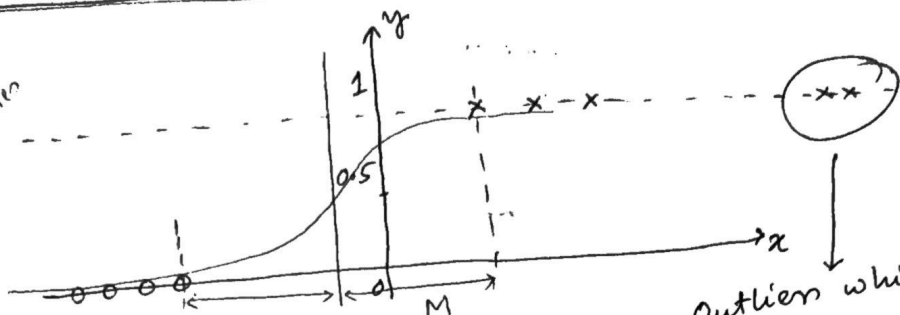


# SUPPORT VECTOR MACHINES - I

faster than  
Logit  
extendable  
to nonlinear  
boundaries



SVM: Give up probabilistic estimates  
for a better decision boundary  
only worry about points  
closest to decision boundary.

Outlier which  
significantly  
distort  
Logistic Regression.

Classification Rule:  $G(x_i) = \text{sgn}[\theta^T x_i]$

$$h(x_i, \theta) = \theta_0 + \theta^T x_i > 0 \Rightarrow \hat{y}_i = +1$$

$$h(x_i, \theta) = \theta_0 + \theta^T x_i < 0 \Rightarrow \hat{y}_i = -1$$

find  $\theta$  such that

$M$ : Margin  
 $y \in \{+1, -1\}$   
 $M$  is maximised for the closest  
points closest to Decision boundary  
hyperplane.

$\vec{r}_N = (\theta^{(1)}, \theta^{(2)})$   
(Normal to the decision boundary)

$$\vec{x}_i = \vec{x}_i^\perp + y_i x_i \frac{\vec{\theta}}{\|\vec{\theta}\|}$$

$$h(x_i, \theta) = \theta_0 + \theta^T x_i = \theta_0 + \vec{\theta} \cdot \vec{x}_i$$

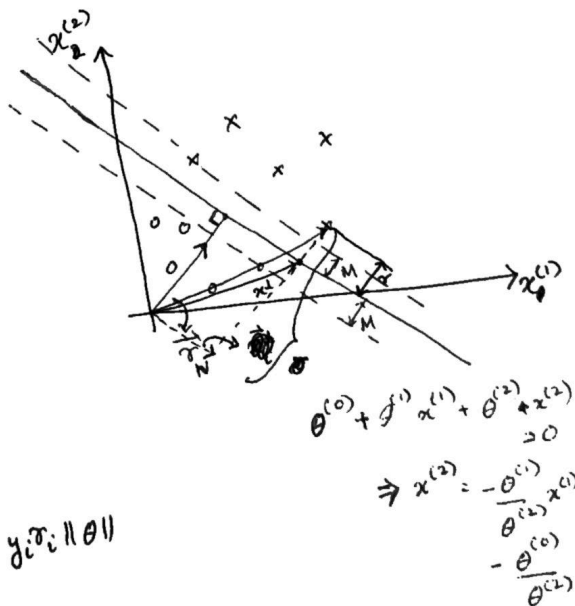
$$= \theta_0 + \vec{\theta} \cdot \left[ \vec{x}_i^\perp + y_i x_i \frac{\vec{\theta}}{\|\vec{\theta}\|} \right]$$

$$= \left( \theta_0 + \vec{\theta} \cdot \vec{x}_i^\perp \right) + y_i x_i \frac{\|\vec{\theta}\|^2}{\|\vec{\theta}\|} = y_i x_i \|\vec{\theta}\|$$

$$\Rightarrow x_i = \frac{h(x_i, \theta)}{y_i \|\vec{\theta}\|}$$

$$= \frac{\theta_0 + \theta^T x_i}{y_i \|\vec{\theta}\|}$$

Objective: Maximise  
 $r_i$  for the closest [Support  
point [Vector]  
while ensuring  
 $y_i h(x_i, \theta) > 0$



## OBJECTIVE

$$\max_{\theta, \theta_0} \min_{x_i} \frac{y_i (\theta_0 + \theta^T x_i)}{\|\theta\|}$$

Note: Scaling  $\theta, \theta_0$  does not change this function.

So, scale them so that  $y_i h_i = 1$  for point closest to decision boundary.

$$\begin{aligned} & \max_{\theta, \theta_0} \frac{1}{\|\theta\|} \quad \text{s.t.} \quad y_i (\theta_0 + \theta^T x_i) \geq 1 \quad \forall i \\ & = \min_{\theta, \theta_0} \frac{1}{2} \|\theta\|^2 \quad \text{s.t.} \quad y_i (\theta_0 + \theta^T x_i) \geq 1 \quad \forall i \end{aligned}$$

QUADRATIC program with N linear constraints [convex]  
 $\Rightarrow$  regularisation  $\Rightarrow$  good generalisation

Lagrange function

Convex  $\hookrightarrow L_D = \frac{1}{2} \|\theta\|^2 - \sum_{i=1}^N \alpha_i [y_i (\theta_0 + \theta^T x_i) - 1]$

$\alpha_i \geq 0 \quad \forall i$   
 $\alpha_i \neq 0$  only for support vectors.

$$\max_{\alpha_i} \min_{\theta, \theta_0} L_D$$

$$\frac{\partial L_D}{\partial \theta_0} = - \sum_{i=1}^N \alpha_i y_i = 0$$

$$\frac{\partial L_D}{\partial \theta_j} = \theta_j - \sum_{i=1}^N \alpha_i y_i x_{ij} = 0 \Rightarrow \theta_j = \sum_{i=1}^N \alpha_i y_i x_{ij}$$

$$\theta_j = \sum_{i=1}^N \alpha_i y_i x_{ij}$$

Support vectors.

Substituting this in  $L_D$ , we get.

$$L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N \alpha_i \alpha_{i'} y_i y_{i'} x_i^T x_{i'}$$

$\hookrightarrow$  lower bound on objective function.

Maximising this is a simpler convex quadratic programming problem.  
 LAGRANGE DUAL OBJECTIVE FUNCTION (WOLFE)

min. of original  $L_D$  is  $-\infty$  !!

## SUPPORT VECTOR MACHINES - II

What if the classes are not linearly separable?

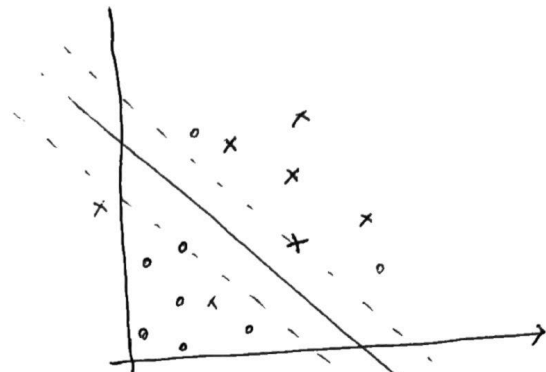
If overlap is small, use SLACK VARIABLES

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

Margin by which points

are allowed to be on ~~over~~ wrong side.

$$\xi_i \geq 0, \quad \sum \xi_i \leq \text{constant} \quad [\text{sort of regularization}]$$



$$\min_{\theta, \theta_0} \frac{1}{2} \|\theta\|^2 \quad \text{s.t.} \quad y_i(\theta_0 + \theta^T x_i) \geq 1 - \xi_i$$

$$L_D = \frac{1}{2} \|\theta\|^2 + c \sum_i \xi_i - \sum_i \alpha_i [y_i(\theta_0 + \theta^T x_i) - (1 - \xi_i)] - \sum_i \mu_i \xi_i$$

↗ trade-off between large margin + small training error

Lagrange (primal) function [minimise w.r.t.  $\theta_0, \theta, \xi_i$ ]

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_{i'} \alpha_i \alpha_{i'} y_i y_{i'} \boxed{x_i^T x_{i'}}$$

Lagrange (Wolfe) Dual Objective  $f_D$ .

Optimisation depends only on dot product of input features.

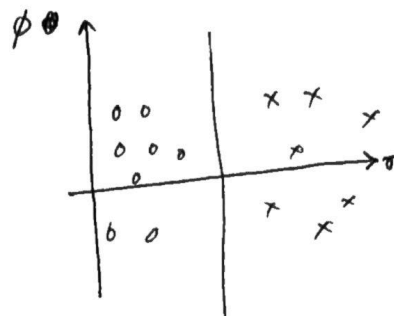
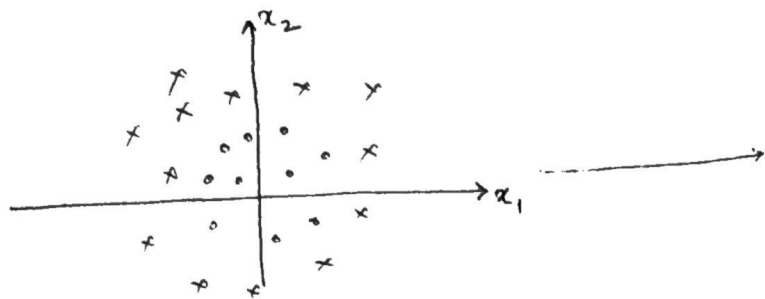
\* What if overlap is large?

Need to transform variables, so that in the new space, the points are linearly separable.

Since dual  $L_D$  only depends on dot products, we don't need the actual transformation expression but only the new dot-product.

⇒ Kernel trick.

$$L_D = \sum_i \alpha_i - \sum_i \sum_{i'} \alpha_i \alpha_{i'} y_i y_{i'} \underbrace{K(x_i, x_{i'})}_{\substack{\text{positive semi-definite} \\ \& \text{symmetric.}}}$$



Popular choices for Kernel.

$$r = \sqrt{x_1^2 + x_2^2}$$

$$\phi = \tan^{-1} x_2/x_1$$

d-th degree polynomial:

$$K(x_i, x_{i'}) = (1 + \langle x_i, x_{i'} \rangle)^d$$

usually best choice.  $\rightarrow$  Radial Basis (RBF) or Gaussian

$$K(x_i, x_{i'}) = \exp(-\gamma \|x_i - x_{i'}\|^2)$$

Most widely used based on Euclidean distance.

Neural Network or sigmoidal

$$K(x_i, x_{i'}) = \tanh(\gamma_1 \langle x_i, x_{i'} \rangle + \gamma_2)$$

In original  $L_0$ , if  $C$  is large, it penalises any  $\beta_i$  & leads to a very wiggly boundary (overfit).

$$h_i = \theta_0 + \theta^T x_i$$

$$\frac{\partial L}{\partial \theta_j} = \theta_j - \sum_i \alpha_i y_i x_{ij} = 0 \Rightarrow \theta_j = \sum_i \alpha_i y_i x_{ij}$$

$$\Rightarrow h_i = \theta_0 + \sum_i \alpha_i y_i x_i^T x_{i'}$$

$$= \theta_0 + \sum_i \alpha_i y_i K(x_i, x_{i'})$$

Curse of Dimensionality

RBF: large  $\lambda$  leads to overfitting since even small distances are penalised.

## SVM-II (contd.)

$$L_D = \frac{1}{2} \|\theta\|^2 + C \sum_i \xi_i - \sum_i \alpha_i [y_i(\theta_0 + \theta^T x_i) - (1 - \xi_i)] - \sum \mu_i \xi_i$$

minimize w.r.t.  $\theta, \theta_0, \xi_i$  → vector

$$\frac{\partial L_D}{\partial \theta_0} = 0 \Rightarrow \sum_i \alpha_i y_i = 0$$

$$\nabla_{\theta} L_D = 0 \Rightarrow \theta = \sum_i \alpha_i x_i y_i$$

$$\frac{\partial L_D}{\partial \xi_i} = 0 \Rightarrow \alpha_i = C - \mu_i$$

Substituting all these in  $L_D$ , we obtain,

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_{i'} \alpha_i \alpha_{i'} y_i y_{i'} x_i^T x_{i'}$$

→ Lagrange Dual objective function

Maximize using Quadratic Programming  
subject to  $0 \leq \alpha_i \leq C$  &  $\sum_i \alpha_i y_i = 0$

To get optimum solution, use Karush-Kuhn-Tucker conditions,

$$\alpha_i [y_i(\theta_0 + \theta^T x_i) - (1 - \xi_i)] = 0 \rightarrow \alpha_i \neq 0 \text{ only for support vectors}$$

$$\mu_i \xi_i = 0 \rightarrow \mu_i \neq 0 \text{ only when } \xi_i \neq 0$$

$$y_i(\theta_0 + \theta^T x_i) - (1 - \xi_i) \geq 0$$

## SVM - Regularisation

$$\min_{\gamma} \|\theta\|^2, \gamma \uparrow$$

$\gamma$  of kernel

$$C \sum \xi_i$$

SVM - Polynomial & Sigmoid Kernel application. [when does RBF not work]  
RBF is slower than polynomial kernel; especially for large data or high dimension  
Polynomial kernel is popular in NLP with  $d=2$  [it  $\Rightarrow$  overfitting]  
RBF is No Free Lunch Theorem  
[numerical instability]

## SVM Multi-class

one vs. one or one vs. rest.

Lagrange Multiplier  
[could be maxima/minima]

high  $\gamma$  RBF overfitting

Logreg  $\rightarrow$  find prob & change threshold.  
true  $\theta_j \Rightarrow$  feature favours  $\hat{y}=1$

$\theta_0$ : intercept  
 $\theta_1$ : coefficients

Missing values  $\rightarrow$  take mean/mode or remove from data

Kernel - Does NOT scale well to large training sets or large feature vectors

~~RBF is slower~~  
Hyper plane - Decision Boundary

Hyper parameters vs. Parameters  
(user tuned) (internal model)

Support vectors  
are not necessarily  
the points closest  
to decision boundary  
for when slack  
variables used.

Generative  
vs.  
Discriminative