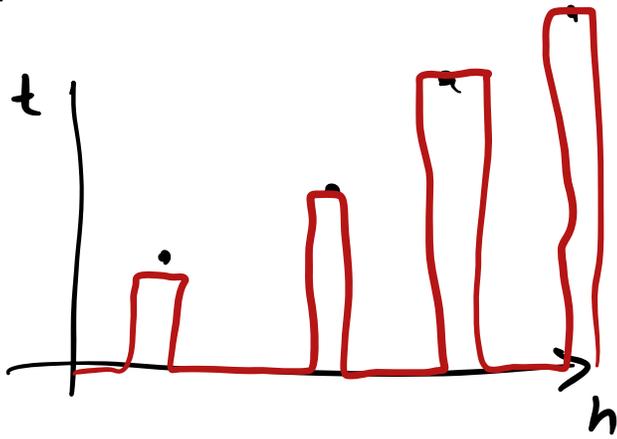


Χρῆνος πῶδης



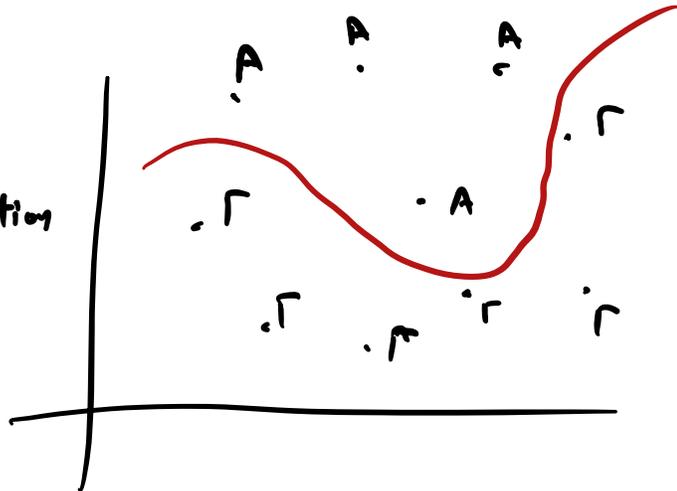
Regression

Φωτογραφίες

Αντρας / Γυναίκα

0.7	0.1
0.3	0.9

Classification



Μοντελοποίηση

$$(x, y) \sim D$$

Features

$$x \in \bar{X} \quad (\text{domain})$$

$$y \in \bar{Y}$$

Στόχος είναι να μάθουμε συνάρτηση

$h: X \rightarrow Y$  ώστε να ελαχιστοποιείται  
το  $E[\rho(h(x), y)]$   
 $(x, y) \sim D$

Για κάποια συνάρτηση κόστους  $\rho$

- Παραδείγματα:

Για regression

$$L(\hat{y}, y) = (\hat{y} - y)^2$$

Για classification

$$L(\hat{y}, y) = \mathbb{1}\{\hat{y} \neq y\} = \begin{cases} 1 & \text{αν } \hat{y} \neq y \\ 0 & \text{αλλιώς} \end{cases}$$

Training set

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$$

όπου  $(x_i, y_i) \sim D$

Κλάση υποθέσεων  $H$

→ min  
h ∈ H

$$E_{(x,y) \sim D} [\rho(h(x), y)]$$

$L_D(h)$

Empirical Risk Minimization

min  
h ∈ H

$$E_{(x,y) \sim S} [\rho(h(x), y)] \rightarrow h_S$$

$L_S(h)$       ERM

$h(x) = \begin{cases} y \\ 0 \end{cases}$

av unäppret  
 $(x,y) \in S$   
 alltid

overfitting

Ανaly περίπτωση

$$L_D(h^*) = 0$$
$$L_D(h) \in [0, 1]$$

$$l(h(x), y) = 1 [h(x) \neq y]$$

Θέλω να βρω  $h \in H$

$$\Pr_{(x, y) \sim D} [h(x) \neq y] \leq \epsilon$$

πραγματικό loss

Έστω  $H_B$  όλες οι συναρτήσεις με

$$\Pr_{(x, y) \sim D} [h(x) \neq y] > \epsilon$$

Για κάθε  $h \in H_B$

Η πιθανότητα να μην την έχω απορριψει μετά από  $m$  δείγματα

είναι το πολύ  $(1-\epsilon)^m \leq e^{-m\epsilon}$

Αν  $m = \frac{\log(1/\delta)}{\epsilon}$  απορρίπτω την κακή  $h$  με πιθανότητα  $\geq 1-\delta$

Αν το  $H$  είναι πεπερασμένο

$\Pr[\exists h \in H_B \text{ που επιβιώνει}]$

$\leq \sum_{h \in H_B} \Pr[h \text{ επιβιώνει}]$

$\leq |H_B| \delta = \underbrace{|H| \delta}_{\delta'}$

$$m = \frac{\log(|H|/\delta')}{\epsilon}$$

τότε η  $h_S \in H \setminus H_B$   
με πιθανότητα  $1 - \delta'$

Γενίκευση με  $\delta$  Uniform Convergence

Στόχος είναι να μάθουμε καλή συνάρτηση:  
 $L_D(h_S) - L_D(h^*) \leq \epsilon$

Uniform Convergence

$$\max_{h \in H} |L_S(h) - L_D(h)| \leq \epsilon$$

$\equiv$  έρουμε  $E_S [L_S(h)] = L_D(h)$   
 $\forall h$

Γιατί είναι καλό το Uniform  
Convergence;

$$\begin{aligned}L_D(h_S) - L_D(h^*) &= L_D(h_S) - L_S(h_S) \\ &\quad + L_S(h_S) - L_S(h^*) \\ &\quad + L_S(h^*) - L_D(h^*)\end{aligned}$$

$$\leq |L_D(h_S) - L_S(h_S)|$$

$$+ \underbrace{L_S(h_S) - L_S(h^*)}_{\leq 0}$$

$$+ |L_S(h^*) - L_D(h^*)|$$

$$\leq 2 \max_{h \in H} |L_D(h) - L_S(h)| \leq 2\varepsilon$$

Περαστική ετάση H

Mε m δείγματα

Για κάθε  $h \in H$

$$\Pr[|L_S(h) - L_D(h)| > \epsilon] \leq e^{-m\epsilon^2}$$

(And Hoeffding Inequality)

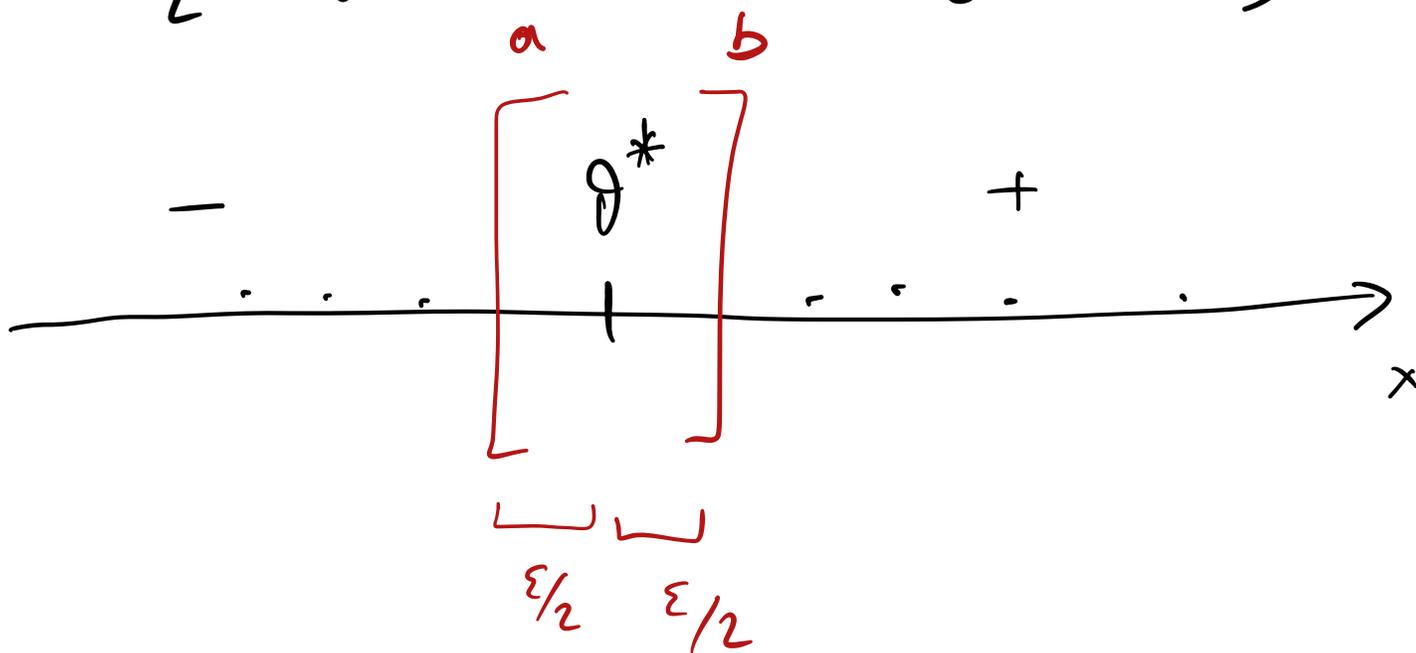
$$\Pr[\exists h \in H : |L_S(h) - L_D(h)| > \epsilon]$$

$$\leq |H| e^{-m\epsilon^2} \quad \text{and} \quad m = \frac{\log(|H|/\delta)}{\epsilon^2}$$

$$\leq \delta$$

# Threshold Ευραστηχοεις

$$H = \left\{ h_\theta : h_\theta(x) = \text{sign}(x - \theta) \right\}$$



a τ.ω.

$$\Pr [ [a, \theta^*] ] = \frac{\epsilon}{2}$$

b τ.ω.

$$\Pr [ [\theta^*, b] ] = \frac{\epsilon}{2}$$

$$\Pr [ x_1, \dots, x_m \notin [a, \theta^*] ] = \left(1 - \frac{\epsilon}{2}\right)^m \leq e^{-m\epsilon/2} \leq \delta/2$$

$$m = \frac{2}{\epsilon} \log(2/\delta)$$

# VC-dimension

Shattering: Ένα σύνολο σημείων  $X$  είναι shattered από μια κλάση  $H$  αν για κάθε labelling  $\gamma$  του  $X$  υπάρχει  $h \in H$  που να δίνει αυτό το labelling  
 SAS  $\forall i \quad h(x_i) = \gamma_i$

VC-dimension: Το μεγαλύτερο <sup>μέγεθος ενός</sup> συνόλου σημείων που είναι shattered από την  $H$ .

Παράδειγμα :

Thresholds

+



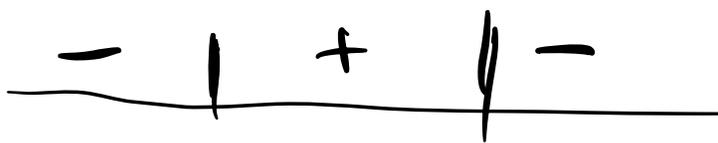
2 σημεία

+	+	✓
-	-	✓
-	+	✓
+	-	✗

VC-dim = 1



Thresholds  $\mu_c$   $\varphi_{op}$   
 $\eta$  intervals

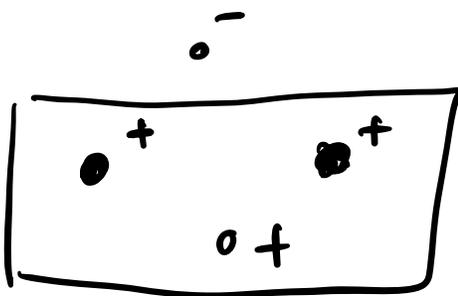
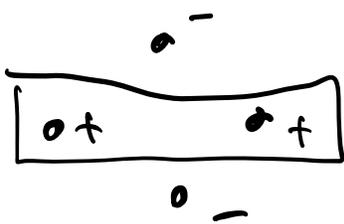


2  $\eta_{\mu_c}$  ✓  
 3  $\eta_{\varphi_{op}}$  ✗

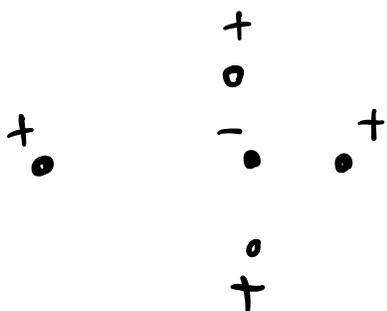


VC-dim = 2

Ορθογώνια (Axis-Aligned)



4  $\eta_{\mu_c}$  ✓



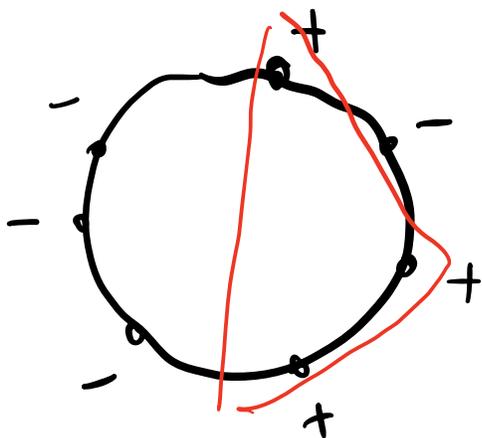
5  $\eta_{\mu_c}$  ✗  
 VC-dim = 4

Βασικό Θεώρημα του Statistical Learning  
Αν έχουμε μια κλάση  $H$  με  $VC\text{-dim} = d$

Τότε για binary classification  
χρειαζόμαστε

δείγματα για uniform convergence.  
$$m = \Theta\left(\frac{d + \log(1/\delta)}{\epsilon^2}\right)$$

$H = \{ \text{Όλα τα κύρια υποσύνολα του } \mathbb{R}^2 \}$



$$VC\text{-dim}(H) = \infty$$

Αν  $H$  έχει μέγεθος  $K$   
τότε

$$VC\text{-dim}(H) \leq \log K$$

Για Realizable Learning

$$m \leq O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon}\right)$$