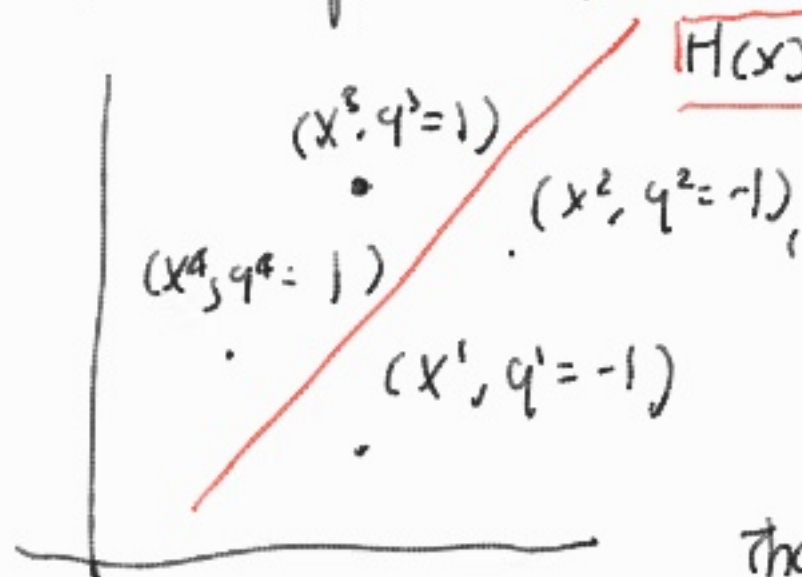


Support Vector Machine.

- given instance vector $X \in \mathbb{R}^n$, $y \in \{-1, 1\}$,
there exists a decision boundary $H = \{\theta^T X + b\}$ (where $\theta \in \mathbb{R}^n$
 b is bias)
that separates the domain where $y = -1$ and the domain where $y = 1$.



$$H(x) = \theta^T x + b$$

- $\theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$

- $x = [x_1, x_2, x_3, \dots, x_n]$

- $H(x) = \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n + b$

Then the distance $d_H(x^{(i)})$ is the distance between a certain datapoint $(x^{(i)}, y^{(i)})$ and $H(x)$

$$d_H(x^{(i)}) = \frac{|\theta^T x^{(i)} + b|}{\|\theta\|_2}$$

The objective is to find an optimal θ that results in a $H(x)$ with largest margin ($d_H(x^{(i)})$) to the closest datapoint $(x^{(i)}, y^{(i)})$

- $\underset{\theta}{\operatorname{argmax}} \left(\min_{\tilde{x}} \left(\frac{|\theta^T x^{(i)} + b|}{\|\theta\|_2} \right) \right)$

$$\bullet \arg \max_{(H)} \left(\min_i \left(\frac{|H^T x^{(i)} + b|}{\|H\|_2} \right) \right)$$

we can rewrite

$$\rightarrow |H^T x^{(i)} + b| \Leftrightarrow y^{(i)} (H^T x^{(i)} + b) \begin{cases} \geq 0 & (y \text{ and } H(x) \text{ have the same sign } \Rightarrow \text{correctly classified}) \\ < 0 & (\text{else}) \end{cases}$$

$$\arg \max_{(H)} \left(\min_i \left(\frac{y^{(i)} (H^T x^{(i)} + b)}{\|H\|_2} \right) \right) = \arg \max_{(H)} \left(\frac{1}{\|H\|_2} \min_i \left(y^{(i)} (H^T x^{(i)} + b) \right) \right)$$

this results in the following optimization problem

$$\text{objective: } \arg \max_{(H)} \frac{1}{\|H\|_2}$$

$$\text{condition: } \min_i y^{(i)} (H^T x^{(i)} + b)$$

$$y^{(i)} (H^T x^{(i)} + b) \begin{cases} \geq 1 \\ < 1 \end{cases} \leftarrow \begin{array}{l} \text{normalize the minimum margin} \\ \text{into 1 by multiplying some constant} \end{array}$$

then

$$\text{obj: } \arg \max_{(H)} \frac{1}{\|H\|_2}$$

$$\text{condition: } \min_i y^{(i)} (H^T x^{(i)} + b) = 1 \Rightarrow \forall_i$$

this indicates that our boundary H can correctly classify EVERY data points $(x^{(i)}, y^{(i)})$ (Hard Margin SVM)

Convert the problem into minimization problem
thus constructing the final form of
Hard Margin SVM

$$\text{obj} : \min_{\mathbf{H}} \frac{1}{2} \frac{1}{\|\mathbf{H}\|_2}$$

$$\text{condition} : y^{(i)} (\mathbf{H}^T x^{(i)} + b) \geq 1, \forall i$$

But in case where it's impossible to correctly classify every data points, (like most cases) we allow errors.

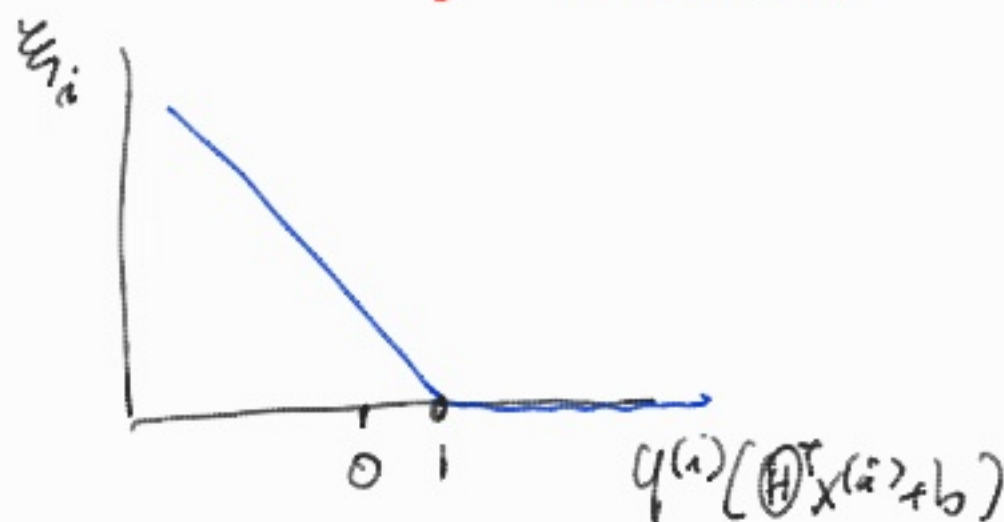
$$y^{(i)}(\theta^T x^{(i)} + b) \leq 0, \exists i \text{ for some } i$$

and introduce an error measurement / slack variable

$$\xi_i = \max(0, 1 - y^{(i)}(\theta^T x^{(i)} + b))$$

$$\xi_i = \begin{cases} 0 & (y^{(i)}(\theta^T x^{(i)} + b) \geq 1, \text{ the data point } (x^{(i)}, y^{(i)}) \text{ is correctly classified}) \\ > 1 & (\text{else}) \end{cases}$$

* the slack variable ξ_i gets larger as the incorrect margin increases.



Then we can rewrite our problem.

$$\text{obj: } \min_{\Theta, \xi} \left(\frac{1}{2} \|\Theta\|_2 + \frac{\lambda}{m} \sum_{i=1}^m \xi_i \right)$$

where λ is regularization hyperparameter
 m is # of instances.

* when λ is small enough, the slack variable loses influence. Else it's the other way around

$$\text{condition: } y^{(i)}(\Theta^T x^{(i)} + b) \geq 1 - \xi_i, \forall i$$
$$\xi_i \geq 0, \forall i$$

* since ξ_i is defined to be
 $\max(0, 1 - y^{(i)}(\Theta^T x^{(i)} + b))$

Hence, the final form of the SVM

$$\text{obj: } \min_{\Theta, \xi} \frac{1}{2} \|\Theta\|_2 + \frac{\lambda}{m} \sum_{i=1}^m \xi_i$$

$$\text{condition: } y^{(i)}(\Theta^T x^{(i)} + b) \geq 1 - \xi_i, \forall i$$
$$\xi_i \geq 0, \forall i$$

In order to solve this problem,

We need to transform our problem into a Convex optimization problem.

$$\text{obj: } \min_{\Theta, b, \xi} \frac{1}{\|\Theta\|_2} + \frac{\lambda}{m} \sum \xi_i$$

$$\text{condition: } (1 - y^{(i)})(\Theta^T x^{(i)} + b) - \xi_i \leq 0, \forall i \\ -\xi_i \leq 0, \forall i$$

Once we get a convex op. form, we then proceed to determine Lagrangian Dual Form

$$L(\Theta, b, \xi, d, \beta)$$

$$= \frac{1}{2} \|\Theta\|_2 + \frac{\lambda}{m} \sum \xi_i + \sum d_i [1 - y^{(i)}(\Theta^T x^{(i)} + b) - \xi_i] + \sum \beta_i (-\xi_i)$$

$$= \frac{1}{2} \Theta^T \Theta + \sum \xi_i \left(\frac{\lambda}{m} - d_i - \beta_i \right) + \sum d_i (1 - y^{(i)}(\Theta^T x^{(i)} + b))$$

* Since both conditions are affine,

the problem has strong duality with an appropriate Θ, b, ξ

then

$$g(d, \beta) = \max_{\Theta, b, \xi} L(\Theta, b, \xi, d, \beta)$$

calculate partial derivative of L

$$\partial_{\Theta} L = 0 \iff W - \sum d_i y^{(i)} x^{(i)} = 0 \quad \therefore \Theta = \sum d_i y^{(i)} x^{(i)}$$

$$\partial_b L = 0 \iff -\sum d_i y^{(i)} = 0 \quad \therefore \sum d_i y^{(i)} = 0$$

$$\partial_{\beta} L = 0 \iff \frac{\lambda}{m} - d_i - \beta \tilde{a} = 0 \quad \therefore d_i + \beta \tilde{a} = \frac{\lambda}{m}$$

Substitute each component with the above result
in $g(d, \beta)$

then the final form of our SVM would be

$$\text{obj: } \max_d \sum d_i - \frac{1}{2} \sum_{i,j} d_i d_j y^{(i)} y^{(j)} x^{(i)T} x^{(j)}$$

$$\text{cond: } \sum d_i y^{(i)} = 0, \quad d_i \in [0, \frac{\lambda}{m}] \rightarrow \text{regularization}$$

the corresponding solution would be

1. solve the above problem to obtain the optimal $d = (d_1, d_2, \dots, d_n, \dots, d_m)$ / dual solutions
2. obtain optimal Θ with $\Theta = \sum d_i y^{(i)} x^{(i)}$