

## Lecture 12 FML April, 21

### Multiclass classification

$x$  input space,  $y$  output space  
 $K$  no of classes

Sample  $S = (x_1, y_1), \dots, (x_m, y_m)$   
 $y_i$  labels

iid  $D$  unknown

$y_i = f(x_i)$  unknown

- we want to know  $f$  from sample  $S$ .

$\mathcal{H}$  = Hypothesis class.

Set of functions from which we search for  $f$

$$\underline{R(h)} = E_D [ \mathbb{1} \{ h(x) \neq f(x) \} ]$$

$\hookrightarrow$  generalization error

$h \in \mathcal{H}$ .

- Aim is to get  $h \in \mathcal{H}$  which has min  $R(h)$ .

- We don't know  $D$   
 $\therefore$  we don't know  $R$

$$\underline{R(h)} = \frac{1}{m} \sum_{i=1}^m \mathbb{1} \{ h(x_i) \neq f(x_i) \}$$

$d_H(h(x_i), f(x_i))$   
 $\hookrightarrow$  distance

Binary classification  
via SVM & kernels

$$\mathcal{Y} = \{-1, 1\}$$

$$h(x) = \text{sgn}(w \cdot x + b)$$

$(w, b)$

$$= \text{sgn} \left( \sum_{i=1}^m \alpha_i y_i (x_i \cdot x + b) \right)$$

$$= \arg \max_{y \in \{-1, 1\}} (y \cdot (w \cdot x + b))$$

kernel.  $k$ .  $\Phi: x \rightarrow \mathcal{H}$ .

$$k(x, x') = \langle \Phi(x), \Phi(x') \rangle$$

$$\langle \Phi(x), \Phi(x') \rangle = k(x, x')$$

$$h(x) = \arg \max_{y \in \{-1, 1\}} \langle w_y, \Phi(x) \rangle$$

$$w_y = \sum_{i=1}^m y_i \alpha_i y_i \Phi(x_i)$$

For  $k \geq 2$  classes

$$h(x) = \arg \max_{y \in \mathcal{Y} = \{1, 2, \dots, k\}} h(x, y)$$

$$h(x, y) = \text{score function}$$

$$h \in \mathcal{H}$$

margin based

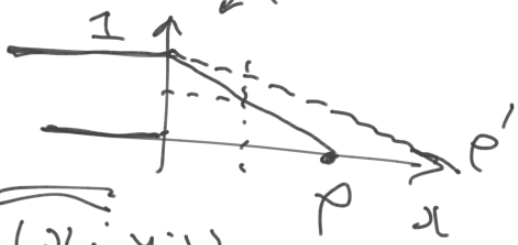
$h \in \mathcal{H}$ ,  $\rho_h(x, y) = h(x, y)$   $\leftarrow$   
 $\max_{y' \neq y} h(x, y')$   
 $(x, y)$   $\leftarrow$   $h$  misclassifies  $x$

if  $\rho_h(x, y) \leq 0$   $\leftarrow$

For a given  $\rho > 0$  define margin loss  $\Phi_\rho(x, y)$

$\hat{R}_\rho(h)$  = empirical  $m$  loss

$= \frac{1}{m} \sum_{i=1}^m \Phi_\rho(\rho_h(x_i, y_i))$



$\Pi_\rho(\mathcal{H}) = \{x \mapsto h(x, y), y \in \mathcal{Y}, h \in \mathcal{H}\}$

Theorem Let  $\mathcal{H}$  be a set of functions  $X \times Y \rightarrow \mathbb{R}$ . Fix  $\rho > 0$ . Then for any  $\delta > 0$ , with prob  $\geq 1 - \delta$ ,  $\forall h \in \mathcal{H}$

$R_\rho(h) \leq \hat{R}_\rho(h) + \frac{2k^2}{\rho} R_m(\Pi_\rho(\mathcal{H})) + \left(\frac{\log 1/\delta}{2m}\right)^{1/2}$

Theorem A family of fns  
 $\mathcal{F} \rightarrow [0, 1]$

$z_1, z_2, z_3, \dots$  i.i.d r.v.s with values in  $\mathcal{Z}$ . Then for any  $\delta > 0$ , with prob  $\geq 1 - \delta$ ,

$$E[g(z)] \leq \frac{1}{m} \sum_{i=1}^m g(z_i) + 2R_m(\mathcal{G})$$

$$E[g(z)] \leq \frac{1}{m} + 2R_S(\mathcal{G}) + 3 \left( \frac{\log \frac{1}{\delta}}{2m} \right)^{1/2}$$

Lemma (Talagrand)  $\Phi: \mathbb{R} \rightarrow \mathbb{R}$ ,  $l$ -lip  
 $\mathcal{H}$  Hypothesis set of real valued fns.  
 Then

$$\hat{R}_S(\Phi \circ \mathcal{H}) \leq l \hat{R}_S(\mathcal{H}) \leq l R_m(\mathcal{H})$$

Proof of Theorem


$$\tilde{\mathcal{H}} = \{ (x, y) \mapsto P_h(x, y), h \in \mathcal{H} \}$$

$$\hat{\mathcal{H}} = \{ \Phi_p \circ f, f \in \tilde{\mathcal{H}} \}$$

From theorem given above with prob  $\geq 1 - \delta$ ,  $\forall h \in \mathcal{H}$

$$R_p(h) = E[\Phi_p(P_h(x, y))] \leq \hat{R}_p(h) + 2R_m(\hat{\mathcal{H}}) + \left( \frac{\log \frac{1}{\delta}}{2m} \right)^{1/2}$$

From Talagrand Lemma

Since  $\Phi_p$  is  $\frac{1}{p}$ -Lip fn. 

$$R_m(\hat{\mathcal{H}}) \leq \frac{1}{p} R_m(\tilde{\mathcal{H}})$$

It is shown in Mohri's Book

$$R_m(\tilde{\mathcal{H}}) \leq K^2 R_m(\Pi_1(\mathcal{H}))$$

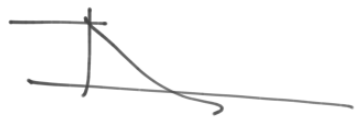
$$R_p(h) \leq \hat{R}_p(h) + \frac{2(K^2)}{p} R_m(\Pi_1(\mathcal{H})) + \dots$$

$R(h)$  = generalization error

$$= E_D \left[ \mathbb{1}_{h(x) \neq f(x)} \right]$$

$$R_p(h) = E_D \left[ \Phi_p(x, y) \right]$$

error occurs when  $\rho_n(x, y) \leq 0$



$$\therefore R(h) \leq R_p(h) \leq \dots$$

$$R_p(h) \leq \hat{R}_p(h) + \frac{2k^2}{\rho} R_m(\Pi_1(H)) + \left( \frac{\log \frac{1}{\delta}}{2m} \right)^{\frac{1}{2}}$$

We compute  $R_m(\Pi_1(H))$  for SVM kernels.

$\Sigma x$  For  $k$  class  
 $K, \Phi$

$$\arg \max_{y \in \mathcal{Y}} w_y \cdot \Phi(x)$$

$$\Phi: \mathcal{X} \rightarrow \mathcal{H}$$

$$\downarrow$$

$$\{-1, 1\}$$

$$\underline{w} = (w_1, w_2, \dots, w_k)$$

$$w_i \in \mathcal{H}$$

For  $p \geq 1$

$$\|w\|_{\mathcal{H}, p} \triangleq \left( \sum_{l=1}^k \|w_l\|_{\mathcal{H}}^p \right)^{1/p}$$

$$\mathcal{H}_{k,p} = \left\{ (x, y) \mapsto \langle w_y, \Phi(x) \rangle, \|w\|_{\mathcal{H}, p} \leq 1 \right\}$$

Prop Let  $k$  be PDS kernel

with  $k(x, x) \leq \sigma^2$  for some  $\sigma > 0$ .  
 Then for  $m \geq 1$

$$R_m(\Pi, (\mathcal{H}_k, p)) \leq \left( \frac{\sigma^2 \lambda^2}{m} \right)^{1/2}$$

Proof:  $S$  sample size

$(x_1, y_1), \dots, (x_m, y_m)$

$$\|W\|_{\mathcal{H}, p} \leq \lambda, \quad \|W\|_{\mathcal{H}, p} = \left( \sum_{l=1}^k \|w_l\|_{\mathcal{H}}^p \right)^{1/p}$$

$$\therefore (\|w_l\|_{\mathcal{H}})^{1/p} \leq \lambda$$

$$R_m(\Pi, (\mathcal{H}_k, p)) = \frac{1}{m} E_{\sigma, S} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^m \sigma_i g(z_i) \right]$$

$$= \frac{1}{m} E_{\sigma, S} \left[ \sup_{\|w\| \leq \lambda} \sum_{i=1}^m \sigma_i \langle w, \Phi(x_i) \rangle \right]$$

$$= \frac{1}{m} E_{\sigma, S} \left[ \sup_{\|w\| \leq \lambda} \langle w, \sum_{i=1}^m \sigma_i \Phi(x_i) \rangle \right]$$

$$\leq \frac{\lambda}{m} E_{\sigma, S} \left[ \sup_{\|w\|_{\mathcal{H}} \leq 1} \left\| \sum_{i=1}^m \sigma_i \Phi(x_i) \right\|_{\mathcal{H}} \right]$$

$$\leq \frac{\lambda}{m} E_{\sigma, S} \left[ \left\| \sum_{i=1}^m \sigma_i \Phi(x_i) \right\|_{\mathcal{H}} \right] \quad (\text{Cauchy-Schwarz})$$

$$\leq \frac{\lambda}{m} E_{\sigma, S} \left[ \left\| \sum_{i=1}^m \sigma_i \Phi(x_i) \right\|_{\mathcal{H}}^2 \right]^{1/2}$$

$$= \frac{\lambda}{m} E_{\sigma, S} \left[ \sum_{i=1}^m \sum_{j=1}^m \sigma_i \sigma_j \langle \Phi(x_i), \Phi(x_j) \rangle \right]$$

$$= \frac{\lambda}{m} E_S \left[ \sum_{i=1}^m \|\Phi(x_i)\|_{\mathcal{H}}^2 \right]^{1/2}$$

$$\leq \frac{\lambda}{m} E_S \left[ \sum_{i=1}^m k(x_i, x_i) \right]^{1/2}$$

$$\leq \frac{\lambda}{m} m^{1/2} \sigma = \left( \frac{\sigma^2 \lambda^2}{m} \right)^{1/2}$$

$E[\|X\|_{\mathcal{H}}^2] \leq E[\|X\|_{\mathcal{H}}^4]$   
 $E[\sigma_i \sigma_j] = 0$  for  $i \neq j$

For kernel SVM multi-class classifier

$$R_e(\hat{h}) \leq \underbrace{\hat{R}_e(h)} + \frac{2k^2}{\rho} \left( \frac{r^2 \hat{\lambda}}{m} \right)^{\frac{1}{2}} + \sqrt{\frac{\log 1/s}{2m}}$$

$$\|W\|_{A,p} \leq \lambda$$

$$k(x_i, x_j) \leq r^2$$

To get a good classifier we solve

$$\min_w \hat{R}_e(h) + \frac{1}{2} \|W\|_{A,p}$$

If the sol is  $w^*$  then the classifier is

$$x \mapsto \arg \max_{y \in \mathcal{Y}} \langle W_{y^*}^* \Phi(x) \rangle$$

- This is convex opt prob
- Primal, dual probs are in fact QP.
- complexity can be large for  $k$  large.

$$\rho \geq 1.$$

Take  $\rho = 2$

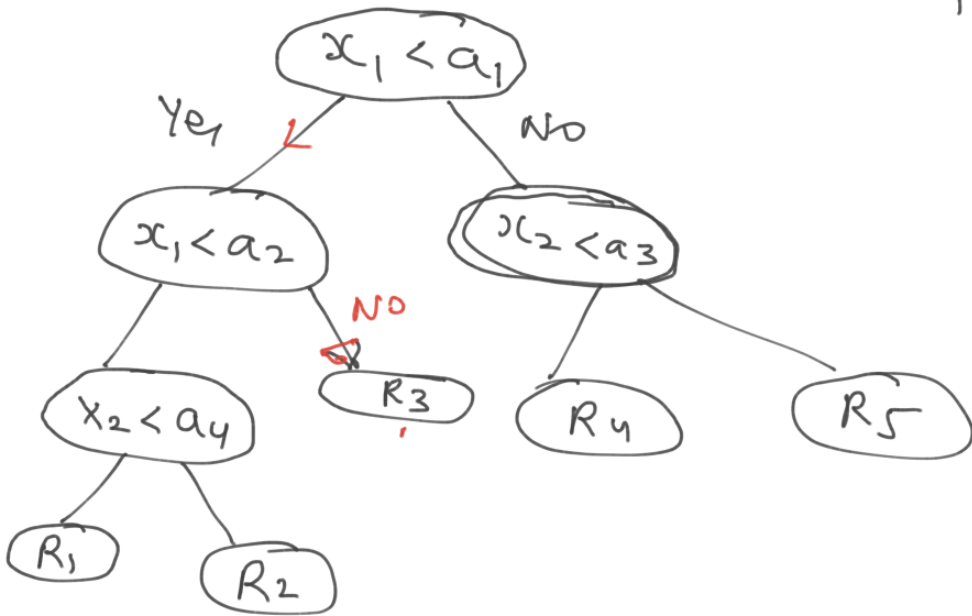
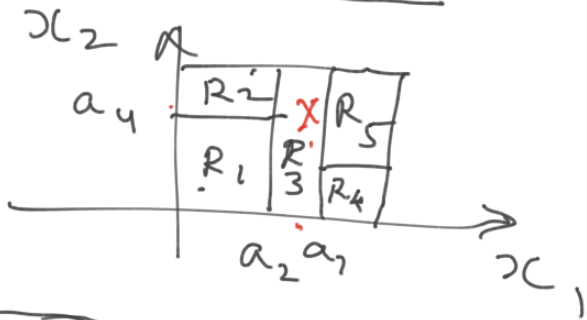
$$\hat{R}_e(h) = \frac{1}{m} \sum_{i=1}^m \xi_i$$

$$\xi_i = \max\{1 - w_{y_i} \cdot \Phi(x_i)\}$$

$$= \max_{y' \neq y_i} \{w_{y'} \cdot \Phi(x_i), 0\}$$

# Decision Tree classifiers

Ex



- Decision tree can be used also for regression, classification
- Performance of decision tree is not the best - SVM, NN can give better performance
  - Can boost performance via Random forest to get the best performing systems.
- Advantages of decision trees:
  - Fast to learn and evaluate
  - can treat numerical and categorical data water ally.



- Scales well with large data set
- Easy to explain.
- Decision tree need not be binary tree. But Binary is preferred.

- Computationally easier to learn
- For other trees the splitting of data may be too fast

Learning Decision tree:

- obtain optimal tree is NP hard.

- VC dim of a binary tree with  $n$  nodes and  $N$  features, ( $\text{dim of } X = N$ ) is  $O(n \log N)$

$\therefore$  VC dim of a general family of decision trees is  $\infty$ .

$n$  nodes,  $N$  features.  
binary tree: with  $\text{prob} \geq 1-\delta$

$$R_D(h) \leq \underbrace{R_S(h)}_{\substack{\text{general} \\ \text{error for} \\ \text{tree } h}} + \underbrace{\sqrt{\frac{(n+1) \log_2(N+3)}{2m} + \log \frac{2}{3}}}_{2m}$$

Greedy algorithm to learn a "good" decision tree