

Μάθηση Κατανομών Πιθανότητας και Ομαδοποίηση

Κώστας Διαμαντάρας
Τμήμα Πληροφορικής
ΤΕΙ Θεσσαλονίκης

Μάθηση κατανομής πιθανότητας

- Σε όλη την ανάλυση μέχρι τώρα έγινε σιωπηρά η παραδοχή ότι γνωρίζουμε την κατανομή πιθανότητας $P(\mathbf{x}|C_i)$ για κάθε κλάση C_i . Πρακτικά όμως, συχνά, αυτό δεν ισχύει.
- Συνήθως διαθέτουμε μόνο K δείγματα $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K$ από τα οποία πρέπει να «εκτιμήσουμε» (βλ. να «μάθουμε») τη συνάρτηση $P(\mathbf{x}|C_i)$. Αντιμετωπίζουμε, δηλαδή, το γενικότερο πρόβλημα της εκτίμησης μιας κατανομής πιθανότητας.
- Αν το \mathbf{x} παίρνει μόνο διακριτές τιμές τότε η εκτίμησή μας είναι σχετικά εύκολη:

$$P(\mathbf{x} = \mathbf{a}|C_i) = \frac{K_a}{K}$$

- K_a = το πλήθος των δειγμάτων όπου $\mathbf{x} = \mathbf{a}$,
- K = το συνολικό πλήθος των δειγμάτων
- Μεγαλύτερο ενδιαφέρον έχει το πρόβλημα όταν το \mathbf{x} μπορεί να πάρει συνεχείς τιμές.

Η μέθοδος του ιστογράμματος


- Ας δούμε κατ' αρχήν την απλή περίπτωση όπου επιθυμούμε να εκτιμήσουμε την κατανομή μιας απλής μεταβλητής x (όχι διάνυσμα \mathbf{x}).
- **Μέθοδος του ιστογράμματος (histogram)**: Χωρίζουμε το διάστημα $[-M, M]$ όπου παίρνει τιμές το x σε B ίσα τμήματα που καλούνται κελιά. Το μέγεθος του κάθε κελιού είναι $V=2M/B$. Μετράμε πόσα δείγματα K_i πέφτουν σε κάθε κελί.
- Η ακολουθία K_1, \dots, K_B λέγεται *ιστόγραμμα*. Η εκτίμησή μας για την πυκνότητα πιθανότητας μέσα στο κελί i είναι σταθερή και ίση με

$$\hat{p}_i = \frac{K_i}{KV}$$

- Πλεονεκτήματα της μεθόδου:
 - Απλότητα
 - Δεν απαιτείται να υποθέσουμε ότι η κατανομή ανήκει σε κάποια οικογένεια κατανομών (πχ. τις Γκαουσιανές, ή άλλες)

Η μέθοδος του ιστογράμματος

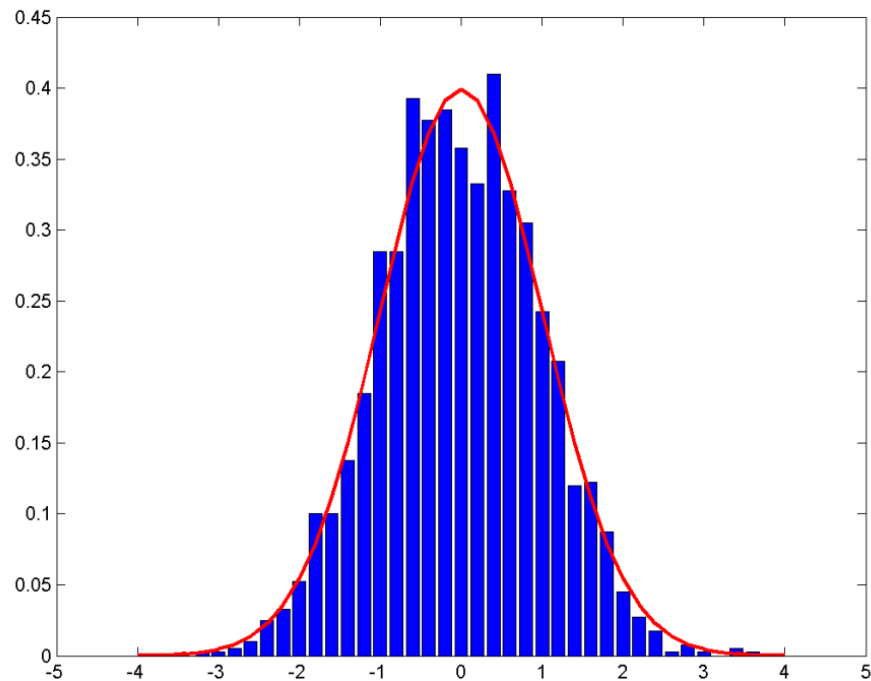
- **Μειονεκτήματα**

- Μη συνεχείς τιμές. Όχι ομαλή μετάβαση από κελί σε κελί
- Κακή ποιότητα εκτίμησης
- Εξάρτηση από το πλάτος των κελιών. Δες παραδείγματα 

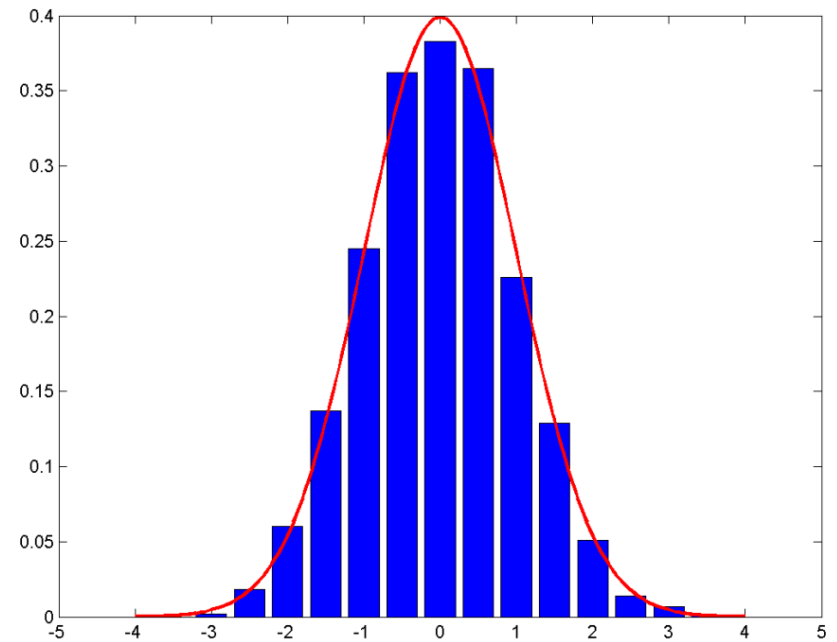
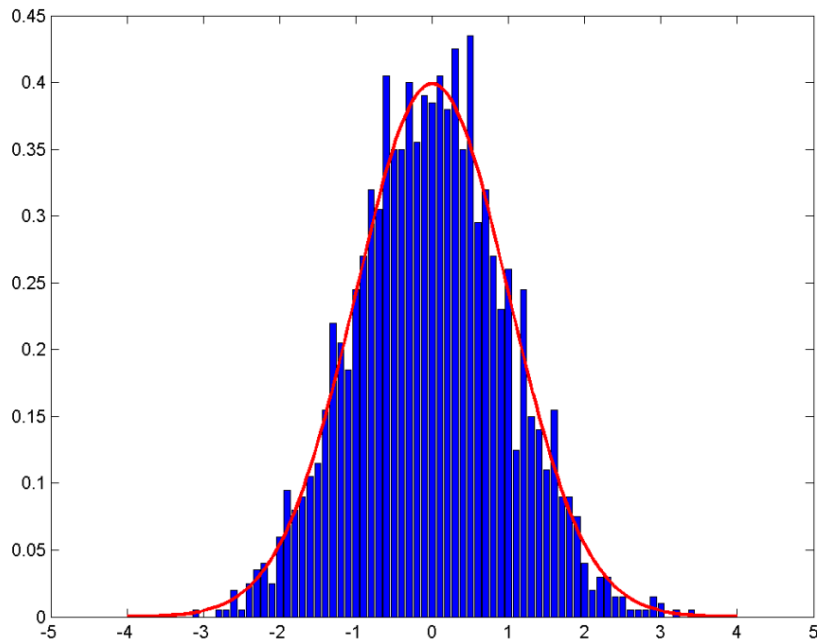
Ιστόγραμμα κανονικής κατανομής από $K=2000$ δείγματα.

 = ιστόγραμμα

 = κανονική κατανομή



Η μέθοδος του ιστογράμματος



Η μέθοδος των παραθύρων

- Συλλέγουμε K δείγματα x_1, x_2, \dots, x_K
- Για κάθε δείγμα x_i τοποθετούμε μια συνάρτηση που καλείται «πανάθυρο», πχ. την Γκαουσιανή, κεντραρισμένη στο σημείο αυτό και με πλάτος h (αυθαίρετο), δηλαδή

$$w_i(x) = \frac{1}{h\sqrt{2\pi}} \exp\left\{-\frac{(x - x_i)^2}{2h^2}\right\}$$

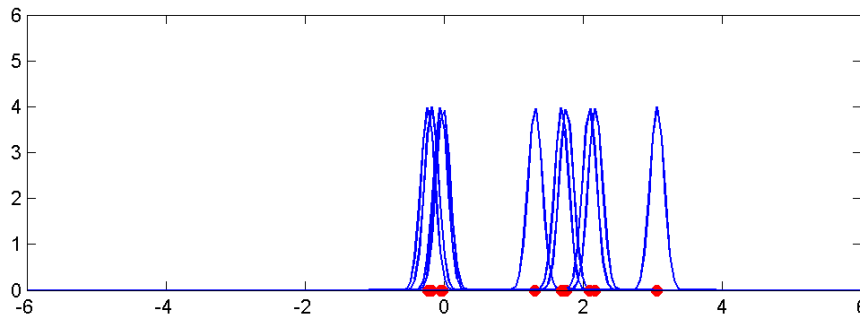
- Συνολικά η εκτίμησή μας για την κατανομή πιθανότητας είναι

$$\hat{p}(x) = \frac{1}{K} \sum_{i=1}^K w_i(x) = \frac{1}{Kh\sqrt{2\pi}} \sum_{i=1}^K \exp\left\{-\frac{(x - x_i)^2}{2h^2}\right\}$$


- Αν τα δείγματα είναι διανύσματα $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K$ τότε απλώς

$$w_i(\mathbf{x}) = \frac{1}{h\sqrt{2\pi}} \exp\left\{-\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{2h^2}\right\}$$

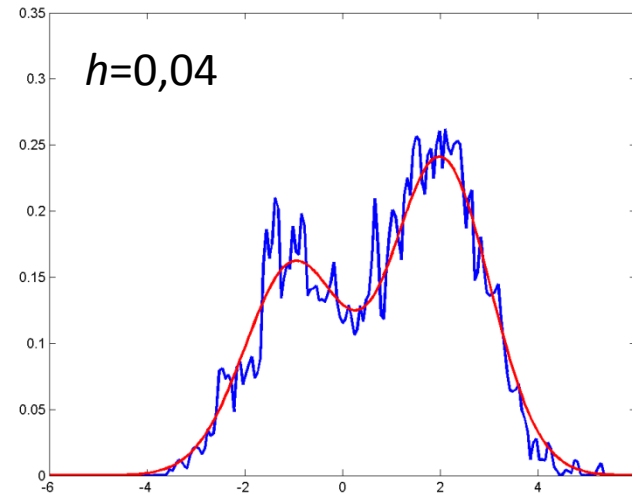
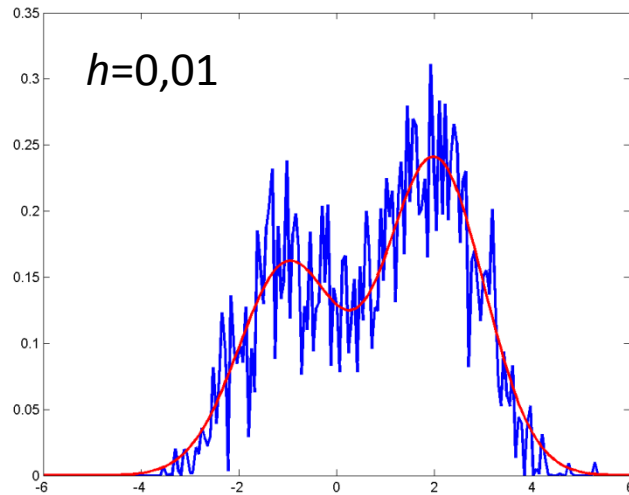
Η μέθοδος των παραθύρων



Για κάθε δείγμα $\bullet x(k)$ βάζουμε ένα Γκαουσιανό παράθυρο \wedge κεντραρισμένο στο δείγμα αυτό.

Στο τέλος αφού βάλουμε όλα τα παράθυρα για όλα τα δείγματα, η εκτίμησή μας είναι ο μέσος όρος, δηλαδή αθροίζουμε όλα τα παράθυρα και διαιρούμε δια K (το πλήθος των δειγμάτων και των παραθύρων επίσης). Το πλάτος των παραθύρων h είναι αυθαίρετο και το ορίζει ο χρήστης. Παίζει σημαντικό ρόλο στην ποιότητα της εκτίμησης και την ομαλότητα της καμπύλης εκτίμησης. Δες παραδείγματα 

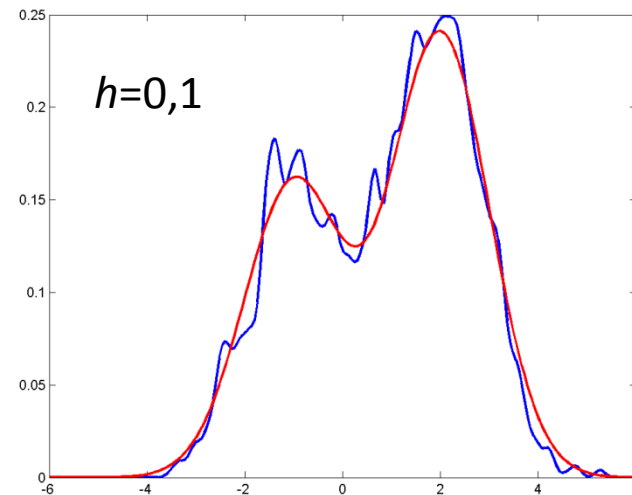
Η μέθοδος των παραθύρων



- = Εκτίμηση
- = Πραγματική κατανομή

Όσο μεγαλώνει το h τόσο πιο ομαλή είναι η καμπύλη εκτίμησης.

Μειονέκτημα: Δυσκολία εντοπισμού του βέλτιστου h .



Αν η κατανομή είναι Γκαουσιανή...

- Αν γνωρίζουμε εκ των προτέρων ότι η κατανομή είναι Γκαουσιανή τότε τα πράγματα είναι αρκετά απλά: μπορούμε να εκτιμήσουμε όλες τις παραμέτρους της κατανομής, δηλαδή τη μέση τιμή και τη διασπορά, μέσα από τα δείγματα.
- Αν x : αριθμός:

- Κατανομή:
$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x - \mu)^2}{2\sigma^2}\right\}$$

- Εκτίμηση μέσης τιμής:
$$\hat{\mu} = \frac{1}{K} \sum_{i=1}^K x(k)$$

- Εκτίμηση διασποράς:
$$\hat{\sigma}^2 = \frac{1}{K} \sum_{i=1}^K (x(k) - \hat{\mu})^2$$

Αν η κατανομή είναι Γκαουσιανή...

- Αν \mathbf{x} : διάνυσμα διάστασης n :

- Κατανομή:
$$p(\mathbf{x}) = \frac{1}{(2\pi \cdot \det(\Sigma))} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$$

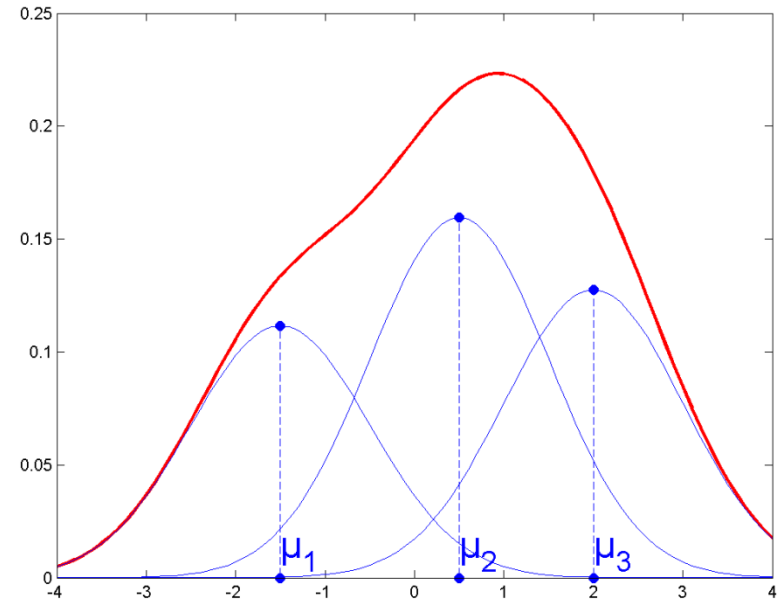
- Εκτίμηση μέσης τιμής:
$$\hat{\boldsymbol{\mu}} = \frac{1}{K} \sum_{i=1}^K \mathbf{x}(k)$$

- Εκτίμηση πίνακα διασποράς:
$$\hat{\boldsymbol{\Sigma}} = \frac{1}{K} \sum_{i=1}^K (\mathbf{x}(k) - \hat{\boldsymbol{\mu}})(\mathbf{x}(k) - \hat{\boldsymbol{\mu}})^T$$

- **Βασικό μειονέκτημα:** συνήθως δεν γνωρίζουμε εκ των προτέρων το είδος της κατανομής. Αν το γνωρίζαμε, τότε υπάρχουν αντίστοιχες φόρμουλες εκτίμησης των παραμέτρων για όλες τις γνωστές κατανομές (Gauss, Poisson, Bernoulli, Εκθετική, Rayleigh, Βήτα, Γάμμα, Maxwell, Δυωνυμική, κλπ)

Αν η κατανομή είναι μίγμα Γκαουσιανών...

- 1 Διάσταση:
- Η κατανομή είναι το άθροισμα m Γκαουσιανών κατανομών.
- Κάθε κατανομή έχει το δικό της κέντρο μ_i και τη δική της διασπορά σ_i^2 .
- Κάθε κατανομή από τις m περιγράφει μια ομάδα (cluster) προτύπων G_i



$$p(x) = \sum_{i=1}^m P(G_i) \underbrace{\frac{1}{\sigma_i \sqrt{2\pi}} \exp\left\{-\frac{(x - \mu_i)^2}{2\sigma_i^2}\right\}}_{P(x|G_i)}$$

Ο αλγόριθμος Expectation Maximization (EM) για 1 διάσταση

Αρχικοποίηση: $\hat{P}(G_1) = \dots = \hat{P}(G_m) = \frac{1}{m}$, $\hat{\mu}_i = \text{τυχαία}$, $\hat{\sigma}_1 = \dots = \hat{\sigma}_m = 1$ (πχ.)

Για κάθε επανάληψη $e=1, \dots, \text{MAX}$ {

Για κάθε πρότυπο $k=1, \dots, K$ {

Για κάθε ομάδα $i=1, \dots, m$ {

}

}

Για κάθε ομάδα $i=1, \dots, m$ {

}

}

$$P(G_i|x_k) = \frac{\frac{1}{\hat{\sigma}_i \sqrt{2\pi}} \exp\left\{-\frac{(x_k - \hat{\mu}_i)^2}{2\hat{\sigma}_i^2}\right\} \hat{P}(G_i)}{\sum_{j=1}^m \frac{1}{\hat{\sigma}_j \sqrt{2\pi}} \exp\left\{-\frac{(x_k - \hat{\mu}_j)^2}{2\hat{\sigma}_j^2}\right\} \hat{P}(G_j)}$$

$$\hat{\mu}_i = \frac{\sum_{k=1}^K P(G_i|x_k) x_k}{\sum_{k=1}^K P(G_i|x_k)}$$

$$\hat{\sigma}_i^2 = \frac{1}{2} \frac{\sum_{k=1}^K P(G_i|x_k) (x_k - \hat{\mu}_i)^2}{\sum_{k=1}^K P(G_i|x_k)}$$

$$\hat{P}(G_i) = \frac{1}{K} \sum_{k=1}^K P(G_i|x_k)$$

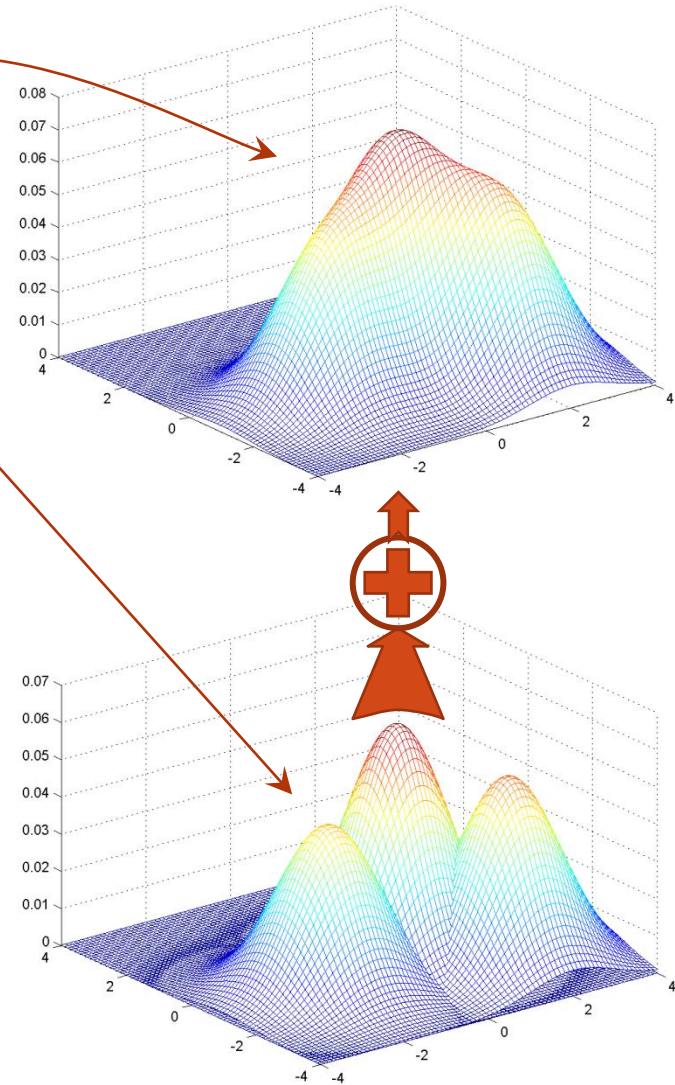
Ο αλγόριθμος EM

- Επαναληπτικός αλγόριθμος. Εκτιμά (μαθαίνει) όλο και καλύτερα τις παραμέτρους $P(G_i)$, μ_i , και σ_i^2 όσο προχωρούν οι επαναλήψεις.
- Δυστυχώς δεν υπάρχει εγγυημένη σύγκλιση στις σωστές παραμέτρους. Μερικές φορές ο αλγόριθμος κάνει λάθος.
- Προϋποθέτει ότι γνωρίζουμε το πλήθος m των Ομάδων.
- Η αρχικοποίηση παίζει σημαντικό ρόλο

Αν η κατανομή είναι μίγμα Γκαουσιανών...

- 2 Διαστάσεις:
- Η κατανομή είναι το άθροισμα m Γκαουσιανών κατανομών.
- Κάθε κατανομή έχει το δικό της κέντρο μ_i και πίνακα διασποράς $\Sigma = \sigma_i^2 \mathbf{I}$.
- Κάθε κατανομή από τις m περιγράφει μια ομάδα προτύπων G_i (cluster)

$$p(\mathbf{x}) = \sum_{i=1}^m P(G_i) \underbrace{\frac{1}{2\pi\sigma_i^2} \exp\left\{-\frac{\|\mathbf{x} - \mu_i\|^2}{2\sigma_i^2}\right\}}_{P(\mathbf{x}|G_i)}$$



Ο αλγόριθμος EM για 2 διαστάσεις

Αρχικοποίηση: $\hat{P}(G_1) = \dots = \hat{P}(G_m) = \frac{1}{m}$, $\hat{\boldsymbol{\mu}}_i = \text{τυχαία}$, $\hat{\sigma}_1 = \dots = \hat{\sigma}_m = 1$ (πχ.)

Για κάθε επανάληψη $e=1, \dots, \text{MAX}$ {

Για κάθε πρότυπο $k=1, \dots, K$ {

Για κάθε ομάδα $i=1, \dots, m$ {

}

}

Για κάθε ομάδα $i=1, \dots, m$ {

}

}

$$P(G_i | \mathbf{x}_k) = \frac{\frac{1}{2\pi\hat{\sigma}_i^2} \exp\left\{-\frac{\|\mathbf{x}_k - \hat{\boldsymbol{\mu}}_i\|^2}{2\hat{\sigma}_i^2}\right\} \hat{P}(G_i)}{\sum_{j=1}^m \frac{1}{2\pi\hat{\sigma}_j^2} \exp\left\{-\frac{\|\mathbf{x}_k - \hat{\boldsymbol{\mu}}_j\|^2}{2\hat{\sigma}_j^2}\right\} \hat{P}(G_j)}$$

$$\hat{\boldsymbol{\mu}}_i = \frac{\sum_{k=1}^K P(G_i | \mathbf{x}_k) \mathbf{x}_k}{\sum_{k=1}^K P(G_i | \mathbf{x}_k)}$$

$$\hat{\sigma}_i^2 = \frac{1}{2} \frac{\sum_{k=1}^K P(G_i | \mathbf{x}_k) \|\mathbf{x}_k - \hat{\boldsymbol{\mu}}_i\|^2}{\sum_{k=1}^K P(G_i | \mathbf{x}_k)}$$

$$\hat{P}(G_i) = \frac{1}{K} \sum_{k=1}^K P(G_i | \mathbf{x}_k)$$

Ο αλγόριθμος EM για 2 διαστάσεις

Αρχικοποίηση: $\hat{P}(G_1) = \dots = \hat{P}(G_m) = \frac{1}{m}$, $\hat{\boldsymbol{\mu}}_i = \text{τυχαία}$, $\hat{\sigma}_1 = \dots = \hat{\sigma}_m = \text{τυχαία}$

Για κάθε πρότυπο $k=1, \dots, K$ {

$$P(G_i | \mathbf{x}_k) = \frac{\frac{1}{2\pi\hat{\sigma}_i^2} \exp\left\{-\frac{\|\mathbf{x}_k - \hat{\boldsymbol{\mu}}_i\|^2}{2\hat{\sigma}_i^2}\right\} \hat{P}(G_i)}{\sum_{j=1}^m \frac{1}{2\pi\hat{\sigma}_j^2} \exp\left\{-\frac{\|\mathbf{x}_k - \hat{\boldsymbol{\mu}}_j\|^2}{2\hat{\sigma}_j^2}\right\} \hat{P}(G_j)}$$

}

Για κάθε ομάδα $i=1, \dots, m$ {

$$\hat{\boldsymbol{\mu}}_i = \frac{\sum_{k=1}^K P(G_i | \mathbf{x}_k) \mathbf{x}_k}{\sum_{k=1}^K P(G_i | \mathbf{x}_k)}$$

$$\hat{\sigma}_i^2 = \frac{1}{2} \frac{\sum_{k=1}^K P(G_i | \mathbf{x}_k) \|\mathbf{x}_k - \hat{\boldsymbol{\mu}}_i\|^2}{\sum_{k=1}^K P(G_i | \mathbf{x}_k)}$$

$$\hat{P}(G_i) = \frac{1}{K} \sum_{k=1}^K P(G_i | \mathbf{x}_k)$$

}

Άσκηση-3

- Υλοποιήστε τον αλγόριθμο EM για μια διάσταση σε MATLAB
- Χρησιμοποιήστε το ακόλουθο [<link>](#) σετ δεδομένων και τρέξτε τον αλγόριθμο
- Προαιρετικά: Υλοποιήστε τον αλγόριθμο EM για δύο διαστάσεις σε MATLAB
- Χρησιμοποιήστε το ακόλουθο [<link>](#) σετ δεδομένων και τρέξτε τον αλγόριθμο

Γιατί είναι σημαντικά τα μίγματα Γκαουσιανών;

- Αποδεικνύεται ότι με αθροίσματα Γκαουσιανών μπορούμε να προσεγγίσουμε οποιαδήποτε συνάρτηση κατανομής με όση ακρίβεια θέλουμε!
- Αρκεί να έχουμε πολλές Γκαουσιανές...
- Πόσες Γκαουσιανές ομάδες απαιτούνται ; Ποια είναι η τιμή του m ; Η απάντηση δεν είναι εύκολη (και βασικά άγνωστη!)
- Μειονεκτήματα:
 - Ο αλγόριθμος EM απαιτεί τη γνώση του m . Πάντως αν χρησιμοποιήσουμε αρκετά μεγάλο m έχουμε καλή πιθανότητα να προσεγγίσουμε οποιαδήποτε συνάρτηση κατανομής.
 - Ο αλγόριθμος EM μπορεί να μην βρεί την καλύτερη δυνατή λύση αλλά κάτι κοντά σε αυτή...
 - Θέμα τιμών αρχικοποίησης. Ποιες τιμές να δώσω;;;

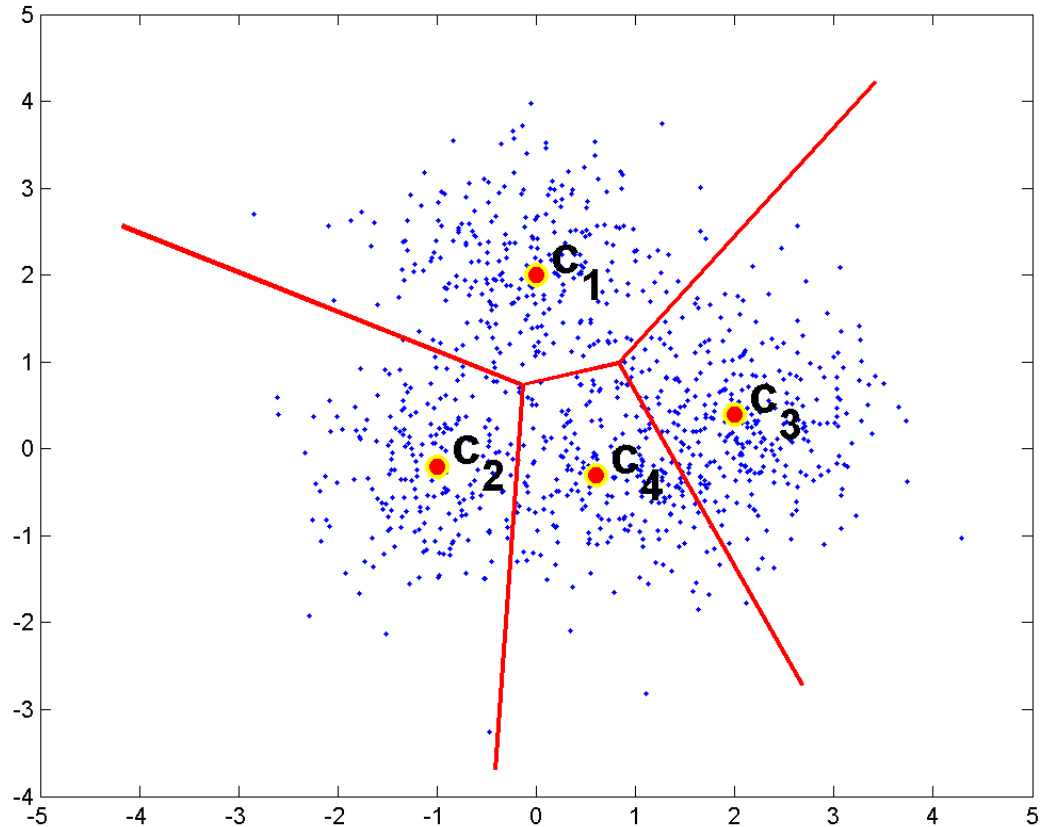
Ομαδοποίηση (Clustering)

- Στη γενική περίπτωση μια κατανομή αποτελείται από μια ή περισσότερες μικρότερες ομάδες κατανομών οι οποίες μπορεί να είναι Γκαουσιανές, αλλά μπορεί και όχι.
- Το πρόβλημα της **ομαδοποίησης (clustering)** αφορά στον εντοπισμό των ομάδων (κέντρα, διασπορές, κλπ) αλλά και στην απόφαση για κάθε πρότυπο σε ποια ομάδα ανήκει.
- Αν οι κατανομές των ομάδων είναι γνωστές, δηλαδή όλες οι ομάδες ακολουθούν, πχ. την κατανομή Poisson, τότε μπορούμε να εφαρμόσουμε τον αλγόριθμο EM (εφαρμόζεται για οποιαδήποτε κατανομή ομάδων, όχι μόνο για Γκαουσιανές ομάδες, αλλά δεν θα μπορούμε σε λεπτομέρειες...)
- Τι γίνεται όμως αν δεν γνωρίζουμε τις κατανομές των ομάδων;

Ο αλγόριθμος K -μέσων (K -means)

- Ο αλγόριθμος αυτός δεν υποθέτει κάποια συγκεκριμένη συνάρτηση κατανομής.
- Θεωρεί γνωστό το πλήθος m των ομάδων (clusters) .
- Θεωρεί ότι κάθε ομάδα G_i αντιπροσωπεύεται από ένα κεντρικό σημείο c_i (κέντρο)
 1. Κάθε πρότυπο x_k ανήκει στην ομάδα G_i της οποίας το κέντρο βρίσκεται πιο κοντά στο x_k
 2. Το κέντρο c_i είναι ο μέσος όρων των προτύπων που ανήκουν στην ομάδα G_i

Ο αλγόριθμος των K -μέσων



Ο αλγόριθμος K -μέσων

Δίνονται: Τα πρότυπα $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_P$, το πλήθος κέντρων K

Έξοδος: K κέντρα $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_K$

Αρχικοποίησε τα $\mathbf{c}_1, \dots, \mathbf{c}_K$ σε τυχαίες τιμές.

Επανάλαβε {

 Για κάθε πρότυπο $i = 1, \dots, P$ {

 Βρες το κοντινότερο κέντρο \mathbf{c}_j στο \mathbf{x}_i

$label(i) \leftarrow j$ }

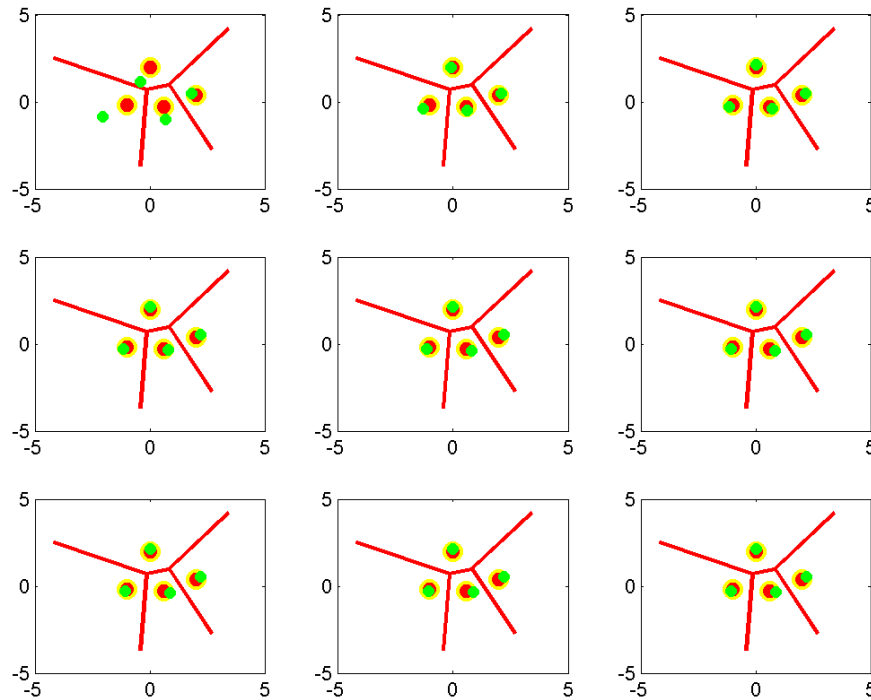
 Για κάθε κέντρο $k = 1, \dots, K$ {

$\mathbf{c}_k \leftarrow$ μέσος όρος των \mathbf{x}_i τα οποία ανήκουν στην κλάση k , δηλαδή για όσα $i: label(i) = k$ }

} Μέχρι να μην υπάρξει καμία αλλαγή στα κέντρα $\mathbf{c}_1, \dots, \mathbf{c}_K$

Σύγκλιση του αλγορίθμου

- Δεν υπάρχει εγγυημένη σύγκλιση αλλά συνήθως τα αποτελέσματα είναι καλά έως πολύ καλά.



Άσκηση-4

- Υλοποιήστε τον αλγόριθμο k-means σε MATLAB
- Χρησιμοποιήστε το ακόλουθ σεντ δεδομένων [<link>](#)