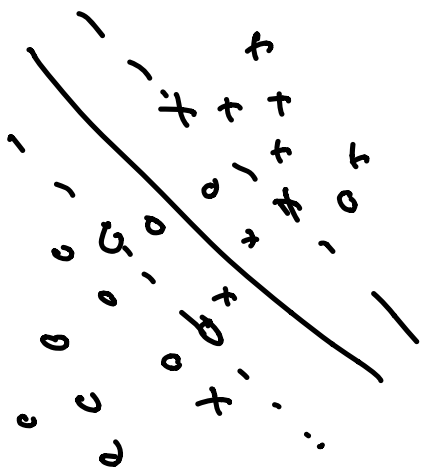
If  $y_i (w^T x_i + w_0) > 1$ , then  $\alpha_i = 0$

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

↓  
 linear combination of support vectors

Recall original separable SVM formulation:

$$\begin{aligned} \max_{w, w_0} \quad & C \\ \text{st} \quad & y_i (w^T x_i + w_0) \geq C \|w\| \end{aligned}$$



Relax above inequality to have a slack, i.e.,

$$y_i (x_i^T w + w_0) \geq C(1 - \xi_i) \|w\|, \quad \xi_i \geq 0$$

$$C \|w\| = 1$$

Non-linearly separable SVM:

$$\text{Primal: } \min_{w, w_0, \xi} \frac{1}{2} \|w\|^2$$

$$y_i (w^T x_i + w_0) \geq 1 - \xi_i, \quad i = 1, 2, \dots, N$$

$$\xi_i \geq 0, \quad i = 1, 2, \dots, N$$

$$\sum_{i=1}^N \xi_i \leq \text{constant}$$

SVM  
Primal

$$\begin{aligned} \min_{w, w_0, \xi} \quad & \frac{1}{2} \|w\|^2 + \gamma \left( \sum_{i=1}^N \xi_i \right) \\ & 1 - \xi_i - y_i (w^T x_i + w_0) \leq 0, \quad i = 1, 2, \dots, N \\ & \xi_i \geq 0, \quad i = 1, 2, \dots, N \end{aligned}$$