# * Support Vector Machine (SVM)

- Supervised algorithm used for classification and regression.

  SVC: Support Vector Classifier

  SVR: Support Vector Regressor

- We plot each of the data points in n-dimensional space (n = no. of features) with value of each feature being the coordinate of each data point.

- Then we try to find hyperplane which seperates data points for classification or try to find hyperplane which has max. no. of data points for regression.

- Advantages:
  1. Effective in high dimensional space
  2. Effective if n > no. of samples
  3. Versatile as different kernels can be used for descision function.

- Disadvantages:
  1. If number of features too much high, i.e no. of dimensions high then overfitting occurs, to overcome this we need to choose our kernel wisely.
  2. Don't provide probability estimate directly and need to use 5 cross validation technique.

[A.] Building Formula by intution.

- Equation of simple line is, (Ref. linear regression)

  $$y = mx + c$$

  or

  $$y = m_1 x_1 + m_2 x_2 + \dots + m_n x_n + c \quad \text{or}$$

  $$h\theta_o(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_2 + \dots + \theta_n x_n$$

---

- Algebrically it is same as, (mul. by constants)

  $$ax + by + c = 0$$

  with,

  $$y = \boxed{-\frac{a}{b}} x - \boxed{\frac{c}{b}}$$

  (arrow to intercept)
  (arrow to coefficient)

  also for multiple features,

  $$a x_1 + b x_2 + d x_3 + \dots + z_n x_n + c = 0$$

  to make it bit generalized replace coeff by $w_n$

  $$\boxed{w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n + c = 0}$$

  Converting this equation to matrix for. ease, where,

  $$W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}_{n \times 1} \qquad X = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}_{1 \times n}$$
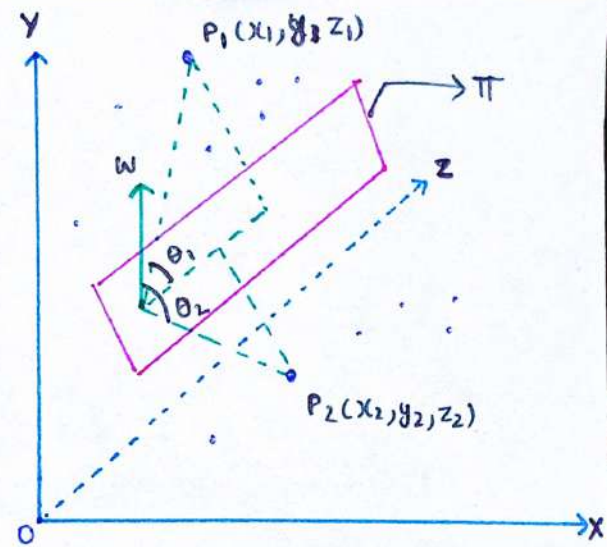
  To put up in equation change the order of w by taking Transpose,
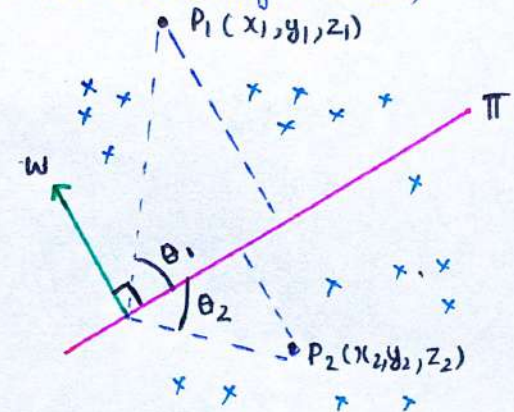
  $$\boxed{W^T X + C = 0}$$

  If intercept $c = 0$, $W^T X = 0$ is eq. of line passing via origin.

[B] Distance of point from plane

- For simplification, Let we have 3 features, i.e 3 dimensions.

- Our features $x_1, x_2, x_3$ encoded as $x, y, z$ for geometric intution

- Their are two data points $P_1(x_1, y_1, z_1)$ and $P_2(x_2, y_2, z_2)$

---



For ease lets bring it to 2D,



- W: unit vector on $\pi$ $(w \perp \pi)$
- $\pi$: a plane
- Points : X

- Distance of $P_1$ from $\pi$ is

  $$\frac{W^T P_1}{\|W\|}$$

- Distance of $P_2$ from $\pi$ is

  $$\frac{W^T P_2}{\|W\|}$$

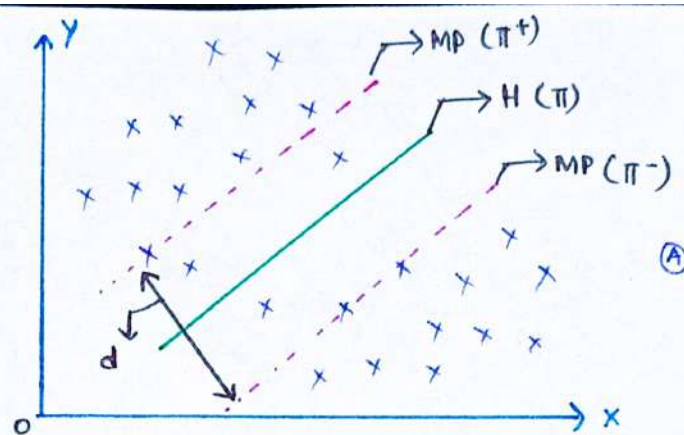- In general, Distance of a point from a plane (d) is

$$d = \frac{W^T P}{\|W\|}$$

or

$$d = \|W\| \cdot \|P\| \cdot \cos\theta$$

- **Observation:**

1. Points above the plane make $\theta_1$ angle with $\hat{w}$ and from formula and fact that $\theta_1$ is always less than $90°$ as $\hat{w} \perp \pi$
   
   ∴ Value of $\cos\theta$ for $\theta < 90$ is always +ve.
   
   ⇒ Points above plane always have +ve 'd'.

2. Points below plane make $\theta_2$ angle with $\hat{w}$ and from formula and fact $\theta_2$ always more than $90°$
   
   ∴ Value of $\cos\theta$ for $\theta > 90$ is always -ve.
   
   ⇒ Points below plane always have -ve 'd'

### C Support Vector Classifier

- **Support Vectors (sv):** Data points that are closer to hyperplane and influence the position and orientation of hyperplane.

- **Marginal Plane (MP):** Planes closer to the sv are marginal plane and help in choosing hyperplane.

- **Margin (d):** Distance between two marginal plane.

- **Hyperplane (H):** Best plane which clearly seperates data points with highest margin.

---


(A)

- Our aim is to find best plane which can clearly seperate data points.

- This plane has maximum margin and called hyperplane ($\pi$).

- Our goal is to maximize the margin (d), Classifier using such methodology is called **maximal margin classifier** and that maximum margin plane is called **margin maximizing plane**.

- Maximizing margin (M) given,

$$c + x_1 + x_2 + \cdots + x_n \qquad \text{(i)}$$

$$y_i(c + w_1 x_{1i} + w_2 x_{2i} + \cdots + w_n x_{ni}) \geq M, \quad i = 1 \ldots n \qquad \text{(ii)}$$

- This equation means define margin M by tunning coefficents of all variables such that margin is maximized(i) and product of predicted (observed) value with equation of respective input features should be greater than margin(ii)

---

- If hyperplane able to clearly seperate data points like fig A then it is called **Hard margin SVC**. In such case,

$$\pi^+ = w^T x_1 + c = +1$$

$$\pi^- = w^T x_2 + c = -1$$

If we add these two we get,

$$d = \frac{w^T(x_1 - x_2)}{\|w\|}$$

$$\therefore \boxed{d = \frac{2}{\|w\|}}$$

as $(\pi^+ + \pi^-) = w^T(x_1 - x_2) = 2$

- So, $d = 2/\|w\|$ is margin in case of hard margin classifier.

- We need to maximize 'd' by changing coefficents of features in matrix X which are present in matrix $W^T$.

- This margin is basically observed values distances under constraint,

$$y_i \begin{cases} 1, & w^T x + c \geq 1 \\ -1, & w^T x + c \leq 1 \end{cases}$$

for all points which is basically error,

- Combining constraints we get,

$$\boxed{y_i(w^T x + c) \geq 1}$$

- As $d = 2/\|w\|$ is our margin with respect to a support vector, which is also error if we inverse it,

$$\therefore \|w\|/2$$

So, we have to maximize our margin $2/\|w\|$ or minimize our loss or error $\|w\|/2$.

- In most of cases data is not linearly seperable by hyperplane and this condition is resolved by **Soft margin svc**.

- For classifying under this case we introduce slack variable $(\xi)(x_i)$ in our equation yielding,

$$y_i(w^T x + c) \geq 1 - \xi_i$$

if $\xi_i = 0$, point is correctly classified else
if $\xi_i > 0$, point incorrectly classified

- Incorrect classification means $\xi$ variable is in incorrect dimension.
- $\xi_i$ is basically an error associated with $\xi$ variable

$$\therefore \text{ Average error} = \frac{1}{n} \sum_{i=1}^{n} \xi_i$$

- Our objective is to minimize cost function, which is:

$$J = \underset{(w,c)}{\text{minimize}} \frac{\|w\|}{2} + c_i \sum_{i=1}^{n} \xi_i$$

where

$c_i$: how many points we can ignore for miss classification

$\xi_i$: summation of incorrect data points from marginal plane

$$\boxed{c_i \sum_{i=1}^{n} \xi_i} \text{ is Hinge loss function.}$$

$\Rightarrow$ $c_i$ and $\xi_i$ are hyperparameters.

D. **Support Vector Regressor**

- As a Regressor it triest to fit a best plane which has maximum number of points.
- In this case our marginal planes equation are updated by introducing marginal error $(\epsilon)$

$$\therefore$$
$$\pi^+ = w^T x + c + \epsilon$$
$$\pi = w^T x + c$$
$$\pi^- = w^T x + c - \epsilon$$

- Thus our cost function updates, insted of $\xi(x_i)$ we use $\xi_i$ (itaafi).

$$\therefore \boxed{J = \underset{(w,c)}{\text{minimize}} \frac{\|w\|}{2} + c_i \sum_{i=1}^{n} \xi_i}$$

- And constraints are;

$$\boxed{|y_i - w^T x| \leq \epsilon + \xi_i}$$

where

$\epsilon$: margin of error (to decide original plane)

$\xi_i$: error above the margin