

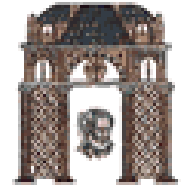
# Αναγνώριση Προτύπων

## Εκτίμηση Παραμέτρων (Parameter Estimation)

*Χριστόδουλος Χαμζάς*

*Τα περιεχόμενα των παρουσιάσεων προέρχονται κυρίως από τις παρουσιάσεις του αντίστοιχου διδασκέου μαθήματος του καθ. Παναγιώτη Τσακαλίδη, Τμ. Επιστήμης Υπολογιστών, Παν. Κρήτης. Το πρωτογενές υλικό βρίσκεται στην σελίδα <http://www.csd.uoc.gr/~hy473/> και βασίζεται στο βιβλίο: "Pattern Classification", R.O. Duda, P.E. Hart, D.G. Stork, Wiley, 2<sup>nd</sup> Ed., 2001*

Τμήμα Ηλεκτρολόγων Μηχανικών & Μηχανικών Υπολογιστών

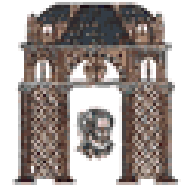


# Εκτίμηση Παραμέτρων

- Ο Μπεϋζιανός ταξινομητής δεν μπορεί να χρησιμοποιηθεί όταν δεν είναι γνωστές οι συναρτήσεις πυκνότητας πιθανότητας  $p(x/\omega_j)$  &  $P(\omega_j)$ .
- Οι κατανομές μπορούν να εκτιμηθούν εάν υπάρχουν επαρκή δεδομένα → Δύσκολο!
- Εάν είναι γνωστή η μορφή της κατανομής, για παράδειγμα κανονική, αλλά όχι οι παράμετροι της, π.χ. η μέση τιμή και η διασπορά της, το πρόβλημα ανάγεται σε αυτό της εκτίμησης των παραμέτρων.



- Υπάρχουν επίσης δύο τεχνικές μάθησης
- 1. Εκπαίδευση με επίβλεψη (supervised learning)
  - Υπάρχει αρχικό δείγμα στο οποίο γνωρίζουμε σε ποια κατηγορία ανήκει το καθένα (π.χ. έχουμε ένα δείγμα από χαρακτηριστικά στο οποίο ξέρουμε **ποια** είναι από λαυράκια και ποια είναι από σολωμούς)
- 2. Εκπαίδευση δίχως επίβλεψη (unsupervised learning)
  - Δεν υπάρχει αρχικό δείγμα εκπαίδευσης απλά ξέρουμε ότι τα δείγματα μας προέρχονται από γνωστές κατηγορίες (π.χ. έχουμε ένα δείγμα από χαρακτηριστικά και απλά ξέρουμε ότι προέρχονται μόνο από λαυράκια και από σολωμούς)
  - Το πρόβλημα αυτό θα το αφήσουμε για άλλο μάθημα.

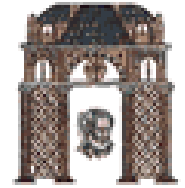


## Εκτίμηση Παραμέτρων (2)

- Υπάρχουν δύο βασικές τεχνικές:
  1. Εκτίμηση Μέγιστης Πιθανοφάνειας (Maximum Likelihood Estimation)
  2. Εκτίμηση παραμέτρων κατά Bayes (Bayesian Parameter Estimation)

Τα αποτελέσματα και των 2 τεχνικών είναι περίπου ίδια αλλά η θεωρητική προσέγγιση του προβλήματος είναι διαφορετική

1. Η εκτίμηση μέγιστης πιθανοφάνειας θεωρεί τις τιμές των αγνώστων παραμέτρων «άγνωστες σταθερές» και προσπαθεί να βρεί ποιες είναι αυτές ώστε να μεγιστοποιείται η πιθανότητα να πάρουμε το συγκεκριμένο «δείγμα» τιμών.
2. Η εκτίμηση παραμέτρων κατά Bayes, αντίθετα θεωρεί ότι οί άγνωστοι παράμετροι είναι τυχαίες μεταβλητές με γνωστή από πριν αρχική πυκνότητα κατανομής (prior) και το δείγμα τιμών που έχουμε λάβει μας οδηγεί σε βελτίωση (posterior) της αρχικής μας κατανομής.

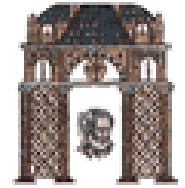


# Εκτίμηση Μέγιστης Αληθοφάνειας

- Υποθέτουμε ότι παίρνουμε τυχαία επιλεγμένα δείγματα από μια δεδομένη κατανομή της οποίας οι παράμετροι είναι άγνωστοι. Συμβολίζουμε το σύνολο των παραμέτρων με το διάνυσμα  $\theta$ .
- Εάν για παράδειγμα γνωρίζουμε ότι η κατανομή είναι κανονική αλλά δεν ξέρουμε τον μέσο και τη διασπορά της, τότε

$$\theta = (\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (\mu_1, \dots, \mu_d, \sigma_1^2, \dots, \sigma_d^2, \text{cov}(x_m, x_n); m, n = 1, \dots, d; m \succ n)$$

$$d + \frac{d(d+1)}{2} \text{ παράμετροι}$$



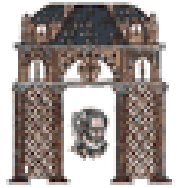
# Πρόβλημα MLE

- Για κάθε μία από τις κλάσεις, υπολόγισε το διάνυσμα παραμέτρων  $\theta$  χρησιμοποιώντας ένα σύνολο δεδομένων εκπαίδευσης  $D^n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  το οποίο έχει  $n$  ανεξάρτητα δείγματα (i.i.d):

$$p(D^n | \theta) = \prod_{k=1}^n p(\mathbf{x}_k | \theta)$$

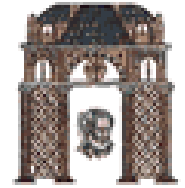
Πιθανοφάνεια του  $\theta$   
αναφορικά με το δείγμα  $D^n$

- Η εκτίμηση μέγιστης πιθανοφάνειας του  $\theta$  είναι εκείνη η τιμή που μεγιστοποιεί την παραπάνω συνάρτηση.
- Διαισθητικά, αντιστοιχεί στην τιμή του  $\theta$  που «συμφωνεί» όσο το δυνατό καλύτερα με τα δείγματα.

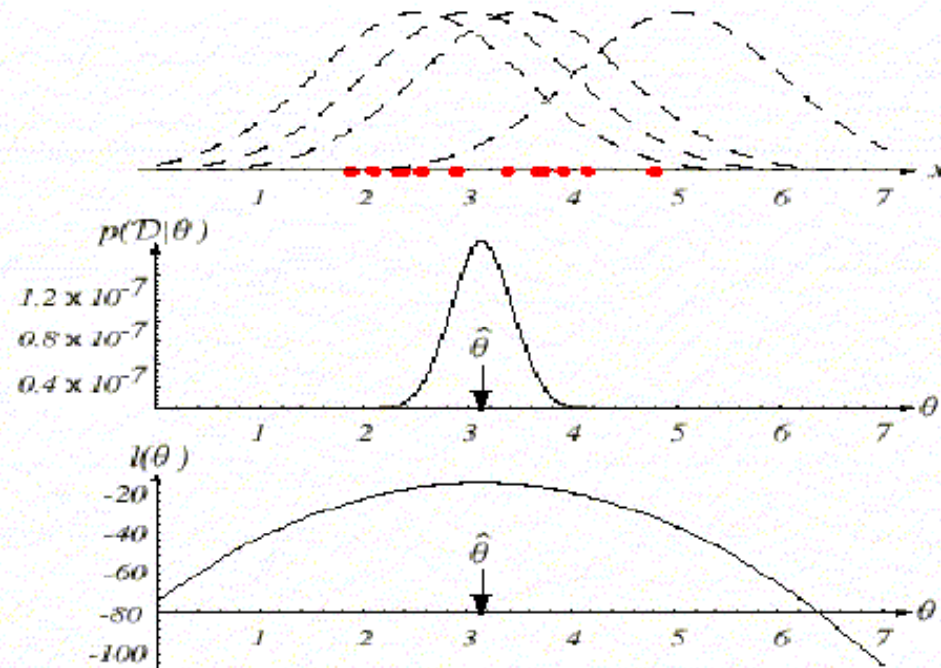


# Παράδειγμα με νόμισμα

- Ρίχνουμε ένα νόμισμα  $N$  φορές και παρατηρούμε το αποτέλεσμα.
- Ποια είναι η εκτίμηση της πιθανότητας να έχουμε «γράμματα» ?

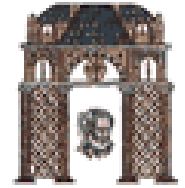


# Συνάρτηση Πιθανοφάνειας



**FIGURE 3.1.** The top graph shows several training points in one dimension, known or assumed to be drawn from a Gaussian of a particular variance, but unknown mean. Four of the infinite number of candidate source distributions are shown in dashed lines. The middle figure shows the likelihood  $p(\mathcal{D}|\theta)$  as a function of the mean. If we had a very large number of training points, this likelihood would be very narrow. The value that maximizes the likelihood is marked  $\hat{\theta}$ ; it also maximizes the logarithm of the likelihood—that is, the log-likelihood  $l(\theta)$ , shown at the bottom. Note that even though they look similar, the likelihood  $p(\mathcal{D}|\theta)$  is shown as a function of  $\theta$  whereas the conditional density  $p(x|\theta)$  is shown as a function of  $x$ . Furthermore, as a function of  $\theta$ , the likelihood  $p(\mathcal{D}|\theta)$  is not a probability density function and its area has no significance. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.





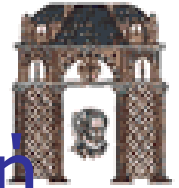
# Λογάριθμος της Πιθανοφάνειας (Log-Likelihood)

- Το  $\theta$  που μεγιστοποιεί την συνάρτηση πιθανοφάνειας, μεγιστοποιεί επίσης και τον λογάριθμό της, με τον οποίο μπορούμε πολλές φορές να δουλέψουμε ευκολότερα:

$$l(\boldsymbol{\theta}) \equiv \ln p(D^n | \boldsymbol{\theta}) = \sum_{k=1}^n \ln p(\mathbf{x}_k | \boldsymbol{\theta})$$

Το  $\theta$  που μεγιστοποιεί αυτή τη συνάρτηση μπορεί να υπολογισθεί θέτοντας την παράγωγο του ως προς  $\theta$  ίση με το μηδέν, και λύνοντας την εξίσωση ως προς  $\theta$ .  $\hat{\boldsymbol{\theta}}_{ML} = \arg \max_{\boldsymbol{\theta}} l(\boldsymbol{\theta})$

$$\nabla_{\boldsymbol{\theta}} l(\boldsymbol{\theta}) = \sum_{k=1}^n \nabla_{\boldsymbol{\theta}} \ln p(\mathbf{x}_k | \boldsymbol{\theta}) = \mathbf{0}$$



## Ειδικές Περιπτώσεις – Κανονική κατανομή

- Περίπτωση κανονικής κατανομής με άγνωστο  $\mu$

$$\hookrightarrow p(\mathbf{x}_k | \boldsymbol{\mu}) \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\ln p(\mathbf{x}_k | \boldsymbol{\mu}) = -\frac{1}{2} \ln[(2\pi)^d |\boldsymbol{\Sigma}|] - \frac{1}{2} (\mathbf{x}_k - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_k - \boldsymbol{\mu})$$

$$\nabla_{\boldsymbol{\theta}} \ln p(\mathbf{x}_k | \boldsymbol{\mu}) = \boldsymbol{\Sigma}^{-1} (\mathbf{x}_k - \boldsymbol{\mu})$$

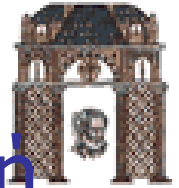
- Ο εκτιμητής μέγιστης πιθανοφάνειας του  $\boldsymbol{\mu}$  ικανοποιεί:

$$\sum_{k=1}^n \boldsymbol{\Sigma}^{-1} (\mathbf{x}_k - \hat{\boldsymbol{\mu}}) = 0$$

- Πολλαπλασιάζοντας από αριστερά με  $\boldsymbol{\Sigma}$  και λύνοντας, παίρνουμε:

$$\hat{\boldsymbol{\mu}}_{ML} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

Δηλαδή είναι ο αριθμητικός μέσος όρος των δειγμάτων εκπαίδευσης!



# Ειδικές Περιπτώσεις – Κανονική κατανομή

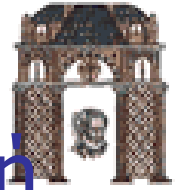
➤ Περίπτωση κανονικής κατανομής με αγνώστους  $\mu$  και  $\sigma^2$ :

$$\Rightarrow \boldsymbol{\theta} = (\theta_1, \theta_2) = (\mu, \sigma^2)$$

$$l(\boldsymbol{\theta}) = \ln P(x_k | \boldsymbol{\theta}) = -\frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2} (x_k - \theta_1)^2$$

$$\nabla_{\boldsymbol{\theta}} l = \begin{pmatrix} \frac{\partial}{\partial \theta_1} (\ln P(x_k | \boldsymbol{\theta})) \\ \frac{\partial}{\partial \theta_2} (\ln P(x_k | \boldsymbol{\theta})) \end{pmatrix} = \mathbf{0}$$

$$\begin{cases} \frac{1}{\theta_2} (x_k - \theta_1) = 0 \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} = 0 \end{cases}$$



# Ειδικές Περιπτώσεις – Κανονική κατανομή

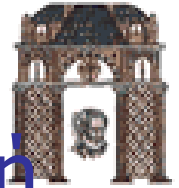
Λαμβάνοντας υπ' όψην όλα τα δείγματα :

$$\left\{ \begin{array}{l} \sum_{k=1}^n \frac{1}{\hat{\theta}_2} (x_k - \hat{\theta}_1) = 0 \end{array} \right. \quad (1)$$

$$\left\{ \begin{array}{l} -\sum_{k=1}^n \frac{1}{\hat{\theta}_2} + \sum_{k=1}^n \frac{(x_k - \hat{\theta}_1)^2}{\hat{\theta}_2^2} = 0 \end{array} \right. \quad (2)$$

Συνδυάζοντας τις (1) και (2), παίρνουμε:

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k \quad ; \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})^2$$



# Ειδικές Περιπτώσεις – Κανονική κατανομή

## ➤ Bias

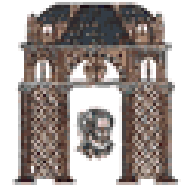
↪ Ο ML εκτιμητής για τη διασπορά  $\sigma^2$  μεροληπτεί (biased)

$$E\left[\frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})^2\right] = \frac{n-1}{n} \sigma^2 \neq \sigma^2$$

↪ Ένας στοιχειώδης αμερόληπτος εκτιμητής για τον πίνακα συνδιασποράς,  $\Sigma$ , είναι:

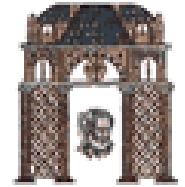
$$\mathbf{C} = \frac{1}{n-1} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu})(\mathbf{x}_k - \hat{\boldsymbol{\mu}})^t$$

*Sample covariance matrix*



## Εκτίμηση κατά Bayes

- Σε αντίθεση με τον MLE, όπου κάναμε την υπόθεση ότι οι άγνωστες παράμετροι έχουν σταθερές τιμές, ο εκτιμητής κατά Bayes (BE) υποθέτει ότι οι άγνωστες παράμετροι είναι τυχαίες μεταβλητές και ακολουθούν μία εκ των προτέρων γνωστή σ.π.π.
- Επομένως, ο BE υπολογίζει μια κατανομή των τιμών του  $\theta$  και όχι τις τιμές αυτές καθ' αυτές. Ο BE παρέχει περισσότερη πληροφορία, όμως συχνά είναι δύσκολο να υπολογισθεί.
- Η ύπαρξη μετρήσεων εκπαίδευσης (training data) επιτρέπει την μετατροπή της εκ των προτέρων πληροφορίας σε εκ των υστέρων σ.π.π. → φαινόμενο της εκμάθησης (Bayesian learning) όπου κάθε νέα παρατήρηση οξύνει την εκ των υστέρων σ.π.π.



# Εκτίμηση κατά Bayes

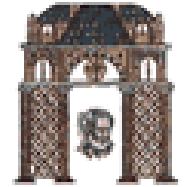
➤ Μπεϋζιανή Εκμάθηση για προβλήματα ταξινόμησης προτύπων.

↳ Ο υπολογισμός των εκ των υστέρων σ.π.π. αποτελεί τη βάση της Μπεϋζιανής ταξινόμησης.

↳ Στόχος: Υπολογισμός των  $P(\omega_i | \mathbf{x}, D)$  δεδομένου του συνόλου δειγμάτων εκπαίδευσης  $D = \{D_1, \dots, D_c\}$ , όπου τα δείγματα στο σύνολο  $D_j$  αντιστοιχούν στην κλάση  $j, j=1, \dots, c$ .

↳ Για κάθε νέο, αταξινομήτο δείγμα,  $\mathbf{x}$ , ο κανόνας του Bayes δίνει:

$$P(\omega_i | \mathbf{x}, D) = \frac{p(\mathbf{x} | \omega_i, D)P(\omega_i | D)}{\sum_{j=1}^c p(\mathbf{x} | \omega_j, D)P(\omega_j | D)}$$



# Εκτίμηση κατά Bayes

➤ Υποθέτουμε ότι

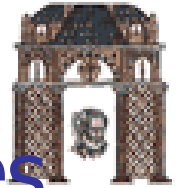
↳ Οι εκ των προτέρων  $P(\omega_i)$  είναι γνωστές, οπότε  $P(\omega_i/D) = P(\omega_i)$ .

↳ Μόνο τα δείγματα του συνόλου  $D_i$  έχουν πληροφορία για τη σ.π.π.  $p(\mathbf{x}/\omega_i, D_i) \rightarrow$  προκύπτουν  $c$  ανεξάρτητα προβλήματα εκτίμησης των  $p(\mathbf{x}/\omega_i, D_i)$ , η οποία μπορεί να γραφεί και ως  $p(\mathbf{x}/D_i)$ .

$$P(\omega_i | \mathbf{x}, D_i) = \frac{p(\mathbf{x} | \omega_i, D_i)P(\omega_i)}{\sum_{j=1}^c p(\mathbf{x} | \omega_j, D_j)P(\omega_j)}$$

Αυτή τη συνάρτηση θέλουμε να εκτιμήσουμε (γράφεται και ως  $p(\mathbf{x}/D_i)$ )





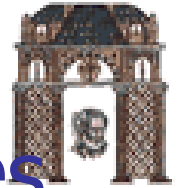
# Γενική μεθοδολογία εκτίμησης κατά Bayes

- Για κάθε κλάση, γνωρίζουμε τη μορφή της σ.π.π.  $p(\mathbf{x}|\boldsymbol{\theta})$ , αλλά η τιμή του διανύσματος παραμέτρων  $\boldsymbol{\theta}$  είναι άγνωστη.
- Έχουμε κάποια αρχική γνώση για το  $\boldsymbol{\theta}$  με τη μορφή της εκ των προτέρων (a priori) σ.π.π.  $p(\boldsymbol{\theta})$ .
- Για κάθε κλάση, διαθέτουμε ένα σύνολο  $D^n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  από  $n$  στατιστικώς ανεξάρτητα δείγματα. Τότε:

$$p(\mathbf{x} / D) = \int_{\boldsymbol{\theta}} p(\mathbf{x} / \boldsymbol{\theta}) p(\boldsymbol{\theta} / D) d\boldsymbol{\theta}$$

Σχέση κλειδί: Συνδέει την δεσμευμένη σ.π.π.  $p(\mathbf{x}/D)$  με την εκ των υστέρων σ.π.π.  $p(\boldsymbol{\theta}/D)$  του διανύσματος των παραμέτρων. Δηλώνει ότι η  $p(\mathbf{x}/D)$  είναι ένας γραμμικός συνδιασμός των  $p(\mathbf{x}/\boldsymbol{\theta})$  με βάρη τις  $p(\boldsymbol{\theta}/D)$ .

Αν η  $p(\boldsymbol{\theta}/D)$  έχει ένα απότομο μοναδικό μέγιστο στο  $\boldsymbol{\theta}^*$  τότε  $p(\mathbf{x}/D) \approx p(\mathbf{x}/\boldsymbol{\theta}^*)$



# Γενική μεθοδολογία εκτίμησης κατά Bayes

- Για τον υπολογισμό της  $p(\theta/D)$ , έχουμε ότι:

$$p(\theta / D) = \frac{p(D / \theta) p(\theta)}{\int_{\theta} p(D / \theta) p(\theta) d\theta}$$

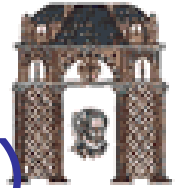
- Λόγω της ανεξαρτησίας των δειγμάτων εκπαίδευσης,

$$p(D / \theta) = \prod_{k=1}^n p(\mathbf{x}_k / \theta)$$

- Αν η  $p(D/\theta)$  είναι επικεντρωμένη γύρω από το  $\theta^*$  με μεγάλη κορυφή σε αυτό το σημείο, και αν η  $p(\theta^*)$  δεν είναι 0, τότε και η  $p(\theta/D)$  έχει μεγάλη κορυφή στο  $\theta^*$ , και επομένως θα είναι

$$p(\mathbf{x}/D) \approx p(\mathbf{x}/\theta^*)$$

- Αλλά το σημείο  $\theta^*$ , όπως περιγράφεται παραπάνω, είναι ο εκτιμητής μέγιστης πιθανοφάνειας του  $\theta$ !!



# Μπεϋζιανή εκμάθηση (Bayesian Learning)

- Ένα θέμα που ενδιαφέρει είναι αυτό του υπολογισμού και της σύγκλισης της ακολουθίας των σ.π.π.  $p(\theta/D^n)$ , όπου επαναφέραμε τον δείκτη  $n$  του αριθμού των δειγμάτων εκπαίδευσης στο σύνολο  $D^n$ .

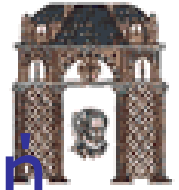
$$p(D^n / \theta) = \prod_{k=1}^n p(\mathbf{x}_k / \theta) = p(\mathbf{x}_n / \theta) p(D^{n-1} / \theta)$$

- Επομένως,

$$p(\theta / D^n) = \frac{p(\mathbf{x}_n / \theta) p(\theta / D^{n-1})}{\int_{\theta} p(\mathbf{x}_n / \theta) p(\theta / D^{n-1}) d\theta}$$

$$p(\theta / D^0) = p(\theta)$$

- Η παραπάνω σχέση δημιουργεί μία ακολουθία σ.π.π.  $p(\theta/x_1)$ ,  $p(\theta/x_1, x_2)$ , ...,  $p(\theta/x_1, \dots, x_n) \rightarrow$  Αναδρομική Μπεϋζιανή μεθοδολογία εκμάθησης, Bayesian recursive (or incremental) learning.



# Μπεϋζιανή εκμάθηση–Κανονική Κατανομή

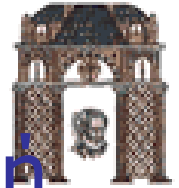
- Πρόβλημα: Υπολογισμός των σ.π.π.  $p(\theta/D^n)$  και  $p(x/D^n)$  όταν υποθέτουμε ότι  $p(x/\theta)=p(x/\mu)=N(\mu,\sigma^2)$  (δηλ.  $\theta=\mu$ ) και  $p(\mu)=N(\mu_0,\sigma_0^2)$ . Το  $\mu_0$  είναι η εκ των προτέρων καλύτερη γνώση μας για το  $\mu$  και το  $\sigma_0^2$  δηλώνει την αβεβαιότητά μας.
- Προκύπτει ότι  $p(\mu/D^n) = N(\mu_n, \sigma_n^2)$ , όπου:
- Το  $\mu_n$  αντιπροσωπεύει την καλύτερη γνώση μας για το  $\mu$  μετά την παρατήρηση  $n$  δειγμάτων εκπαίδευσης και το  $\sigma_n^2$  μετράει την αβεβαιότητά μας.
- Επίσης,  $p(x/D^n) = N(\mu_n, \sigma^2 + \sigma_n^2)$ .

$$\mu_n = \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0,$$

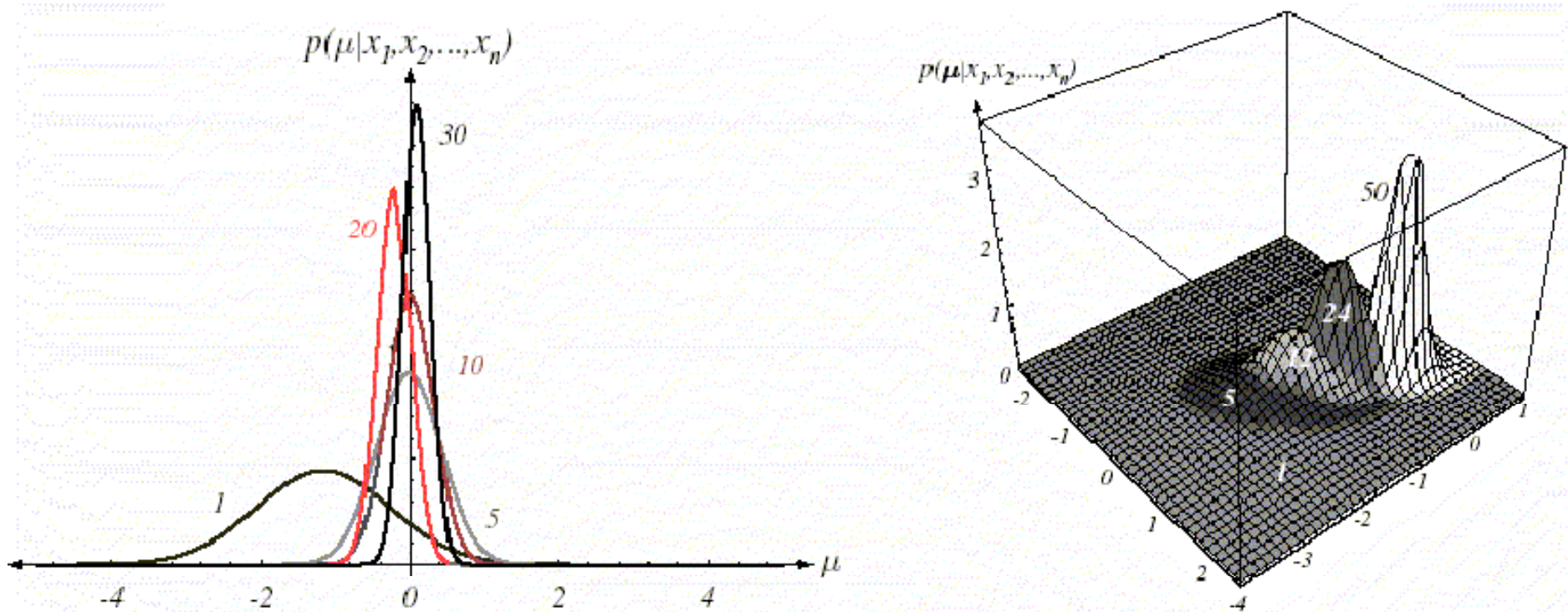
$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2},$$

$$\hat{\mu}_n = \frac{1}{n} \sum_{k=1}^n x_k$$

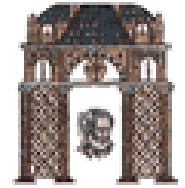
- ✓ Καθώς  $\sigma_0^2 \rightarrow \text{inf}$ , έχουμε  $\mu_n = \hat{\mu}_n$  για κάθε  $n$ , δηλαδή επιστρέφουμε στην εκτίμηση μέγιστης πιθανοφάνειας.
- ✓ Καθώς  $n \rightarrow \text{inf}$ , έχουμε  $\sigma_n^2 \approx \sigma^2/n$  δηλαδή για αρκετά μεγάλο αριθμό δειγμάτων εκπαίδευσης, η ακρίβεια της εκτίμησης του  $\mu$  δεν εξαρτάται από την αβεβαιότητα της εκ των προτέρων γνώσης μας,  $\sigma_0^2$ .



# Μπεϋζιανή εκμάθηση–Κανονική Κατανομή



Καθώς  $n \rightarrow \infty$ , η εκ των υστέρων σ.π.π.  $p(\mu|D^n)$  γίνεται όλο και περισσότερο συγκεντρωμένη γύρω από το μέσο της. Το φαινόμενο αυτό ονομάζεται Μπεϋζιανή εκμάθηση (Bayesian learning).



# Παράδειγμα : Εκτίμηση Παραμέτρων

Πετάμε ένα νόμισμα  $N$  φορές και παρατηρούμε  $k$  φορές κεφάλι. Ποια είναι η πιθανότητα  $\theta$  να φέρουμε κεφάλι;

## Μεγίστη Πιθανοφάνεια

Για να πάρουμε την συγκεκριμένη σειρά δεδομένων με κάποιο  $\theta$  έχουμε πιθανότητα  $P(D^N | \theta) = \theta^k (1-\theta)^{N-k}$   
 και συνεπώς  $P(D^N | \theta) = \theta^k (1-\theta)^{N-k}$ ,  $l(\theta) = \ln(P(D^N | \theta)) = k \ln(\theta) + (N-k) \ln(1-\theta)$ ,  $\frac{\partial l(\theta)}{\partial \theta} = \frac{k}{\theta} - \frac{N-k}{1-\theta} = 0$

Άρα  $\hat{\theta} = \frac{k}{N}$

## Εκτιμητής Bayes (Αναδρομική Εκτίμηση)

$$p(\theta | D^n) = \frac{p(x_n | \theta) p(\theta | D^{n-1})}{\int p(x_n | \theta) p(\theta | D^{n-1}) d\theta}, \quad p(\theta | D^0) \text{ γνωστό (apriori)}$$

ομως  $p(x_n | \theta) = \begin{cases} \theta & \text{εάν κεφάλι} \\ (1-\theta) & \text{εάν γράμματα} \end{cases}$ , και συνεπώς έχουμε

$$p(\theta | D^N) = \frac{\theta^k (1-\theta)^{N-k} p(\theta | D^0)}{\int \theta^k (1-\theta)^{N-k} p(\theta | D^0) d\theta}$$

Εκ των προτέρων (a-priori) Πιθανότητα

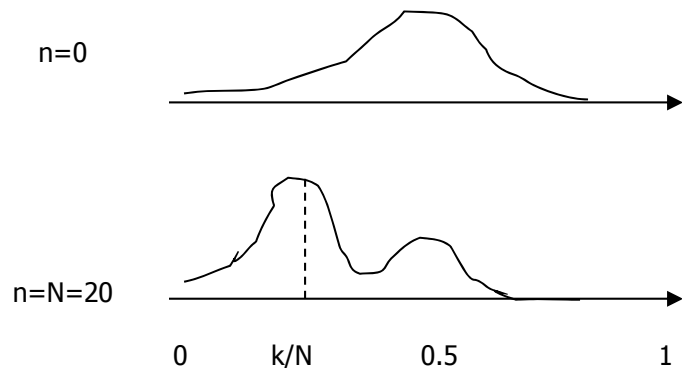
$$p(\theta | D^0) = 1 \text{ για } 0 \leq \theta \leq 1 \text{ ή}$$

$$p(\theta | D^0) = A e^{-\frac{(\theta-0.5)^2}{0.08}} (u(\theta) - u(\theta-1))$$

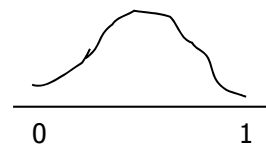
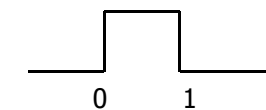
Πλήρης άγνοια, ομοιόμορφη κατανομή στο  $[0,1]$

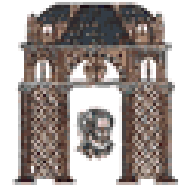
Ή

Αναμένεται να είναι δίκαιο νόμισμα, Κανονική στο  $[0,1]$



$p(\theta | D^{20})$  για  $K=4, N=20$





## Παράδειγμα : Εκτίμηση Παραμέτρων

Έχουμε δεδομένα από μία ομοιόμορφη κατανομή  $U(0, \theta)$ , δηλαδή,  $p(x|\theta) = U(0, \theta) = 1/\theta$  για  $0 \leq x \leq \theta$ . Τα δεδομένα που έχουμε είναι  $D = \{4, 7, 2, 8\}$  και  $0 \leq \theta \leq 10$  θέλουμε να εκτιμήσουμε πρώτα το  $\theta$  και μετά το  $p(x|\theta)$ . Εκτίμηση του  $\theta$

### Μεγίστη Πιθανοφάνεια

Για να πάρουμε την συγκεκριμένη σειρά δεδομένων με κάποιο  $\theta$  έχουμε πιθανότητα  $p(D^n | \theta) = \prod_{k=1}^n p(x_k | \theta) = (1/\theta)^n$   $8 \leq \theta \leq 10$  και συνεπώς  $\hat{\theta} = \max\{D_n\} = 8$

### Εκτιμητής Bayes (Αναδρομική Εκτίμηση)

$$p(\theta | D^n) = \frac{p(x_n | \theta) p(\theta | D^{n-1})}{\int_0^{\theta} p(x_n | \theta) p(\theta | D^{n-1}) d\theta}, \quad p(\theta | D^0) = p(\theta) = 0.1, \quad 0 \leq \theta \leq 10$$

$$p(x | \theta) \approx u(0, \theta) = \begin{cases} 1/\theta & 0 \leq x \leq \theta \\ 0 & \text{άλλού} \end{cases}$$

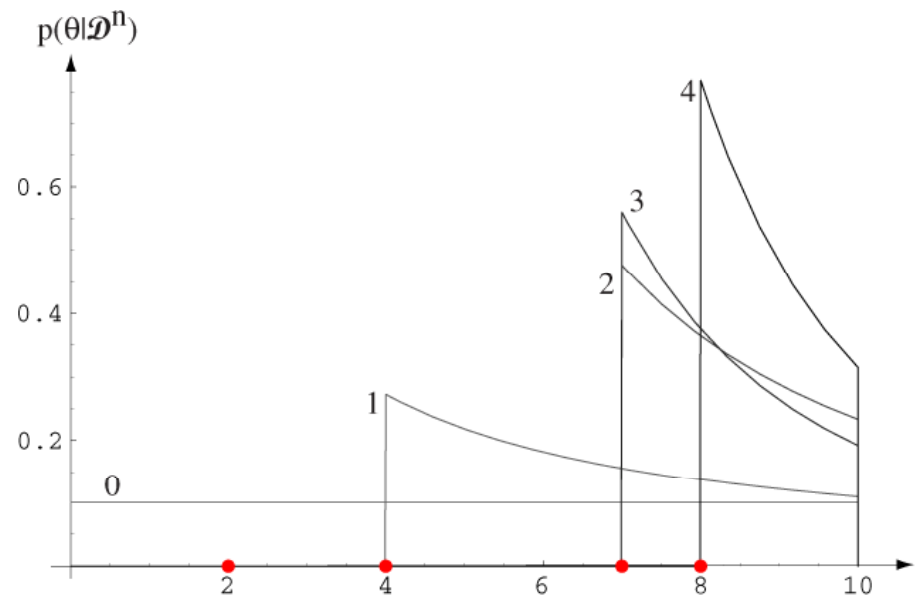
$$p(\theta | D^1) \propto p(x | \theta) p(\theta | D^0) = \begin{cases} 1/\theta & \text{για } 4 \leq \theta \leq 10 \\ 0 & \text{άλλού} \end{cases}$$

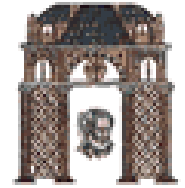
$$p(\theta | D^2) \propto p(x | \theta) p(\theta | D^1) = \begin{cases} 1/\theta^2 & \text{για } 7 \leq \theta \leq 10 \\ 0 & \text{άλλού} \end{cases}$$

$$p(\theta | D^3) \propto p(x | \theta) p(\theta | D^2) = \begin{cases} 1/\theta^3 & \text{για } 7 \leq \theta \leq 10 \\ 0 & \text{άλλού} \end{cases}$$

$$p(\theta | D^4) \propto p(x | \theta) p(\theta | D^3) = \begin{cases} 1/\theta^4 & \text{για } 8 \leq \theta \leq 10 \\ 0 & \text{άλλού} \end{cases}$$

$$p(\theta | D^n) \propto p(x | \theta) p(\theta | D^{n-1}) = \begin{cases} 1/\theta^n & \text{για } \max[D^n] \leq \theta \leq 10 \\ 0 & \text{άλλου} \end{cases}$$

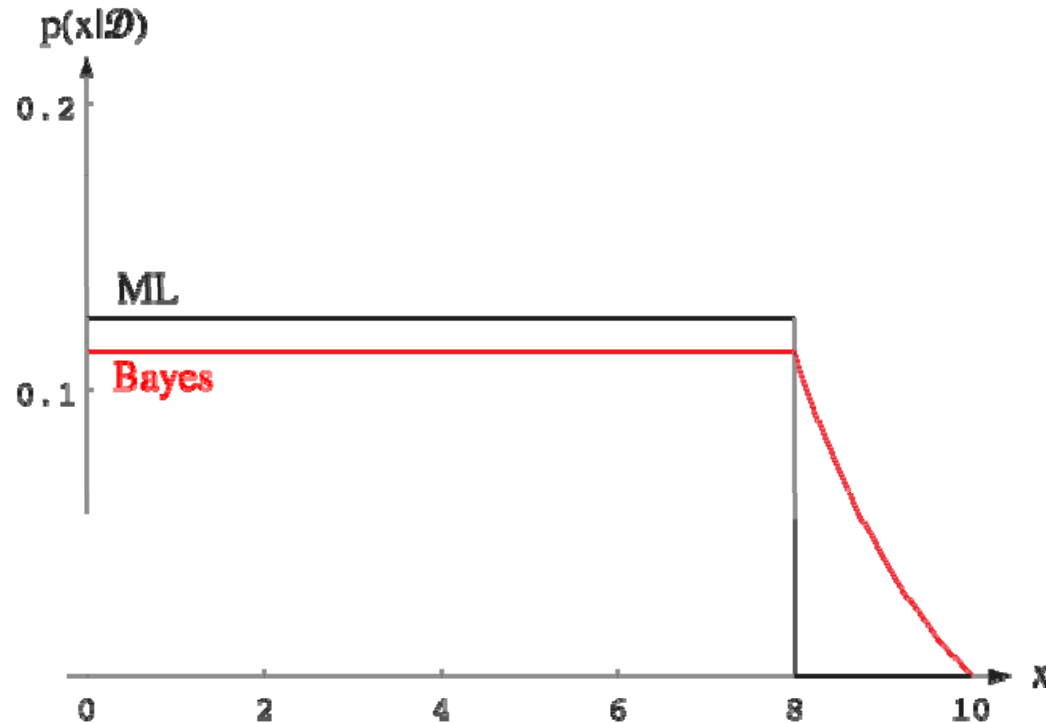




# Παράδειγμα 1 (συν.): Αναδρομική Εκτίμηση Παραμέτρων

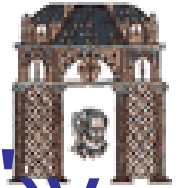
Άρα έχουμε με τα 4 σημεία ότι ο εκτιμητής του  $p(x|D)$  της μέγιστης πιθανοφάνειας είναι **η τιμή  $\theta=8$** ,  
 άρα  $p(x|D) = U(0,8)$ ,

ενώ ο εκτιμητής Bayes (κόκκινη γραμμή) του  $p(x|D)$  είναι  $p(x|D) = \int p(x|\theta)p(\theta|D)d\theta$



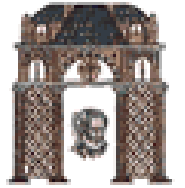
Διότι για  $0 \leq x \leq 8$ ,  $p(x|\theta) \sim 1/8$ , για  $8 \leq x \leq 10$ ,  $p(x|\theta) \sim 1/\theta$ .





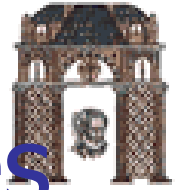
# Διαφορές ML και Bayes εκτιμητών

- Υπολογιστική πολυπλοκότητα
- ML πιο κατανοητό αποτέλεσμα
- Bayes δεν οδηγεί απαραίτητα στο παραμετρικό μοντέλο που υποθέσαμε (δες τελευταίο παράδειγμα)
- Bayes χρησιμοποιεί καλύτερα τα δεδομένα εκπαίδευσης
- Bayes χρειάζεται και την εκ των προτέρων κατανομή (συνήθως δεν είναι γνωστή)



# Curse of Dimensionality

Εάν για την **καλή** εκτίμηση των παραμέτρων μιας μονοδιάστατης pdf χρειάζονται  **$N$**  δείγματα, για την εκτίμηση της pdf σε  **$d$**  διαστάσεις χρειάζονται  **$N^d$**  δείγματα



# Απλοποιημένος Εκτιμητής Bayes

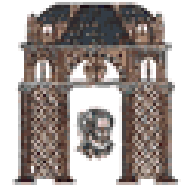
Στην περίπτωση αυτή υποθέτουμε ότι τα  $d$  χαρακτηριστικά ενός προτύπου  $\mathbf{x}=(x_1,x_2,\dots,x_d)$  είναι στατιστικά ανεξάρτητα

Άρα

$$p(\mathbf{x} | \omega_i) = \prod_{n=1}^d p(x_n | \omega_i)$$

Η τεχνική αυτή εφαρμόζεται όταν το  $d$  είναι μεγάλο (curse of dimensionality) και ο αριθμός  $N$  των δεδομένων εκπαίδευσης είναι σχετικά μικρός.

Συνήθως υποθέτουμε ότι τα pdf των  $x_1, x_2, \dots, x_d$  είναι κανονικές κατανομές και προσπαθούμε να εκτιμήσουμε την μέση τιμή και την διασπορά τους ( $\mu_n, \sigma_n^2$ )



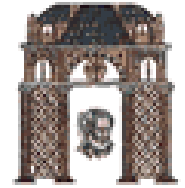
# Παράδειγμα

Έστω δεδομένα που δημιουργούνται από την παρακάτω pdf

$$p(\mathbf{x} | \omega_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_i|^{\frac{d}{2}}} \exp\left(-(\mathbf{x} - \mu_i)^T \Sigma_i^{-1} (\mathbf{x} - \mu_i)\right)$$

Με  $i=2$ ,  $d=5$  και  $\mu_1 = [0 \ 0 \ 0 \ 0 \ 0]^T$ ,  $\mu_2 = [1 \ 1 \ 1 \ 1 \ 1]^T$

$$\Sigma_1 = \begin{bmatrix} 0,8 & 0,2 & 0,1 & 0,05 & 0,01 \\ 0,2 & 0,7 & 0,1 & 0,03 & 0,02 \\ 0,1 & 0,1 & 0,8 & 0,02 & 0,01 \\ 0,05 & 0,03 & 0,02 & 0,9 & 0,01 \\ 0,01 & 0,02 & 0,01 & 0,01 & 0,8 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 0,9 & 0,1 & 0,05 & 0,02 & 0,01 \\ 0,1 & 0,8 & 0,1 & 0,02 & 0,02 \\ 0,05 & 0,1 & 0,7 & 0,02 & 0,01 \\ 0,02 & 0,02 & 0,02 & 0,6 & 0,02 \\ 0,01 & 0,02 & 0,01 & 0,02 & 0,7 \end{bmatrix}$$



# MIXTURE MODELS

Μία τυχαία pdf μπορεί να παρασταθεί σαν άθροισμα γνωστών pdf.

$$p(x) = \sum_{j=1}^J P_j p(x | j) \text{ όπου}$$
$$\sum_{j=1}^J P_j = 1, \quad \int_{x=-\infty}^{\infty} p(x|j) dx = 1$$

Μία συνήθης επιλογή για τις  $p(x/j)$  είναι η κανονική  $N(\mu_j, \Sigma_j)$

Άρα δεδομένου ενός δείγματος τιμών προσπαθούμε να εκτιμήσουμε τους συντελεστές  $P_j$  και τις παραμέτρους των  $p(x/j)$

Ένας δημοφιλής αλγόριθμος που υπολογίζει τις παραμέτρους μέσα από μία επαναληπτική διαδικασία είναι ο EM (δες Θεοδωρίδη, 4-η έκδοση, κεφ, 2.5.5)