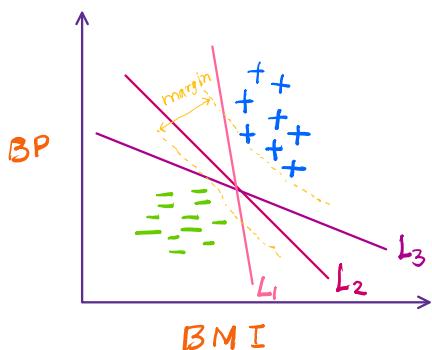
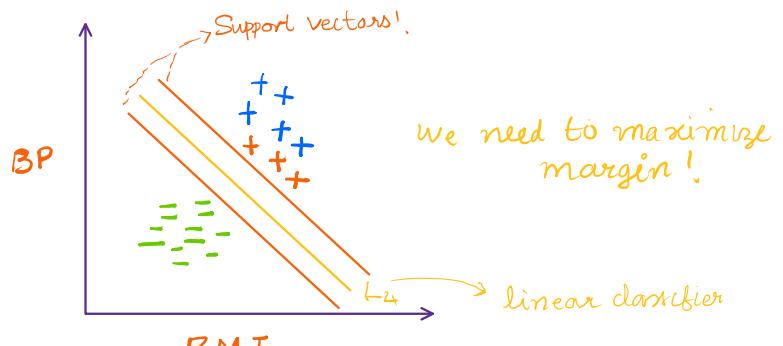


Consider a health data set shown below; where + indicates the existence of a disease and - , the absence

of the three classifiers L_1, L_2, L_3 . L_2 is better as it allows some boundary (margin)



margin based classifier.



Let our linear classifier has weight W

and bias b such that all points on the classifier have $W^T x + b = 0$.
we classify points +1,-1 according to $W^T x + b > 0$ or $W^T x + b < 0$.

Two Variants of SVM:

1. Hard margin SVMs:

allows no misclassification.
(does not exist in all cases!)

Dataset, $D = \{(x_i, y_i), \dots, (x_n, y_n)\}$

where $x_i \in \mathbb{R}^d$, d is large
 $y_i \in \{+1, -1\}$ $i, d > n$

2. Soft margin SVMs: allows misclassification, but to a limited extend.

HARD MARGIN SVM's:

Formal description:

Pick a point on the decision boundary, say x .

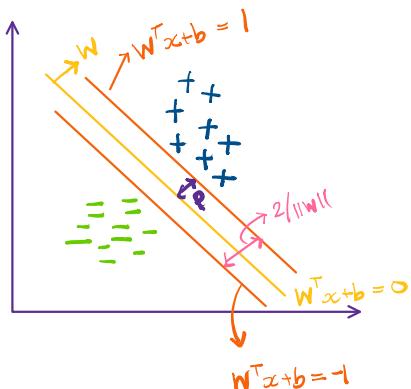
$$\therefore W^T x + b = 0$$

Move ' p ' units in direction on W to reach the margin

$$\therefore W^T (x + p \frac{W}{\|W\|}) + b = 1$$

$$\Rightarrow W^T x + p \frac{\|W\|^2}{\|W\|} + b = 1$$

$$\Rightarrow p \|W\| = 1 \text{ i.e., } p = \frac{1}{\|W\|}$$



Larger the margin larger the

Larger the margin larger the robustness/ error proofness

$$\Rightarrow P \parallel w \parallel = 1 \text{ i.e., } P = \frac{1}{\parallel w \parallel}$$

Thus, since we need a larger margin, we've

$$\max \frac{2}{\parallel w \parallel} \Rightarrow \min \left(\frac{\parallel w \parallel^2}{2} \right)$$

We also have following constraints:

$$y_i = +1 \Rightarrow w^T x_i + b \geq 1$$

$$y_i = -1 \Rightarrow w^T x_i + b \leq -1$$

$$\text{i.e., } y_i (w^T x_i + b) \geq 1, \forall i=1,\dots,n$$

Optimization problem:

$$\min \frac{1}{2} \parallel w \parallel^2$$

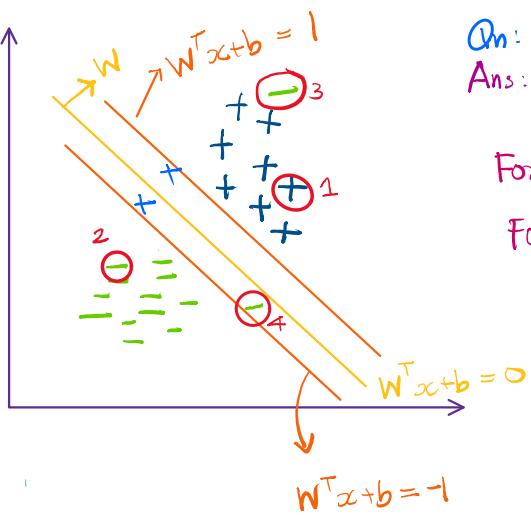
$$\text{s.t. } y_i (w^T x_i + b) \geq 1$$

$$\forall i=1,\dots,n$$

(Convex optimization problem.)

Remark: for $i \in$ support vectors, $y_i (w^T x_i + b) = 1$

SOFT MARGIN SVM's:



Data points are not linearly separable as shown in fig.

Qn: How to capture the degree of misclassification?

Ans: We use a loss function for this.

For a correct classification, $y_i \cdot (w^T x_i + b) \geq 1$

For misclassification: $y_i (w^T x_i + b) < 1$

Observe the quantity $1 - y_i (w^T x_i + b) \rightarrow 0$ for correct classification.

This is ≤ 0 for correct classification, and largely positive for great misclassification.

$$\text{Our loss fn: } L_i(w, b) = \max \{0, 1 - y_i (w^T x_i + b)\}$$

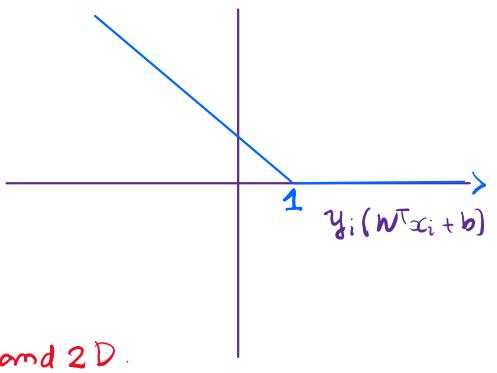
(Hinge Loss)

Our new optimization problem:

$$\min_{w, b} \frac{1}{n} \sum_{i=1}^n L_i(w, b) + \lambda \parallel w \parallel^2$$

hyperparameter.

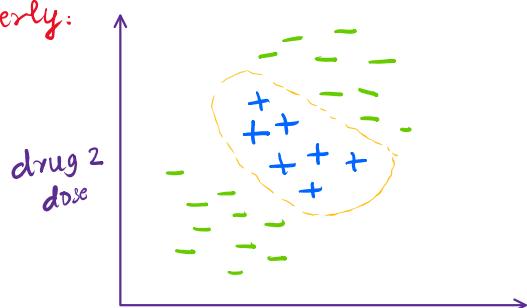
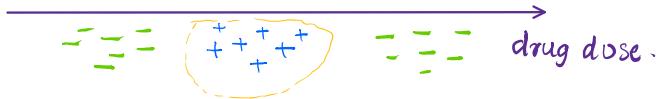
λ determines tradeoff b/w misclassification and margin maximization.



Consider the following classification problems in 1D and 2D.

Ordinary SVMs fails to classify them properly:

Ordinary SVMs fails to classify them properly:



As usual we can use basis fns $x_i \mapsto \phi(x_i)$

But since $x_i \in \mathbb{R}^d$ and $d \gg n$, computation of ϕ is expensive. Thus we adopt KERNALISATION

KERNALISATION: A method to calculate the higher dimension transformation in a computationally efficient manner.

Background of Kernelization:

$$\text{Hard Margin SVM: } \min \frac{1}{2} \|w\|^2 \text{ st: } y_i(w^T x_i + b) \geq 1, \forall i=1, \dots, n.$$

PRIMAL

In general, our objective is..

$$\min f_0(x) \quad \text{PRIMAL}$$

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

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

Domain, D where feasible x's live

$$D = \{x : f_i(x) \leq 0, \forall i; h_j(x) = 0 \forall j\}$$

Lagrangian:

$$\mathcal{L}(x, \lambda, \gamma) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^p \gamma_j h_j(x)$$

Consider the problem,

$$\max_{\lambda, \gamma} \left[\min_{x \in D} \mathcal{L}(x, \lambda, \gamma) \right]$$

This forces us to follow constraints of Primal Problem, otherwise Lagrange becomes too large.

$$g(\lambda, \gamma) = \min_{x \in D} \mathcal{L}(x, \lambda, \gamma) \leq p^* \quad \begin{matrix} \text{primal} \\ \text{optimal} \end{matrix}$$

$$\max_{\lambda \geq 0, \gamma} g(\lambda, \gamma) \quad \begin{matrix} \text{Convex opt:} \\ \text{dual opt} = p^* \end{matrix}$$

EQUIVALENT TO PRIMAL

called LAGRANGE DUAL

$$\text{Lagrange for HSVM: } \mathcal{L}(w, b, \lambda) = \frac{1}{2} \|w\|^2 - \sum \lambda_i (y_i(w^T x_i + b) - 1) \quad \lambda_i \geq 0$$

$$g(\lambda) = \min_{W, b} \mathcal{L}(W, b, \lambda)$$

$$\therefore \frac{\partial \mathcal{L}}{\partial W} = 0 \Rightarrow W = \sum_{i=1}^n \lambda_i y_i x_i \quad \text{---} \textcircled{1}$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \lambda_i y_i = 0 \quad \text{--- (2)}$$

Using ① and ②:

$$g(\lambda) = \sum \lambda_i - \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j x_i^T x_j$$

Dual Problem: $\max_{\substack{\lambda_i \geq 0 \\ \forall i}} g(\lambda)$

Note $d \gg n$. But we only need $x_i^T x_j$ inner product information for this problem. ($n \times n$)
 $nd \gg n^2$

Why we've used dual problem?

- computationally efficient, as $nd \gg n^2$
- Kernel friendly.