

13/3/2025

## Παρασκευή 14/3 : Αξιούσα Γ εργασία

### Classification Problems

$$Y \in \{g_1, g_2, \dots, g_K\}$$

$X = (X_1, \dots, X_p)$  ανεξάρτητες μεταβλητές

$(X, Y)$  ανόικον οικονομική κατανομή

$$P(Y=g_k | X=x) \quad k=1, 2, \dots, K$$

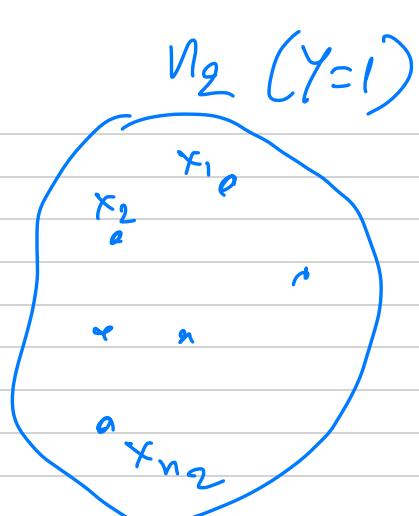
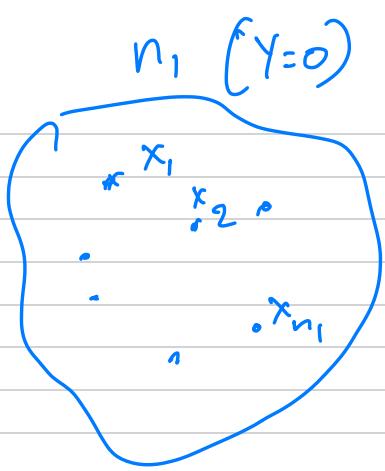
$$F(X | Y=g_k) = P(X \leq x | Y=g_k) \quad \left( \begin{array}{l} \text{για } p=1, X \text{ ουσιώδης} \\ f(x|g_k) = \delta_{x \leq g_k} \end{array} \right)$$

π.χ.  $X = (X_1, \dots, X_p)$  (κερδήστε ανθεκτικό/ανθεκτικό)

$$Y = \begin{cases} 1, & \text{εγκεκριθείσα} \\ 0, & \text{νησί} \end{cases}$$

$P(Y=1 | X=x) \leftarrow$  διατάξεις για προστασία  
δεδομένων των περιπτώσεων  
(δεξιά να το επεξιστεύει  
 $\Leftrightarrow$  "πρόβλεψη")

$F(X | Y=0) : \text{κατανομή } X \text{ για νησί}$       }  
 $F(X | Y=1) : \text{.. } X \text{ .. ανθεκτικός.}$       }  
}  $\Rightarrow$  Κανονικές  
δικτυακές  
(case-control)

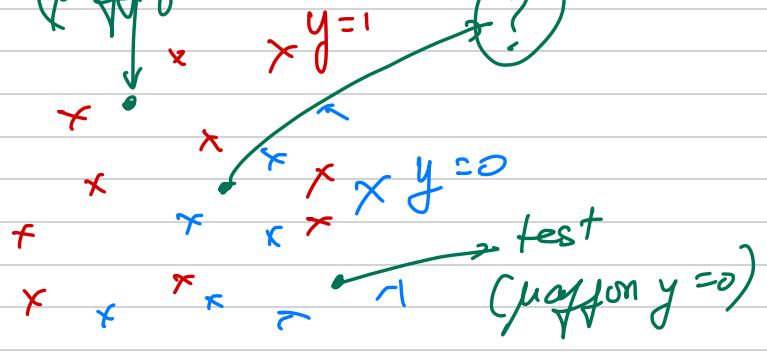


$$X = (X_1, X_2)$$

$$x_2 \uparrow$$



(μ<sub>effektiv</sub>=)



$$\text{Es gilt } \left. \begin{array}{l} P(Y=1) = p \\ P(Y=0) = 1-p \end{array} \right\} \text{nicht wahr } (Y)$$

$$\text{Bayes} \quad P(Y=1 | X=x) = \frac{P \cdot f(x | y=1)}{P f(x | y=1) + (1-p) f(x | y=0)}$$

## Kriptio / Enkription Anwesen

Teukis morfo  $Y \in \{g_1, \dots, g_K\}$  (Κεδονοφορια)

Training set  $(x_1, y_1), \dots, (x_n, y_n) \Rightarrow \dots \Rightarrow$

$\hat{G}(x)$  : noplazentiki ovdixymon }  
 $\hat{G}(x) \in \{g_1, \dots, g_K\}$  } noo tafii ?  
 eivai u  $\hat{G}(x)$

Test Set  $\{(x_0, y_0), \dots\}$

$\hat{y}_i = \hat{G}_i(x_i) \quad i \in \text{TestSet}$

Kriptio (Enkription anwesen) =  $L(y, \hat{G})$

Eow  $L(y_i, \hat{y}_i) = \begin{cases} 1, & y_i \neq \hat{y}_i \\ 0, & y_i = \hat{y}_i \end{cases} = 1 - I(y_i = \hat{y}_i)$

$L(y, \hat{G}) = \sum_{i \in \text{Test}} L(y_i, \hat{G}(x_i))$  = % κασ  
 | Test set | = ταξινομων  
 (misclassification rate)

MPE =  $E_{(X,Y)}(L(Y, \hat{G}(X))) \leftarrow$

Tia  $L = 1 - I(y_i = \hat{y}_i)$

$$\min_{\hat{G}} \left[ E \left[ L(Y, \hat{G}(x)) \right] \right] = E_x \left[ E_{Y|x} \left[ L(Y, \hat{g}(x)) | x \right] \right]$$

$$= E_x \left[ \sum_{k=1}^K P(Y=g_k|x) \cdot L(g_k, \hat{G}(x)) \right]$$

Bayes risk

Infrazi bedeckung

$$\forall x : \min_{\hat{G}(x)} \sum_{k=1}^K P(Y=g_k|x) \cdot L(g_k, \hat{G}(x))$$

$$\hat{G} \in \{g_1, \dots, g_K\}$$

$$\text{Or } L(y, \hat{y}) = 1 - I(y = \hat{y})$$

$$\text{Toze} \sum_{k=1}^K P(Y=g_k|x) L(g_k, \hat{G}(x)) =$$

$$= \sum_{g_k \neq \hat{G}(x)} P(Y=g_k|x) \cdot 1 = \boxed{1 - P(Y=\hat{G}(x)|x)}$$

Bézirki Próbalegy  $\boxed{X=x}$

$$\min_{\hat{G}} \left\{ 1 - P(Y=\hat{G}|x) \right\} \Leftrightarrow \boxed{\max_{\hat{G}} P(Y=\hat{G}|x)}$$

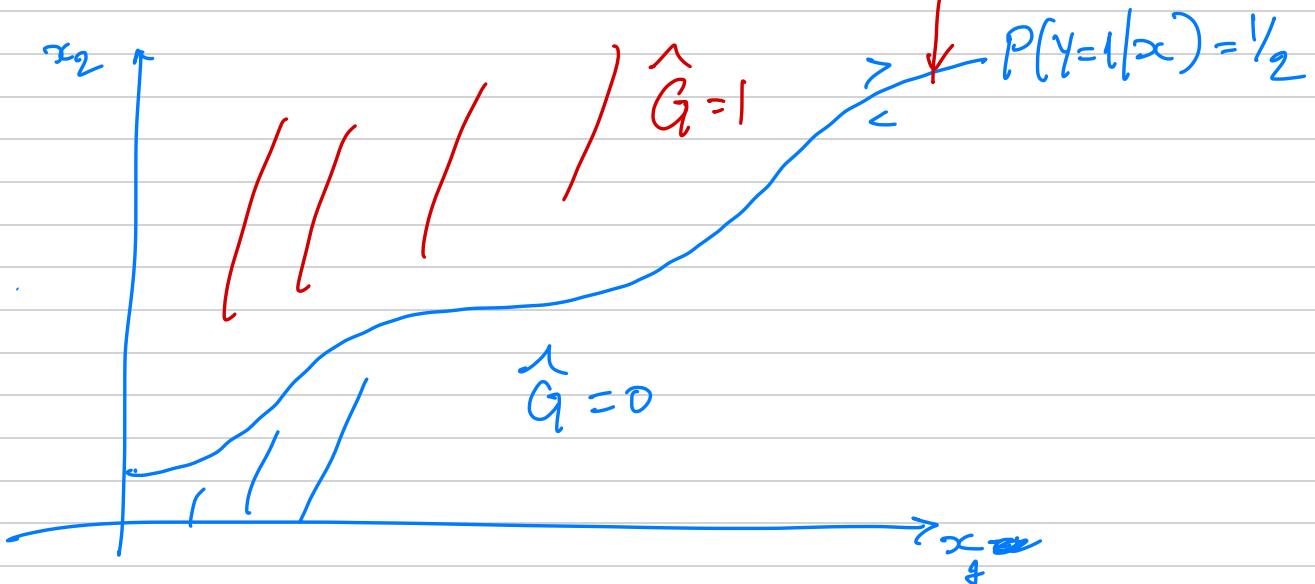
$$\Rightarrow \hat{G}(x) = \operatorname{argmax}_{g \in \{g_1, \dots, g_K\}} P(Y=g|x)$$

Bayes classifier

Für  $X = (X_1, X_2)$ ,  $Y \in \{0, 1\}$

$$\hat{G}(x) = \begin{cases} 0, & P(Y=1|X=x) < \frac{1}{2} \\ 1, & P(Y=1|X=x) > \frac{1}{2} \end{cases}$$

*Grünes Bayes  
Blaues Bayes  
(Bayes boundary)  
ANNO STO??!*



### Parametrische Modelle

1) Logistic Regression ( $Y \in \{0, 1\}$ )  
 $\rightarrow P(Y=1|x) = p(x)$

$$\log \frac{P(x)}{1-P(x)} = b_0 + b_1^T \cdot x$$

mit  $p(x)$

2) Generative Models (Discriminant Analysis)

$$F(x|0) \sim \dots \text{ a priori } \text{ a priori}$$

$$F(x|1) \sim \dots \text{ " } \text{ a priori}$$

Training set  $\Rightarrow \dots \left. \begin{array}{l} \hat{F}(x|0) \\ \hat{F}(x|1) \end{array} \right\} \Rightarrow \text{Bayes} \cdot \hat{P}(Y=1|x=x)$

n.x. or  $X = X_1$

$$X|Y=0 \quad X \sim \mathcal{N}(\mu_0, \sigma^2)$$

$$X|Y=1 \quad X \sim \mathcal{N}(\mu_1, \sigma^2)$$

if  $(\mu_1 > \mu_0)$

$$f(x|y=0) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu_0)^2}{2\sigma^2}}$$

$$f(x|y=1) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu_1)^2}{2\sigma^2}}$$

$$P(Y=1|X=x) = \frac{p f(x|1)}{p f(x|1) + (1-p) f(x|0)}$$

$$P(Y=0|X=x) = \frac{(1-p) f(x|0)}{\dots}$$

Kalibrier Bayes classifier

$$\hat{G}_1 = 1 \quad \text{or} \quad p f(x|1) > (1-p) f(x|0) \Rightarrow$$

$$\Leftrightarrow \frac{f(x|1)}{f(x|0)} > \frac{1-p}{p} \Leftrightarrow$$

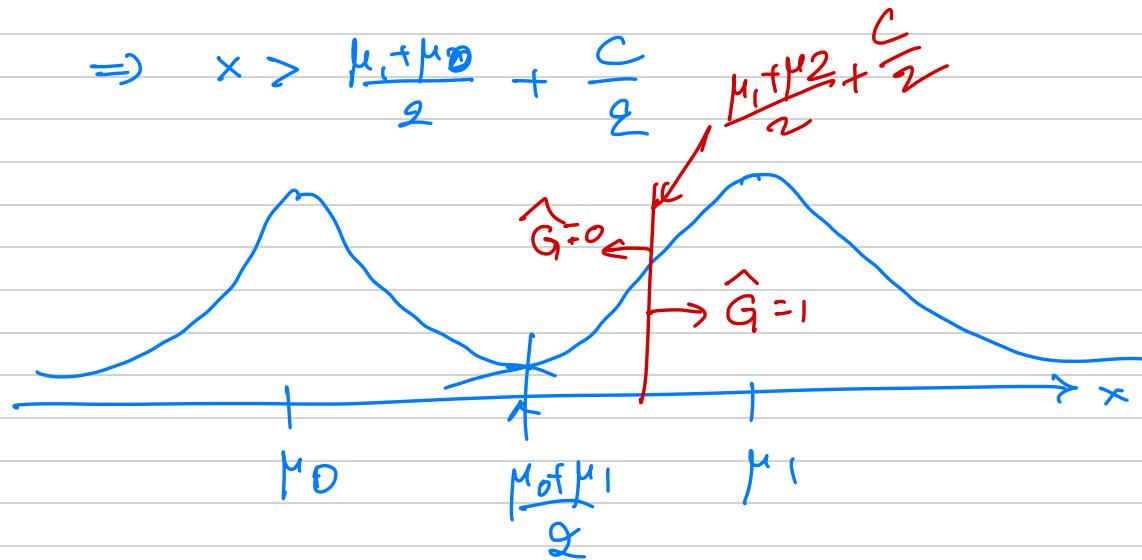
$$\log \frac{f(x|1)}{f(x|0)} > \log \frac{1-p}{p} = A$$

$$\log f(x|1) - \log f(x|0) > A \Leftrightarrow$$

$$\Leftrightarrow \frac{(x-\mu_0)^2}{2\sigma^2} - \frac{(x-\mu_1)^2}{2\sigma^2} > A \Leftrightarrow$$

$$\underbrace{(\mu_1 - \mu_0)}_{>0} (2x - (\mu_1 + \mu_0)) > 2A\sigma^2 \Rightarrow$$

$$\Rightarrow 2x - (\mu_1 + \mu_0) > \frac{2A\sigma^2}{\mu_1 - \mu_0} = C \Rightarrow$$

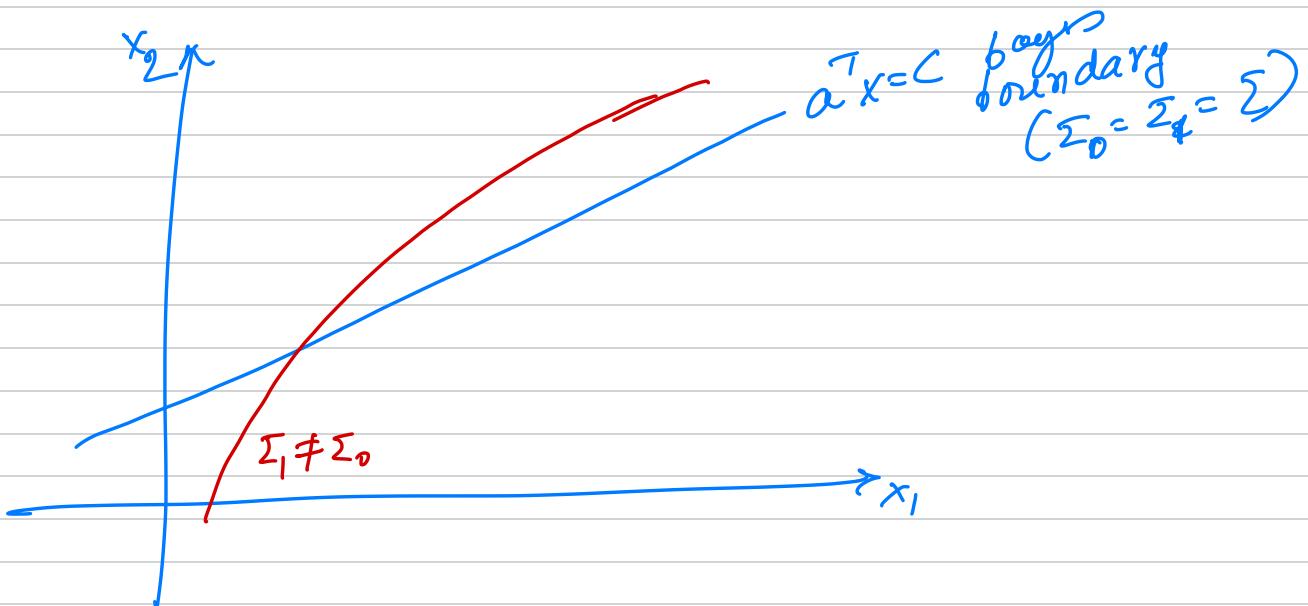


$$\text{Av } X = (X_1, X_2)$$

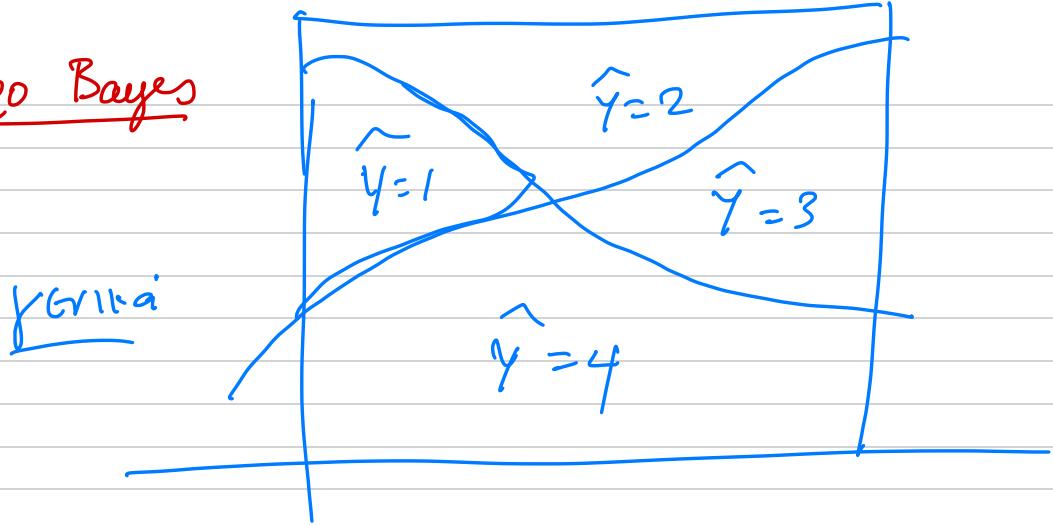
$$X | Y=0 \sim \mathcal{N}(\mu_0, \Sigma) \quad \left. \right\} \Rightarrow \dots$$

$$X | Y=1 \sim \mathcal{N}(\mu_1, \Sigma)$$

$$P(Y=1 | X=x) > P(Y=0 | X=x) \Leftrightarrow \boxed{a^T x \geq C}$$

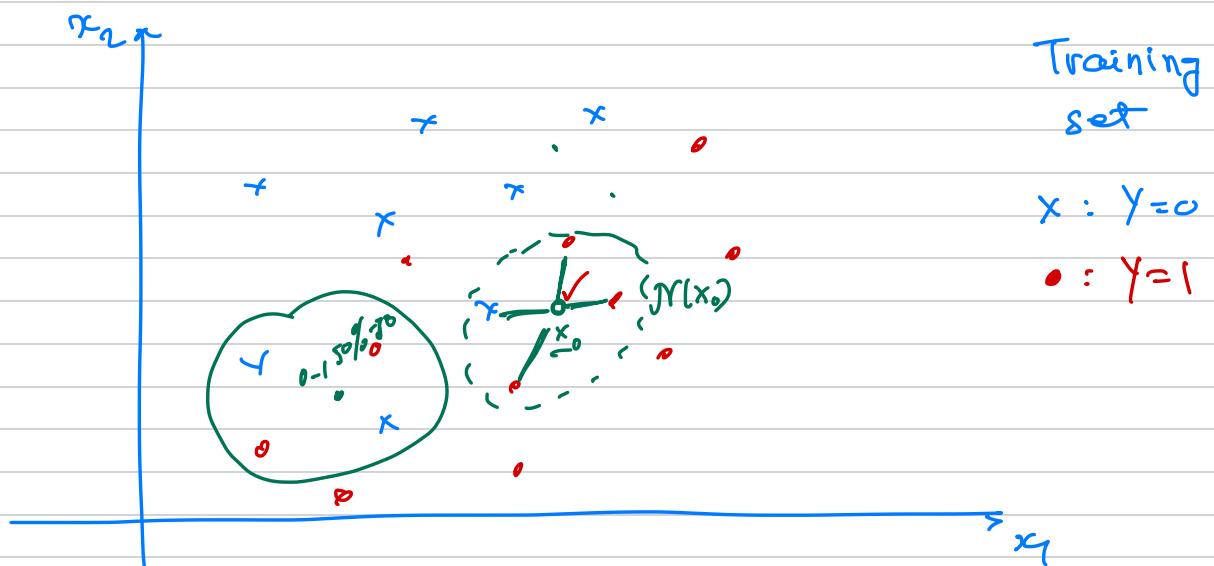


Σύρπο Βάιες



Μέθοδος kNN (k nearest neighbor)

Εκτίμηση  $P(Y=j | X=x_0)$  με παραγενότητα



Εσω  $k=4$

$N_0(x_0) = \{4 \text{ κοντινότερα σημεία του training set σε } x_0 \text{ (ws npō x)}\}$

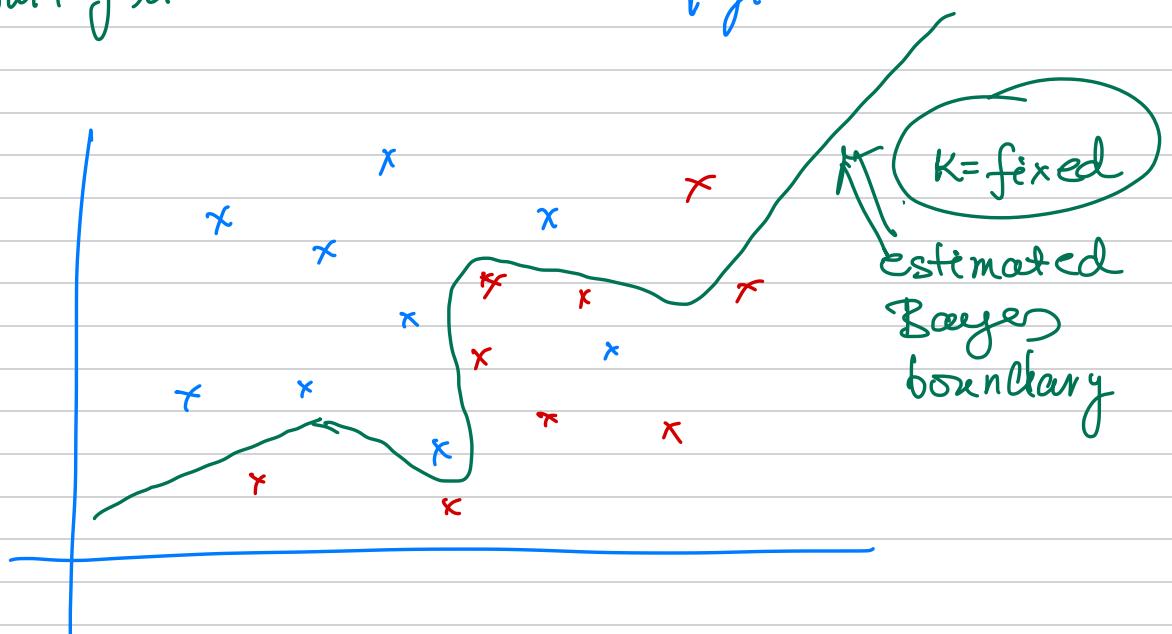
κοντινότητα :  $\|x - x_0\|$

$$|N_0(x_0)| = k \Leftrightarrow \|x_i - x_0\| < \|x_j - x_0\| \quad \forall i \in \mathcal{N}_0, j \notin \mathcal{N}_0.$$

$$\hat{G}(x_0) = g_i \text{ or } \frac{1}{K} \sum_{j=1}^K I(y_j = g_i) \leftarrow \max_i$$

↓  
Estimated  
and training set

% nápr. obojí  
not sivá  $y=g_i$



$$K=1 \quad n' \quad K \rightarrow n \quad (n = |\text{Training set}|)$$

eufiga  
+ -

I Ozav  $K \rightarrow n$  (negado)

$$+x_0 : \% (y=1) = \text{ocetip} = \% (y=1) \text{ obojí} = q_1$$

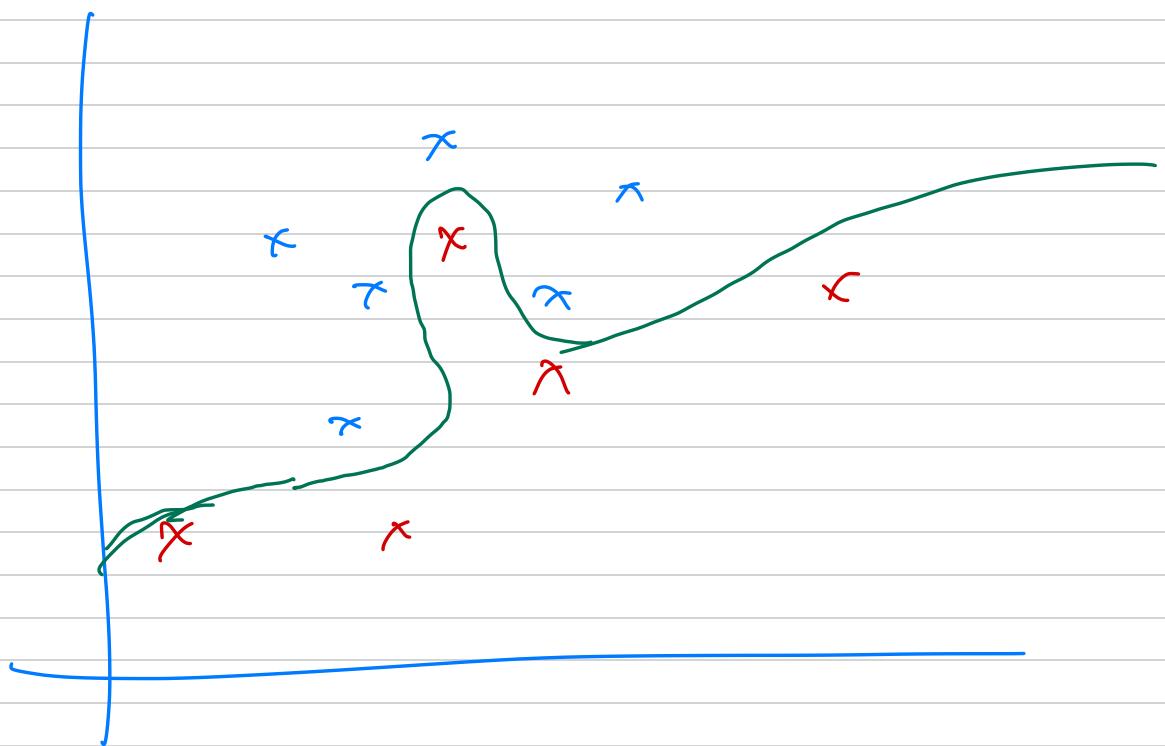
training set

$$\text{Ozav } q_1 > \frac{1}{2} \Rightarrow \hat{G}(x_0) = 1 + x_0$$

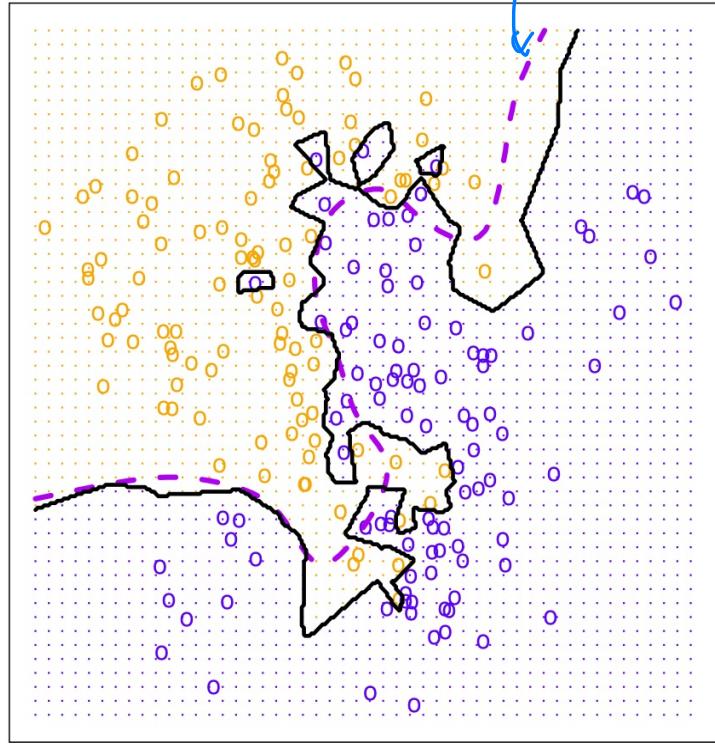
$$q_1 < \frac{1}{2} \quad \hat{G}(x_0) = 0$$

} exazioni  
eufiga

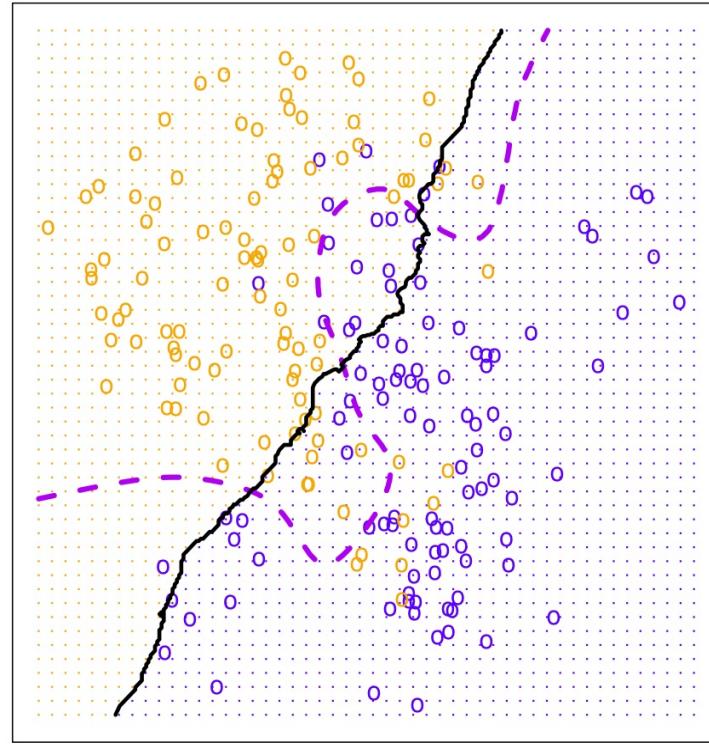
2 Ozav  $k \rightarrow 1$



KNN:  $K=1$



KNN:  $K=100$



Dawen  
Niralo  
Bayer

## Прибаўшча

$$Y \in \{0, 1\}$$

$$X = (X_1, X_2)$$

$$(X_1, X_2) \text{iid} \sim U(0,1)$$

$$P(Y|_1 | X = (x_1, x_2)) = p(x_1, x_2)$$

$$Y((X_1=x_1, X_2=x_2)) \sim \text{Ber}(p(x_1, x_2))$$

## Праоцэсівон

Чнодзір аздырві овідрум:

$$p(x_1, x_2) = \begin{cases} \frac{1}{2} \left( \frac{x_2}{b(x_1)} \right)^\alpha, & x_2 < b(x_1) \\ \frac{1}{2} \left( \frac{1-x_2}{1-b(x_1)} \right)^\alpha, & x_2 \geq b(x_1) \end{cases}$$

$$\text{бір} \quad b(x_1) = 0.7 + 0.1 x_1 - 0.8 x_1^2 + 4x_1^3 - 2.8x_1^4$$

Bayes boundary :  $\{(x_1, x_2) : p(x_1, x_2) = \frac{1}{2}\}$

## Алгебраічны Праоцэсівон

Есно  $n=1000$  (негэдээс training set)

$$\alpha = 2$$

$$\kappa = 10 \quad (\text{тун})$$

- i) Дзялкоўшы  $N$ -training sets па  
Праоцэсівону і з'ю аднаўш  
калькуляція.  $(N \rightarrow \infty)$

2) Διαφορετική Λ test set γεγενήσιμης  
 $n_{test} = 100$  από την ίδια παρασκευή.

3) Τια τις ~~τις~~ training set  $r=1, 2, \dots, N$

④ Εγκατεύθυντες kNN ( $k=10$ ) σε  $r$

⑤ βρίσκω προβλέψεις  $\hat{y}_j$ ,  $j \in \text{Test set}$

⑥ υπολογίζω ~~τις~~ % misclassification

% λαχανικ. test set στοντος  $\hat{y} = y$

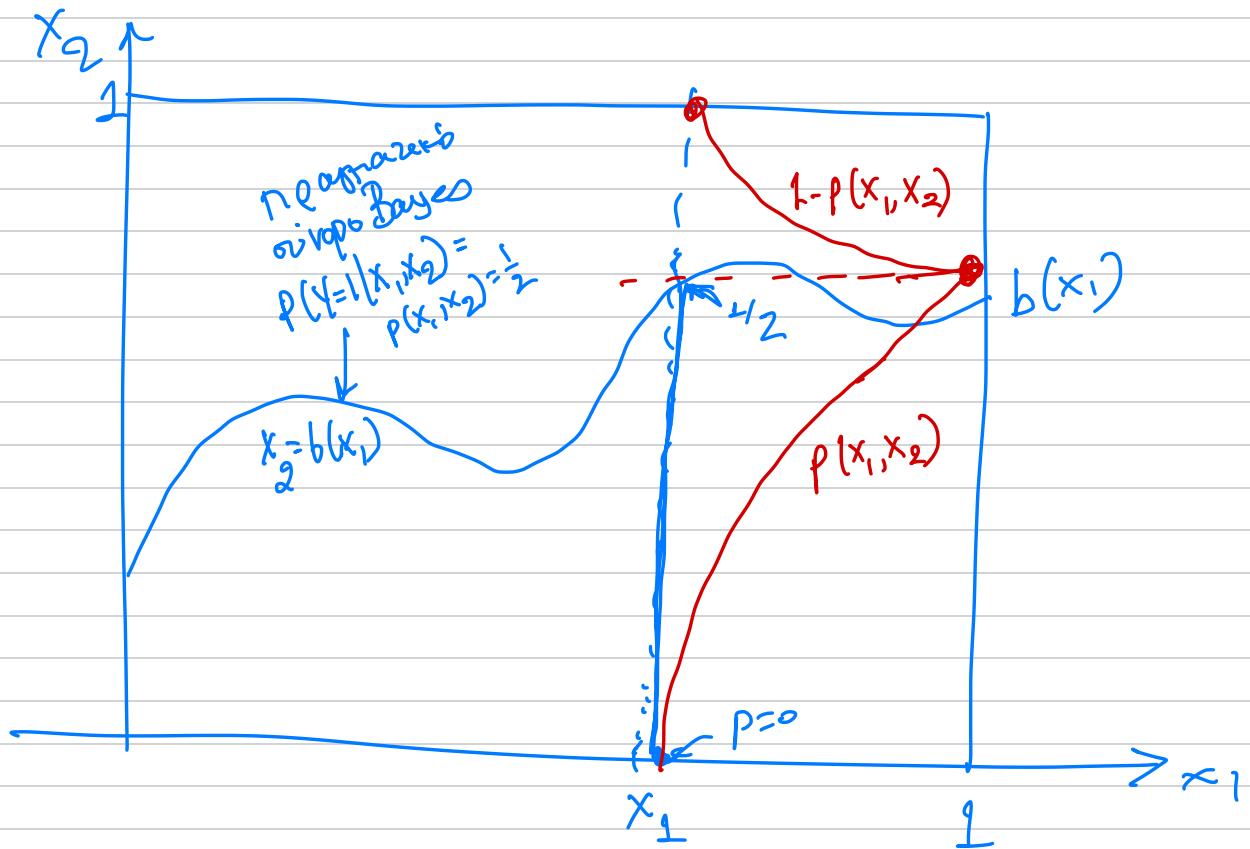
MPE<sub>r</sub>

$$4) \hat{MPE} = \frac{1}{N} \sum_{r=1}^N MPE_r$$

$MPE_r$ ,  $r=1, 2, \dots$  iid

$$\hat{MPE} \xrightarrow[N \rightarrow \infty]{} MPE (\mu, n, 1)$$

(Monte Carlo Simulation)



$$p(x_1, x_2) = \begin{cases} \left(\frac{x_2}{b(x_1)}\right)^2 / 2 & x_2 \leq b(x_1) \\ 1 - \left(\frac{1-x_2}{1-b(x_1)}\right)^2 / 2 & , x_2 > b(x_1) \end{cases}$$