

## 08 SVM

### (1) Separating Hyperplane.

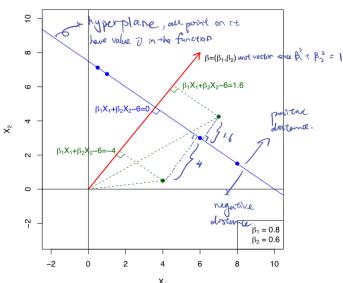
A hyperplane:  $\beta_0 + \mathbf{x}^T \beta = 0$

$\Rightarrow$  points on one side will be negative;  
points on the other side will be positive.

$\Rightarrow$  A classification rule:

$$G(\mathbf{x}) = \text{sign}[\beta_0 + \mathbf{x}^T \beta].$$

Hyperplane in 2 Dimensions



$\Rightarrow$  The best hyperplane is decided by Maximal Margin Classifier: (MMC).

$$\begin{aligned} & \text{maximize } M \\ & \text{subject to } \sum_{j=1}^p \beta_j = 1, \\ & Y_i(\beta_0 + \mathbf{x}_i^T \beta) \geq M \text{ for all } i=1, \dots, n \end{aligned}$$

$f(x_i) = 1$  (i.e. all the points are at  $M$  distance from the hyperplane)

$$\Leftrightarrow \text{maximize } \frac{1}{\|\beta\|} \text{ subject to } Y_i(\beta_0 + \mathbf{x}_i^T \beta) \geq 1, \quad i=1, \dots, n$$

$\Rightarrow$

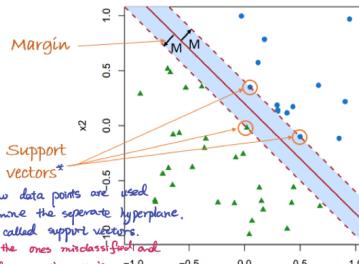
$$\text{minimize } \|\beta\|^2 \text{ subject to } Y_i(\beta_0 + \mathbf{x}_i^T \beta) \geq 1, \quad i=1, \dots, n$$

can be rewritten as:

$$\begin{aligned} & \text{maximize } C \\ & \text{subject to } \frac{Y_i}{\|\beta\|} (\beta_0 + \mathbf{x}_i^T \beta) \geq C, \quad i=1, \dots, n \end{aligned}$$

\*only a few data points are used to determine the separate hyperplane. These are called support vectors. including the ones misclassified and the ones lying on the margin.

We are looking for a hyperplane that is the farthest from the training observations, i.e. has the largest margin.



Using Lagrangian equivalent:  $\beta_j = \sum_{i=1}^n \alpha_i Y_i \mathbf{x}_i$  (4)  $\Rightarrow$  maximise  $L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j Y_i Y_j \mathbf{x}_i^T \mathbf{x}_j$  (5) subject to  $\alpha_i \geq 0; \sum_i \alpha_i Y_i = 0$

note: if  $\alpha_i > 0$ , then  $Y_i(\mathbf{x}_i^T \beta + \beta_0) = 1 \Rightarrow$  observation  $i$  is active.

if  $\alpha_i = 0$ , then  $Y_i(\mathbf{x}_i^T \beta + \beta_0) > 1 \Rightarrow$  observation  $i$  is not active.

$\Rightarrow \beta_j$  only depends on active observations, i.e.  $\beta_j = \sum_{i: \alpha_i > 0} \alpha_i Y_i \mathbf{x}_i$

### (2) Support Vector Classifier (soften separating hyperplane)

Classification rule:  $G(\mathbf{x}) = \text{sign}(\mathbf{x}^T \hat{\beta} + \hat{\beta}_0)$ .

soften MMC: (essentially a regularization method, so it's important to standardize all points first).

maximize  $M$  subject to

$$\sum_{j=1}^p \beta_j = 1; \quad (\text{except for } \beta_0)$$

$Y_i(\beta_0 + \mathbf{x}_i^T \beta) \geq M(1 - \varepsilon_i)$  → allow some points to be misclassified

$\varepsilon_i \geq 0; \sum_i \varepsilon_i \leq t$  (budget twy parameter) → the more points included, the more stable the model is.

$\downarrow$   
minimize  $\|\beta\|^2$  subject to

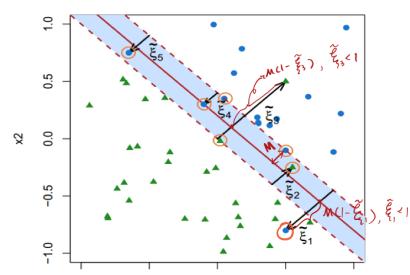
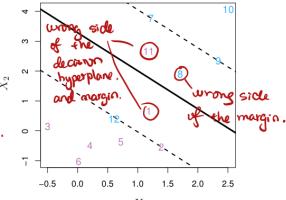
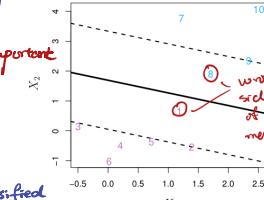
$$Y_i(\mathbf{x}_i^T \beta + \beta_0) \geq 1 - \varepsilon_i,$$

$$\varepsilon_i \geq 0, \quad \sum_i \varepsilon_i \leq t.$$

$\downarrow$

$$\text{minimize } \frac{1}{2} \|\beta\|^2 + C \sum_i \varepsilon_i \quad \text{subject to } Y_i(\mathbf{x}_i^T \beta + \beta_0) \geq 1 - \varepsilon_i, \quad \varepsilon_i \geq 0$$

can use CV to choose  $C$  & margin  $\|\beta\|$



Note that in the above diagram, the margin is  $C = 1/\|\beta\|$ , so the distances  $\varepsilon_i$  marked on the diagram equal  $\varepsilon_i/\|\beta\|$ .

### (3) Support Vector Machine. (enrich the selection space other than linear boundaries).

Classification rule:  $G(\mathbf{x}) = \text{sign}(\beta_0 + \sum_{i \in S} \alpha_i k(\mathbf{x}, \mathbf{x}_i))$  + set of support vector.

where  $k(\mathbf{x}, \mathbf{x}')$  is known as SVM kernels, and it can be

$$K(\mathbf{x}, \mathbf{x}') = (C + \langle \mathbf{x}, \mathbf{x}' \rangle)^d \quad (\text{Degree of polynomials}).$$

$$K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2/C) \quad (\text{Radial basis function})$$

$$K(\mathbf{x}, \mathbf{x}') = 1/(1 + \exp(-C \langle \mathbf{x}, \mathbf{x}' \rangle + C_0)) \quad (\text{Sigmoid function})$$

and  $\alpha_i = \alpha_i Y_i$  (recall that  $\beta_j = \sum_{i=1}^n \alpha_i Y_i \mathbf{x}_i \Rightarrow \sum_{j=1}^p \beta_j Y_j = \sum_{i=1}^n \alpha_i Y_i \mathbf{x}_i^T \mathbf{x}_j = \sum_{i=1}^n \alpha_i Y_i \mathbf{x}_i^T \mathbf{x}_j = \sum_{i=1}^n \alpha_i Y_i \langle \mathbf{x}_i, \mathbf{x}_j \rangle$  where  $\alpha_i = \alpha_i Y_i$ )

so  $\alpha_i = 0 \Rightarrow \alpha_i Y_i = 0 \Rightarrow$  observation  $i$  is non-active and vice versa.

$\hookrightarrow$  can be  $K(\mathbf{x}, \mathbf{x}')$ .

the optimisation expression is :

$$\text{minimize } \|\beta\|^2 + \gamma \sum_{i=1}^n \epsilon_i \text{ subject to } \begin{cases} Y_i (h(x_i)^T \beta + \beta_0) \geq 1 - \epsilon_i, \epsilon_i \geq 0 \\ h(x_i)^T \beta + \beta_0 \leq 1 + \epsilon_i \end{cases}$$

penalisation on complex model.  $\gamma \uparrow$  wiggly model  
 $\gamma \downarrow$  simple model.

where  $h(x_i)$  is the basic function.

\* need to do CV on  $\gamma$  as well as parameters in the kernel, e.g. C in radial function.

Again we can use Lagrangian equation, the 2 most important expressions are:

$$(6) L = \sum_i \ell_i - \frac{1}{2} \sum_i \sum_j \ell_i \ell_j K(x_i, x_j)$$

$$(6) \text{ becomes } \beta_j = \sum_i \ell_i Y_i h(x_i) \Rightarrow h(x)^T \beta + \beta_0 = \sum_i h(x_i) \beta_i + \beta_0 = \sum_i \ell_i Y_i \langle h(x), h(x_i) \rangle + \beta_0$$

where  $\langle h(x), h(x') \rangle = K(x, x')$   $\Rightarrow$  so we can calculate  $\langle x, x' \rangle$  first and put it into  $K(x, x')$  without having to calculate  $\langle h(x), h(x') \rangle$ .  $K(x, x')$  is a symmetric and positive function  
e.g.  $K(x, x') = (C + \langle x, x' \rangle)^2$

#### (4) SVM and penalised loss function.

$$\underset{\beta, \beta_0}{\text{min}} \sum_{i=1}^n [1 - Y_i (\underbrace{h(x_i)^T \beta + \beta_0}_f)]_+ + \frac{\lambda}{2} \|\beta\|^2 \text{ is equivalent to SVM}$$

$\downarrow$   
hinge loss                      ridge penalty. (shinking towards 0).

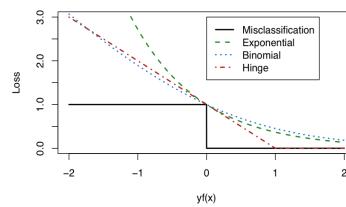
$\Rightarrow$  SVM can be thought of as fitting a basic expansion model  $f(x) = h(x)^T \beta + \beta_0$  to predict  $y_i$ , using hinge loss  $(1 - Y_i f(x))_+$  and a ridge penalty.

We may add the following to the family of loss functions for binary classification. Here  $Y_i \in \{-1, 1\}$  and  $f(X_i)$  are predictions with  $\text{sign}\{f(X_i)\}$  being the classifier.

- Misclassification loss:  $\mathbb{I}[Y_i \neq \text{sign}\{f(X_i)\}]$ .
- Binomial loss:  $-\log \left\{ \frac{1}{1+e^{-2Y_i f(X_i)}} \right\}$ .
- Exponential loss:  $\exp\{-Y_i f(X_i)\}$ .
- SVM loss:  $[1 - Y_i f(X_i)]_+$ .
- Squared loss:  $\left\{ Y_i - \frac{e^{2f(X_i)}}{1+e^{2f(X_i)}} \right\}^2$ .

For a comprehensive discussion of the different loss functions see pp.426-428 in the book "Elements of statistical learning".

We can view the various loss functions against  $y_i f(X_i)$ :



We give results for a variety of binary classification, with performance. Most models have standard error of around 0.5% on the misclassification rate

Model	Misclass. Rate
Single tree	10.4%
Bagged trees	8.4%
Logistic regression*	8.3%
Random forests	5.6%
Adaboost	5.5%
GBM	5.9%
Neural Net	6%
SVM	7.8%

#### (5) SVM in regression.

$$\sum_{i=1}^n L(y_i, h(x_i)^T \beta + \beta_0) + \frac{\lambda}{2} \|\beta\|^2.$$

In principle any loss function  $L(\cdot)$  could be used.

e.g.  $L_\varepsilon(y, f(x)) = \begin{cases} 0 & \text{if } |y - f(x)| < \varepsilon \\ |y - f(x)| - \varepsilon & \text{otherwise} \end{cases}$   $\Rightarrow$  there is a margin  $\varepsilon$  that as long as the observed values lie within, there is no loss. Outside this the loss grows linearly.