

22-12-2021

$d_j(\cdot), j=1, \dots, P$

$$D(x_i, x_{i'}) = \sum_{j=1}^P w_j d_j(x_{ij}, x_{i'j}), \quad i, i' \in \{1, \dots, N\}$$

$$\begin{aligned} \bar{D} &= \frac{1}{N} \sum_{i=1}^N \frac{1}{N} \sum_{i'=1}^N \sum_{j=1}^P w_j d_j(x_i, x_{i'}) = \frac{1}{N^2} \sum_{i=1}^N \sum_{i'=1}^N \sum_{j=1}^P w_j d_j(x_i, x_{i'}) \\ &= \sum_{j=1}^P w_j \underbrace{\left[\frac{1}{N^2} \sum_{i, i'} d_j(x_i, x_{i'}) \right]}_{\bar{d}_j} = \sum_{j=1}^P w_j \bar{d}_j = \bar{D} \end{aligned}$$

Για να έχουν στις οι μεταβλητές παρόμοια επίδραση στην επίλυση της απομείωσης πρέπει $w_j \sim \frac{1}{d_j}$

$$X_1 \in [11, 12] \quad \bar{d}_1 < \bar{d}_2 \Rightarrow w_1 > w_2$$

$$X_2 \in [2, 20]$$

Για το \bar{d}_j , αν $d_j(x_i, x_{i'}) = (x_{ij} - x_{i'j})^2$

$$\text{Τότε } \bar{d}_j = \frac{1}{N^2} \sum_{i=1}^N \sum_{i'=1}^N (x_{ij} - x_{i'j})^2 =$$

$$= \frac{1}{N^2} \sum_{i, i'} \left[(x_{ij} - \bar{x}_j) - (x_{i'j} - \bar{x}_j) \right]^2 =$$

$$= \frac{1}{N^2} \sum_i \sum_{i'} \left[(x_{ij} - \bar{x}_j)^2 + (x_{i'j} - \bar{x}_j)^2 - 2(x_{ij} - \bar{x}_j)(x_{i'j} - \bar{x}_j) \right]$$

$$= \frac{1}{N^2} \left[\underbrace{\sum_{i'} \sum_i (x_{ij} - \bar{x}_j)^2}_{N \sum_{i=1}^N (x_{ij} - \bar{x}_j)^2} + \underbrace{\sum_i \sum_{i'} (x_{i'j} - \bar{x}_j)^2}_{N \sum_{i'=1}^N (x_{i'j} - \bar{x}_j)^2} - 2 \sum_i (x_{ij} - \bar{x}_j) \cdot \sum_{i'} (x_{i'j} - \bar{x}_j) \right]$$

=0 =0

$$= \frac{2}{N} \sum_{i=1}^N (x_{ij} - \bar{x}_j)^2 = \frac{2}{N} (N-1) \text{Var}(x_j) \approx 2 \text{Var}(x_j)$$

$$w_i \sim \frac{1}{\text{Var}_j}$$

Παρατήρηση

μεταβλητών

Αν κάνουμε ενοποίηση των

$$\tilde{X}_j = \frac{X_j - \bar{X}_j}{\sqrt{\text{Var}_j}}$$

είναι ισοδύναμο με $w_j \sim \frac{1}{\text{Var}_j}$

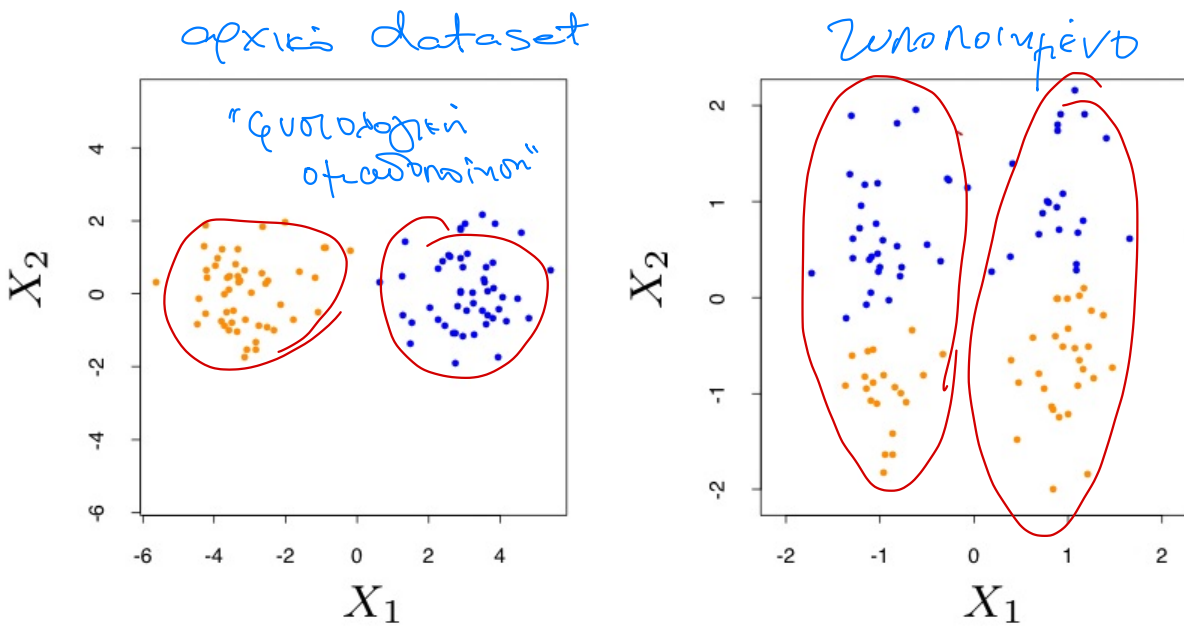


FIGURE 14.5. *Simulated data: on the left, K -means clustering (with $K=2$) has been applied to the raw data. The two colors indicate the cluster memberships. On the right, the features were first standardized before clustering. This is equivalent to using feature weights $1/[2 \cdot \text{var}(X_j)]$. The standardization has obscured the two well-separated groups. Note that each plot uses the same units in the horizontal and vertical axes.*

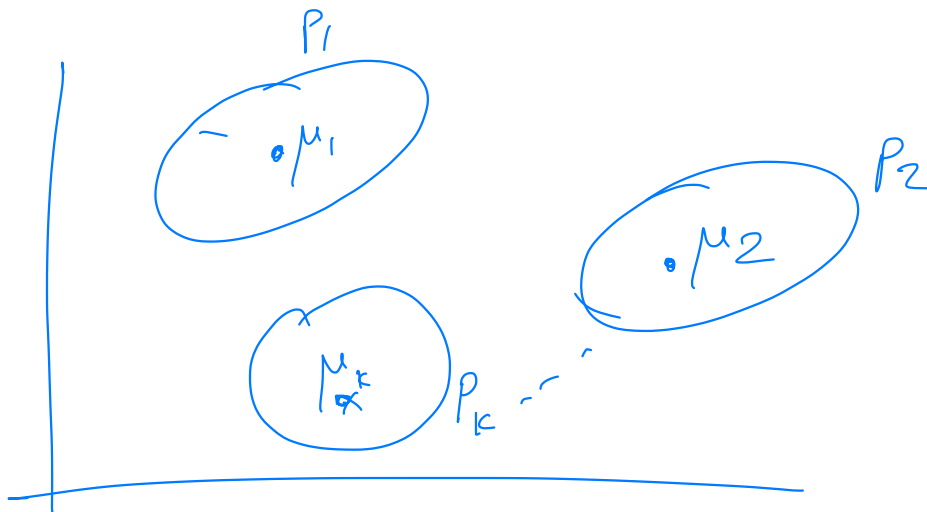
Clustering Algorithms

Χωρίζω το dataset σε ομάδες (συστάδες-clusters) έτσι ώστε η απόσταση μεταξύ παρατηρήσεων στην ίδια ομάδα να είναι μικρή ή σε διαφορετικές ομάδες μεγάλη

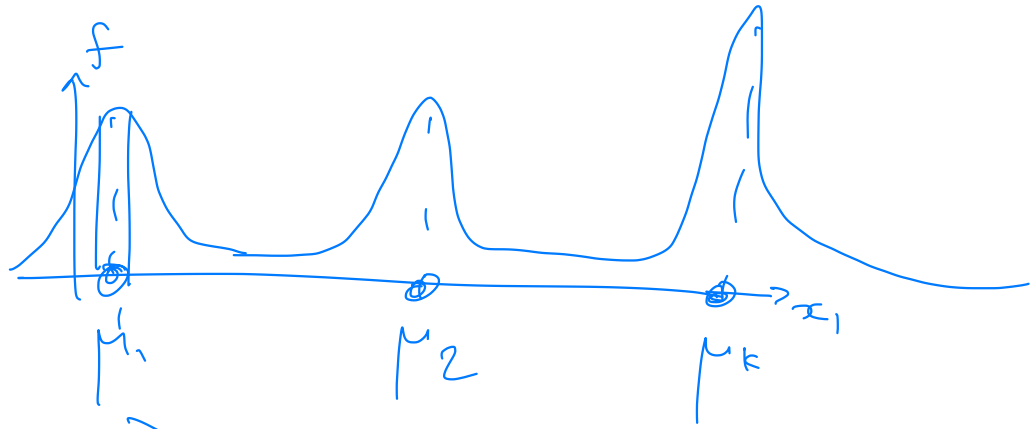
Συνδυαστικοί αλγόριθμοι Δεν κάνουν υπόθεση πιθανοθεωρητικού τύπου

Μειξωμοί Υπόθεση ότι οι μεταβλητές ακολουθούν δύο κοινές κατανομές (μειξωμοί ή κίβριες κατανομές)

$$\underline{X} = (X_1, \dots, X_p) \begin{matrix} \rightarrow F_1 & \text{με πιθανότητα } P_1 \\ \rightarrow F_2 & \text{" " } P_2 \\ \vdots & \vdots \\ \rightarrow F_k & \text{" " } P_k \end{matrix}$$



Αναζήτηση κορυφών
(mode-seeking)



(modes) ζοηικά πεδία εν f

αναζήτηση modes εν από κοινού πεδία εν
χωρίς παρατηρητική υνδύση για το πεδίο πεδίου

Συνδυαστικοί αλγόριθμοι

Θεωρούμε $k =$ αριθμός ομάδων (εξ αρχής σταθερή)

Συνάρτηση κωδικοποίησης (encoder function)

$C(i) =$ ομάδα στην οποία κατατάσσεται η παραρ i
 $i=1, \dots, N$

$$C : \{1, \dots, N\} \rightarrow \{1, \dots, k\}$$

"Loss function" : κριτήριο

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{i: C(i)=k} \sum_{i': C(i')=k} \overbrace{d(x_i, x_{i'})}^{d_{ii'}}$$

απόσταση ανάμεσα ανά δύο των παραρ. που ανήκουν στην ομάδα k .

= within cluster point scatter

Θέλουμε το $W(C)$ μικρό

$B(C)$: between cluster scatter

$$= \frac{1}{2} \sum_{k=1}^K \sum_{i: C(i)=k} \sum_{i': C(i') \neq k} d(x_i, x_{i'})$$

$$T(C) = W(C) + B(C) = \frac{1}{2} \sum_{k=1}^K \sum_{i: C(i)=k} \left[\sum_{i': C(i')=k} d_{ii'} + \sum_{i': C(i') \neq k} d_{ii'} \right]$$
$$= \sum_{i=1}^N d_{ii'}$$

$$\Rightarrow T(C) = \left[T = \frac{1}{2} \sum_i \sum_{i'} d_{ii'} \right] = \frac{\text{average cost}}{!!}$$

$$\Rightarrow \underline{W(C) + B(C)} = T \quad \forall C$$

$$\min W(C) \iff \max B(C)$$

K-means clustering algorithm

$$[d_j(x_i, x_{i'}) = (x_{ij} - x_{i'j})^2]$$

Εναλλακτικός αλγόριθμος

$$d_{i,i'} = \|x_i - x_{i'}\|^2$$

Θέλουμε C , και τα κέντρα κάθε ομάδας m_1, \dots, m_K .

Μπορούμε να δείξουμε ότι για κάθε C

$$W(C) = \sum_{k=1}^K N_k \sum_{i: C(i)=k} \|x_i - \bar{x}_k\|^2$$

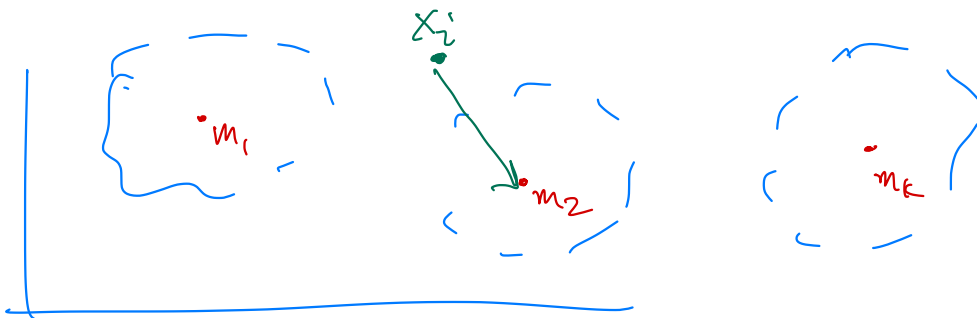
N_k : αρ. παραρ. πού $C(i)=k$

↑ centroids

① αυθαίρετο C : αυθαίρετη κατανομή σε ομάδες

② Δεδομ. του C : βρούμε τα m_1, \dots, m_K : $W(C)$ ελάχιστο

απόδειξη $m_k = \bar{x}_k, k=1, \dots, K$

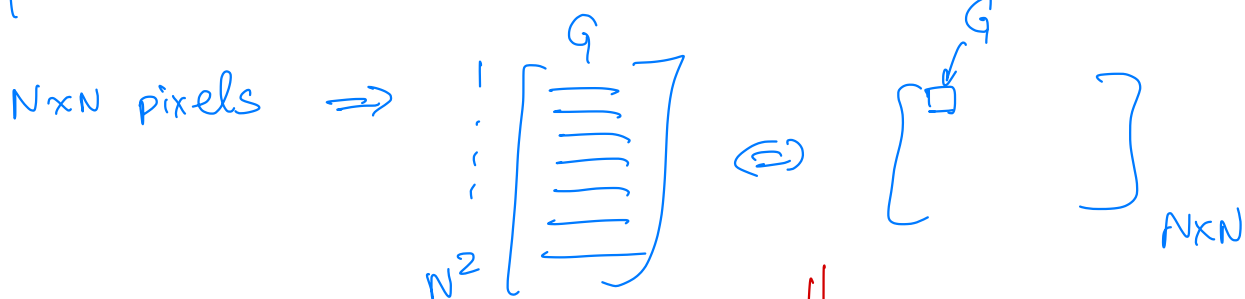
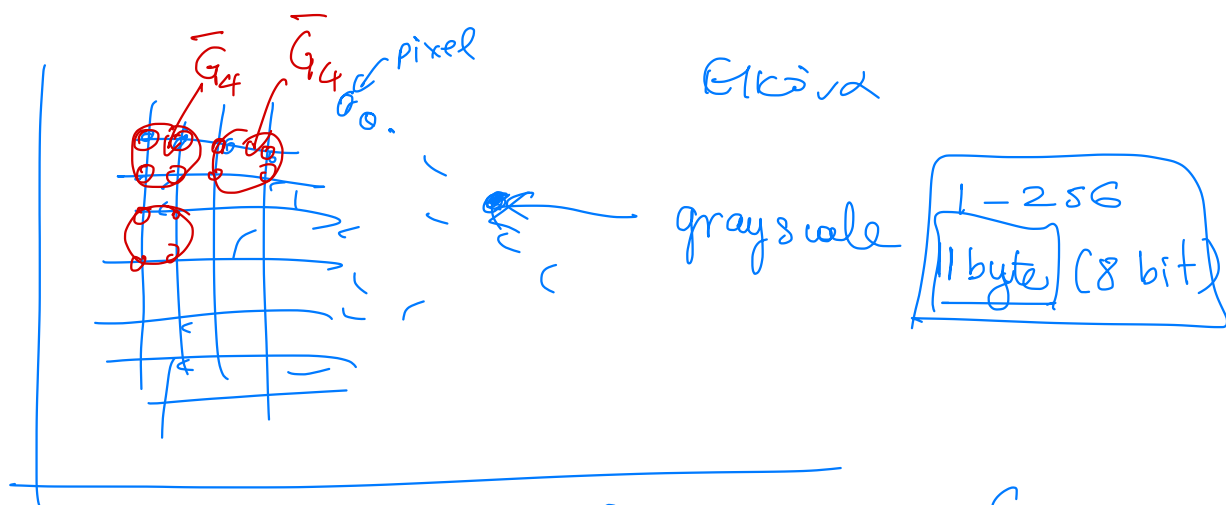


③ Επιστρέφουμε στο 1: Με δεδομένα m_1, \dots, m_K νέο C' : $\forall i=1, \dots, N: C'(i)=k$ αν

$$\|x_i - m_k\|^2 = \min_{k'} \|x_i - m_{k'}\|^2$$

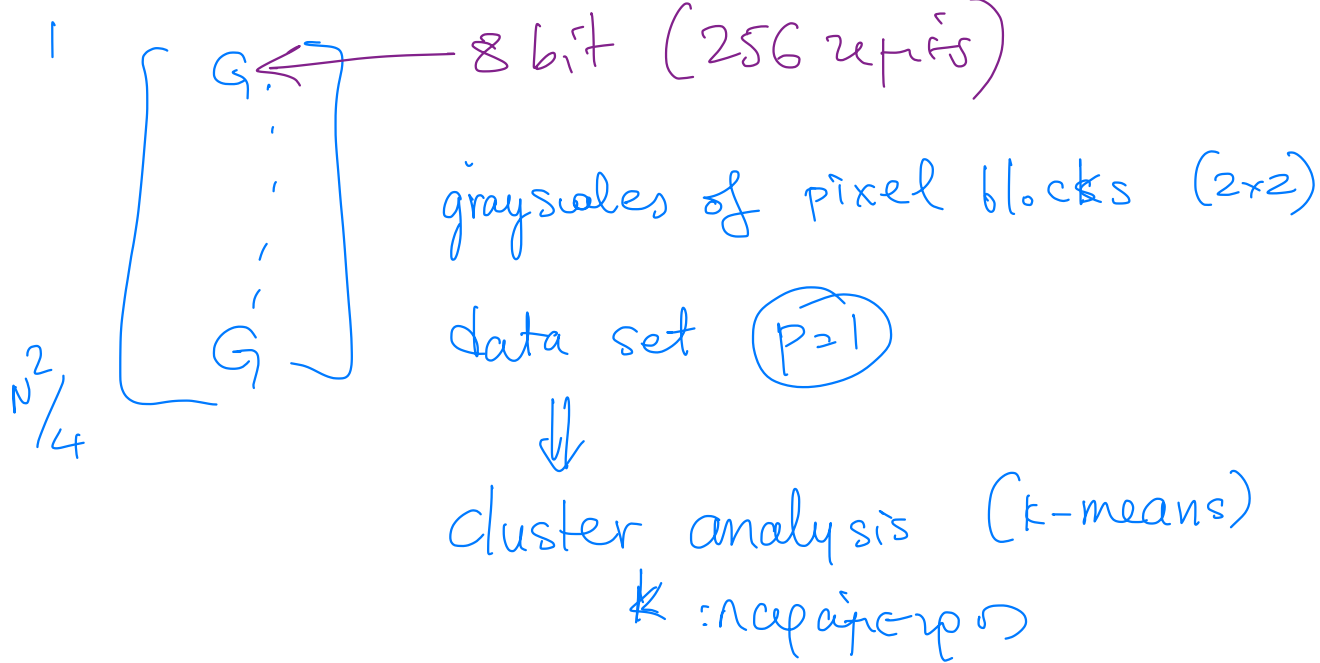
Εφαρμογή

Vector quantization για απεικόνιση γκρι.

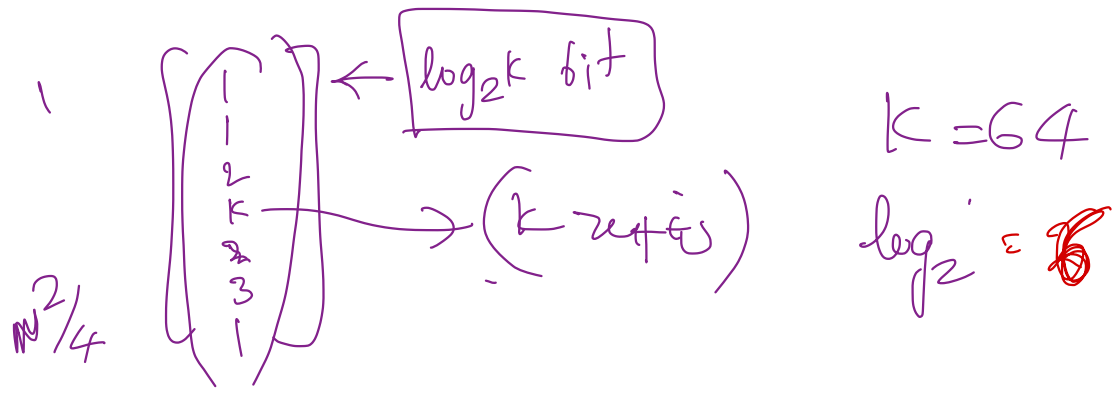
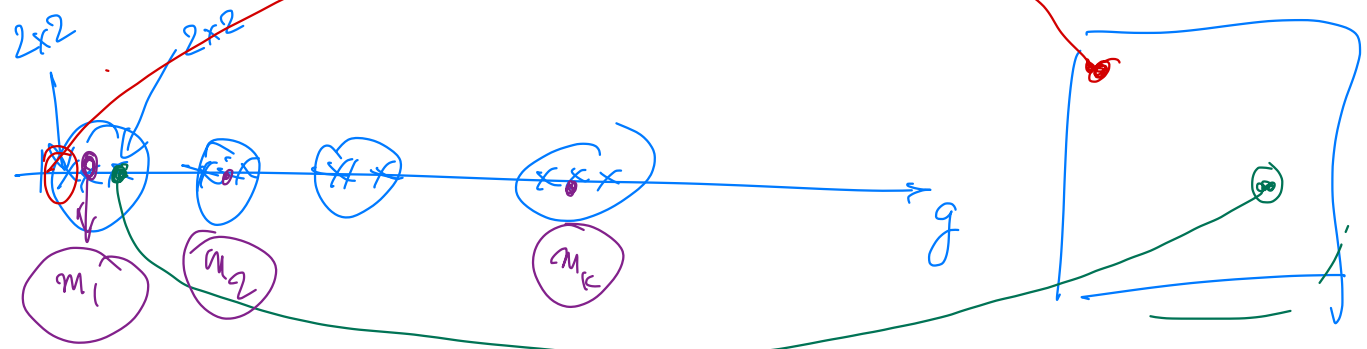


μέγεθος $N^2 \cdot 8$ bits.

$$\left[\bar{G}_4 \quad \bar{G}_4 \right] \quad \frac{N}{2} \times \frac{N}{2} \quad \frac{N^2}{4} \cdot 8 \text{ bit}$$



of clusters vs graylevels of k of codes



4/2/21 exercises

Μέχρι να σταματήσει ο αγγιστικός
σε κάθε βήμα $W(c)$ εφαρμόζεται.

Σταματάει σε τονικό φάσμα $WCC)$ ως προς C .

Επιπαραβάνουμε 2,3 έως 5
σε δύο διαδοχικά βήματα 3 δε υπάρχουν αλλαγές
ως ομάδες