

Support Vector Machine

①

- To find an optimal separating hyperplane which maximizes the margin of the training data.
- A supervised learning model, i.e. have labeled data.
- Used for classification & regression.
- Learns a linear model (hyperplane)

[Hyperplane]: A generalization of a plane $\vec{w}^T \vec{x} = 0$
1 dim (point), 2-dim (line), 3-dim (plane)

Objective: To find ^{an optimum} separating hyperplane

[Margin]: Twice the distance between a given hyperplane and the closest point.

[Optimum separating hyperplane] will be the one with biggest margin, so that it correctly classifies the training data & generalizes the unseen data better.

BASIC-I

Orthogonal projection of \vec{x} onto \vec{y}

$$\vec{OA} = \vec{x}, \quad \vec{OB} = \vec{y} \quad \text{and} \quad \vec{OC} = \vec{z}$$

$$\cos \theta = \frac{\vec{x} \cdot \vec{y}}{\|\vec{x}\| \|\vec{y}\|} \quad (\text{dot product})$$

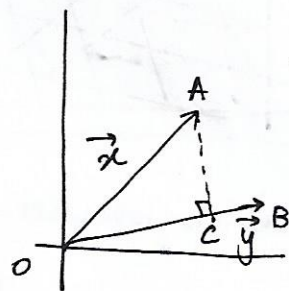
$$\text{Let } \vec{u} = \frac{\vec{y}}{\|\vec{y}\|} \quad (\text{unit vector in direction of vector } \vec{y})$$

$$\text{Also } \cos \theta = \frac{\|\vec{z}\|}{\|\vec{x}\|} \Rightarrow \|\vec{z}\| = \|\vec{x}\| \cos \theta = \frac{\vec{x} \cdot \vec{y}}{\|\vec{y}\|} = \vec{x} \cdot \vec{u}$$

Since \vec{z} is in the same direction as \vec{y} .

& \vec{u} is the unit vector i.e. $\vec{u} = \frac{\vec{z}}{\|\vec{z}\|} \Rightarrow \vec{z} = \|\vec{z}\| \vec{u}$

The vector $\boxed{\vec{z} = (\vec{u} \cdot \vec{x}) \vec{u}}$ is the projection of \vec{x} onto \vec{y}



Orthogonal vector: Two vectors are orthogonal if their dot product is zero. 18

Hyperplane $\vec{\beta}^T \vec{x} = 0$ [Affine subspace of dimension $p-1$ in p -dimensional space]

Consider a 2-dim, hyperplane is a line, say $y = ax + b$

$$\Rightarrow y - ax - b = 0 \Rightarrow \vec{\beta} = \begin{pmatrix} -b \\ -a \\ 1 \end{pmatrix} \quad \vec{x} = \begin{pmatrix} x \\ y \end{pmatrix}$$

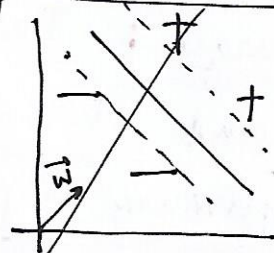
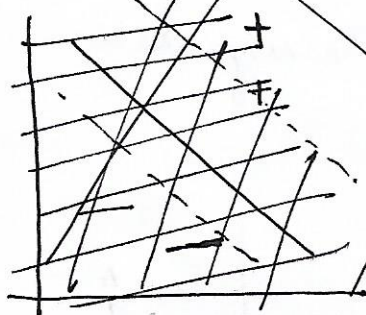
$$\text{Thus } \vec{\beta}^T \vec{x} = y - ax - b$$

~~vector~~ $\vec{\beta}$ is ~~not~~ normal to hyperplane $\vec{\beta}^T \vec{x} = 0$

By definition a hyperplane is defined, we suppose that we have a vector that is orthogonal to the hyperplane.

Normal: A line or vector \perp to a given object.

Decision boundary



MAXIM MARGIN CLASSIFIER

Decision boundary

Consider a two dimensional space, hyperplane is a line.
say $\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$ ① for parameters β_0, β_1 & β_2

Eqⁿ ① defines the hyperplane means that $\vec{x} = (x_1, x_2)^T$ lies on hyperplane eqⁿ ① holds.

In p -dimensional setting, hyperplane is

$$\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p = 0 \quad \text{--- ②}$$

$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$ and $x = (x_1, x_2, \dots, x_p)^T$, then hyperplane is given by $\beta^T x = 0$

If X doesn't satisfy (2) $\Rightarrow \beta^T X > 0$ or $\beta^T X < 0$ (2)
Thus, hyperplane can be ~~be seen as~~ ~~thought of as a~~
thought of as a subspace dividing a p -dim space in 2 halves

Classification using hyperplane

Suppose we have a $n \times p$ data matrix X that consists of
 n training observations in a p -dimensional space

$$x_1 = (x_{11} \ x_{12} \ \dots \ x_{1p})^T$$

$$x_2 = (x_{21} \ x_{22} \ \dots \ x_{2p})^T$$

and these observations fall into 2 classes,

$$y_1, y_2, \dots, y_n \in \{-1, 1\} \text{ where } -1 \text{ and } 1 \text{ represent}$$

two classes respectively.

Let $x^* = (x_1^* \ \dots \ x_p^*)^T$ be a p -dimensional test vector

Goal is to classify x^* using feature measurement

Suppose it is possible to construct a hyperplane that
separates the training observations perfectly to their class
labels. Then a separating hyperplane has a property that

$$\begin{aligned} \beta^T x_i &> 0 \quad \text{if } y_i = 1 \\ \beta^T x_i &< 0 \quad \text{if } y_i = -1 \end{aligned}$$

Equivalently we can say that $y_i(\beta^T x_i) > 0 \quad \forall i = 1, \dots, n$

$$\text{Let } f(x^*) = \beta^T x^*.$$

If $f(x^*) > 0$ classify x^* in class 1

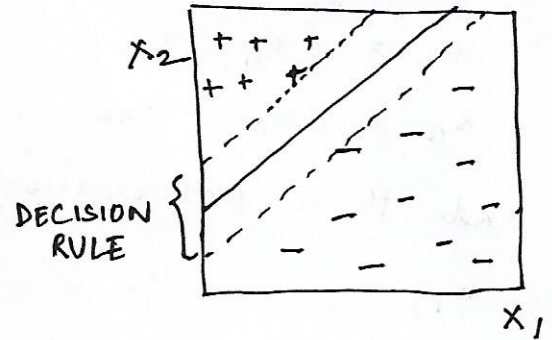
otherwise $f(x^*) < 0$ it ~~does~~ belongs to class -1

More is the magnitude of $f(x^*)$, more certain we are
certain about assignment

optimal separating hyperplane OR Maximal margin classifier 2B

- The maximal margin hyperplane is the separating hyperplane for which margin is largest i.e. hyperplane that has the farthest minimum distance to the training observations. This when used to classify test observations is called maximal margin classifier.

Maximal margin classifier depends only a small subset of observations (support vectors) and not entire training set.



So support vectors are vectors in n -dimensional space which "support" maximal margin classifier, that is, if one moves these, classifier moves

Hyperplane : strong line
Margin = distance between two dotted lines

Note can result in overfitting
BASIC-2 compute the distance from a point to hyperplane

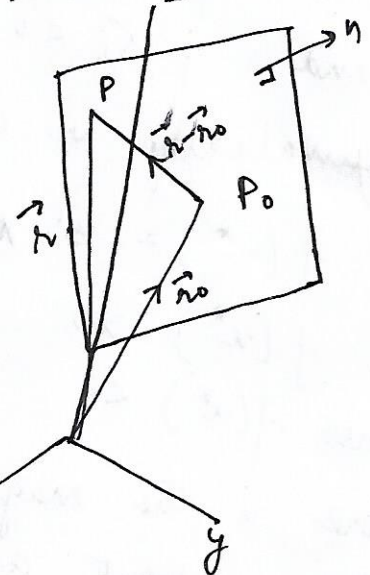
- ① Consider hyperplane $\beta^T x = 0$ and point x_0 .
- ② Let x lie on the plane satisfying $\beta^T x = 0$.
- ③ Consider construct $x_0 - x$ which points x from x_0 and project it on β ($\because \beta$ is normal)

④ distance

$$d = \|\text{proj}_{\beta} (x_0 - x)\|$$

$$= \left\| \frac{(x_0 - x) \cdot \beta}{\beta \cdot \beta} \beta \right\| = \frac{|x_0 \cdot \beta - x \cdot \beta| \|\beta\|}{\|\beta\|^2} = \frac{|x_0 \cdot \beta - x \cdot \beta|}{\|\beta\|}$$

= width of separating plane lines



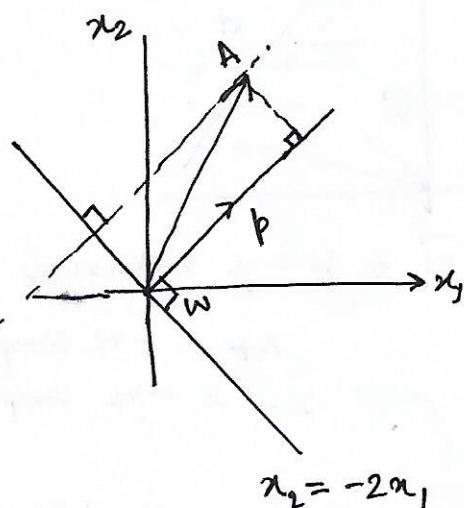
Ex: Consider a hyperplane $x_2 = -2x_1$

$$\beta^T x = 0 \Leftrightarrow \beta^T = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \text{ and } x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

Example

To compute distance b/w point $A(3,4)$ and the hyperplane

This is the distance b/w A and its projection onto hyperplane



$$\beta^T = (2, 1) \quad \vec{a} = (3, 4) \quad (\text{Vector from origin to } A)$$

$$\|\vec{\beta}\| = \sqrt{2^2 + 1} = \sqrt{5}$$

Let \vec{u} be in direction of $\vec{\beta}$

$$\vec{u} = \left(\frac{2}{\sqrt{5}}, \frac{1}{\sqrt{5}} \right)$$

\vec{p} is the orthogonal projection of \vec{a} onto $\vec{\beta}$

$$\vec{p} = (\vec{u} \cdot \vec{a}) \vec{u} = \left(\frac{6}{\sqrt{5}} + \frac{4}{\sqrt{5}} \right) \vec{u} = \frac{10}{\sqrt{5}} \vec{u}$$

$$= \left(\frac{20}{5}, \frac{10}{5} \right) = (4, 2)$$

$$\|\vec{p}\| = \sqrt{4^2 + 2^2} = \sqrt{20} = 2\sqrt{5}$$

$$\text{Margin} = 2\|\vec{p}\| = 4\sqrt{5}$$

Let \vec{x}_+ and \vec{x}_- be two points on the dotted lines

$$\beta^T \vec{x}_+ = 1$$

$$\text{and } \beta^T \vec{x}_- = -1$$

$$\Rightarrow \beta_1 x_{1+} + \beta_2 x_{2+} = 1 - \beta_0$$

$$\text{and } \beta_1 x_{1-} + \beta_2 x_{2-} = -1 - \beta_0$$

$$\Rightarrow \beta_1 (x_{1+} - x_{1-}) + \beta_2 (x_{2+} - x_{2-}) = 2$$

$$\Rightarrow \frac{\vec{w}}{\|\vec{w}\|} (x_+ - x_-)$$

$$\beta^T \vec{x}_+ = \beta_0 + \beta_1 x_{1+} + \beta_2 x_{2+}$$

$$\text{Let } \vec{w} = (w_1, \dots, w_p)$$

$$x_+ = x_{1+}, \dots, x_{p+}$$

$$x_- = x_{1-}, \dots, x_{p-}$$

Consider decision rule

$$\vec{w} \cdot \vec{u} + b \geq 0$$

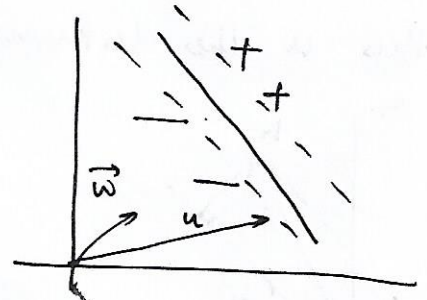
(A)

$$\vec{w} \cdot \vec{u} \geq c$$

- \vec{w} is perpendicular to hyperplane

- \vec{u} is a point for which we have to classify

WLOG $\vec{w} \cdot \vec{u} + b \geq 0$ Then + else -
else.



\vec{w} & b are unknown

$x_+ = +ve$ samples
 $x_- = -ve$ samples

So we put constraints to find \vec{w} & b

$$(1) \quad \vec{w} \cdot \vec{x}_+ + b \geq 1 \quad \text{and} \quad \vec{w} \cdot \vec{x}_- + b \leq -1$$

define y_i st $y_i = \begin{cases} +1 & +ve \text{ samples} \\ -1 & -ve \text{ samples} \end{cases}$

① & ② reduce to $y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1$

$$(B) \quad y_i (\vec{w} \cdot \vec{x}_i + b) - 1 = 0 \quad \text{for } \vec{x}_i \text{ for all samples on the dotted lines}$$

Step B : Say \vec{x}_+ and \vec{x}_- are points on dotted lines
width $\vec{x}_+ = \vec{x}_-$

$$(\vec{x}_+ - \vec{x}_-) \frac{\vec{w}}{\|\vec{w}\|} = \text{width of separating lines}$$

Using (B), we have $\vec{w} \cdot \vec{x}_+ + b - 1 = 0$
& $(\vec{w} \cdot \vec{x}_- + b) - 1 = 0$

$$\Rightarrow \vec{w} \cdot \vec{x}_+ = 1 - b \quad \text{and} \quad \vec{w} \cdot \vec{x}_- = -1 - b$$

$$\Rightarrow \vec{w} \cdot \vec{x}_+ - \vec{w} \cdot \vec{x}_- = \vec{w} \cdot (\vec{x}_+ - \vec{x}_-) \quad 1 - b + 1 + b = 2$$

$$\Rightarrow \vec{w} \cdot \vec{x}_+ - \vec{w} \cdot \vec{x}_- = 2$$

multiplying both sides by $\frac{1}{\|\vec{w}\|}$

$$\Rightarrow \frac{1}{\|\vec{w}\|} \vec{w} \cdot (\vec{x}_+ - \vec{x}_-) = \frac{2}{\|\vec{w}\|}$$

Thus, width of separating hyperplane is $\frac{2}{\|\vec{w}\|}$

To get the widest separating hyperplanes, maximize

$$\frac{2}{\|\vec{w}\|} \propto \text{minimize } \|\vec{w}\|$$

This is equivalent to ~~maximize~~ ^{minimize} $\frac{1}{2} \|\vec{w}\|^2$. (4)

which is a constraint quadratic programming problem with constraints 4 is solved using Lagrange's multiplier.

$$L = \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^n \alpha_i [y_i (\vec{w} \cdot \vec{x}_i + b) - 1]$$

$$\frac{\partial L}{\partial \vec{w}} = \vec{w} - \sum_{i=1}^n \alpha_i y_i \vec{x}_i = 0$$

$$\begin{aligned} \|\vec{w}\|^2 &= \vec{w} \cdot \vec{w} \\ \frac{\partial \|\vec{w}\|^2}{\partial \vec{w}} &= \vec{w} + \vec{w} \\ &= 2\vec{w} \end{aligned}$$

$$\Rightarrow \boxed{\vec{w} = \sum_{i=1}^n \alpha_i y_i \vec{x}_i}$$

ie \vec{w} is a linear sum of samples.

$$\frac{\partial L}{\partial b} = -\sum \alpha_i y_i = 0 \Rightarrow \sum \alpha_i y_i = 0$$

$$L = \frac{1}{2} \|\vec{w}\|^2 - \vec{w} \cdot \vec{w} + \sum \alpha_i \Rightarrow \boxed{L = -\frac{1}{2} \|\vec{w}\|^2 + \sum_{i=1}^n \alpha_i}$$

$$\frac{\partial^2 L}{\partial \vec{w}^2} = 1 \quad \frac{\partial^2 L}{\partial \vec{w} \partial b} = 0 \quad \frac{\partial^2 L}{\partial b^2} = 0 \quad \nabla^2 L = \begin{vmatrix} 1 & 0 \\ 0 & 0 \end{vmatrix}$$

$$L = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \vec{x}_i \cdot \vec{x}_j + \sum_{i=1}^n \alpha_i \quad (4)$$

Thus the decision rule is only dependent of \vec{x}_i & \vec{x}_j

In case points are not linearly separable, transform to higher dimension $\phi(\vec{x})$ and we need to maximize

$$\phi(\vec{x}_i) \phi(\vec{x}_j)$$

Kernels

let $k(\vec{x}_i, \vec{x}_j)$ be the kernel funcⁿ such that-

$$k(\vec{x}_i, \vec{x}_j) = \phi(\vec{x}_i) \phi(\vec{x}_j)$$

Thus using kernels one can produce non linear boundary by constructing a linear boundary in a large transformed version of the feature space