

ΑΝΑΓΝΩΡΙΣΗ ΠΡΟΤΥΠΩΝ

ΠΡΟΑΙΡΕΤΙΚΗ ΕΡΓΑΣΙΑ 1

Σταματόπουλος Βασίλειος – 1115201400188

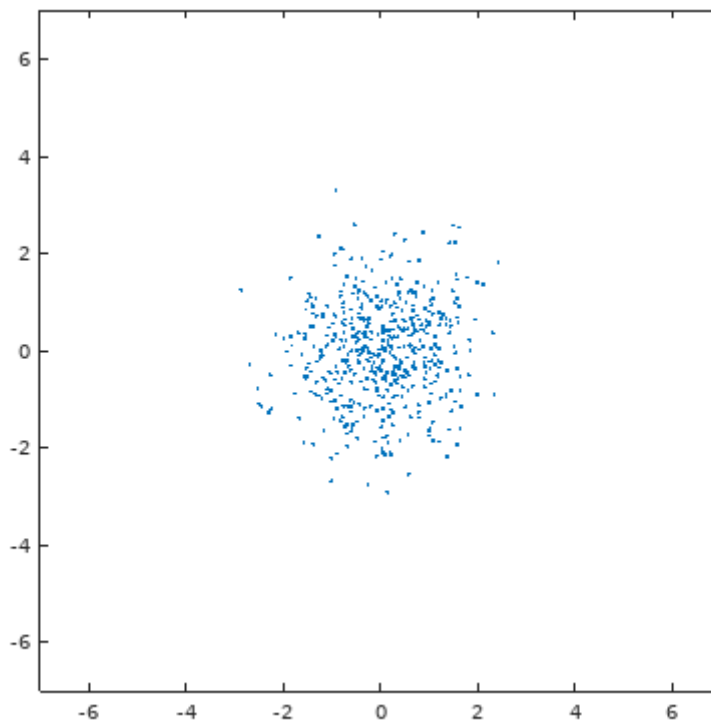
1.

Στην άσκηση αυτή χρησιμοποιήθηκε ο κώδικας που δίνεται, για κάθε ένα από τα παραδείγματα.

Στη συνέχεια αναλύονται τα σχήματα που προκύπτουν για κάθε ομάδα. Αφού η μέση τιμή είναι μηδενική, τα δεδομένα θα κατανομούνται κανονικά γύρω από το (0,0). Συνεπώς, οι διαφορές κάθε σχήματος είναι ανάλογες τις συνδιασποράς που χρησιμοποιείται. Οι τιμές των σ , δίνουν στο γράφημα την πυκνότητα και την κατεύθυνση των στοιχείων. Το ιδιοδιάνυσμα που αντιστοιχεί στην μεγαλύτερη ιδιοτιμή, ορίζει την κατεύθυνση κατά την οποία τα στοιχεία παρουσιάζουν τη μέγιστη διασπορά. Το αμέσως επόμενο, είναι κάθετο στο πρώτο και αντιστοιχεί στη κατεύθυνση στην οποία τα δεδομένα παρουσιάζουν την αμέσως επόμενη μέγιστη διασπορά. Έστω, $S = \begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix}$

α)

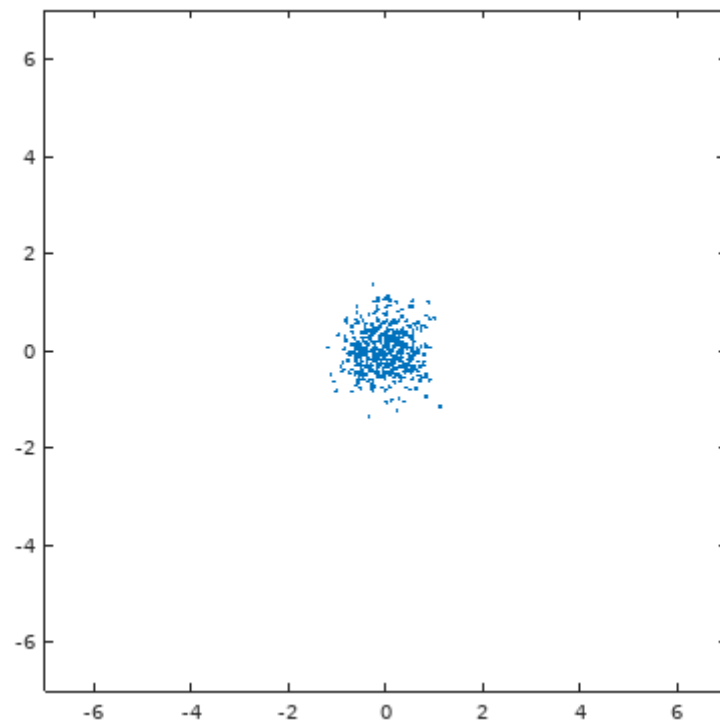
$$m = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad S = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



Στην περίπτωση αυτή, καθώς $\alpha = \delta = 1$ και $\beta = \gamma = 0$ δηλαδή ο πίνακας είναι διαγώνιος με ίσα διαγώνια στοιχεία, τα δεδομένα έχουν ομοιόμορφη κατανομή.

β)

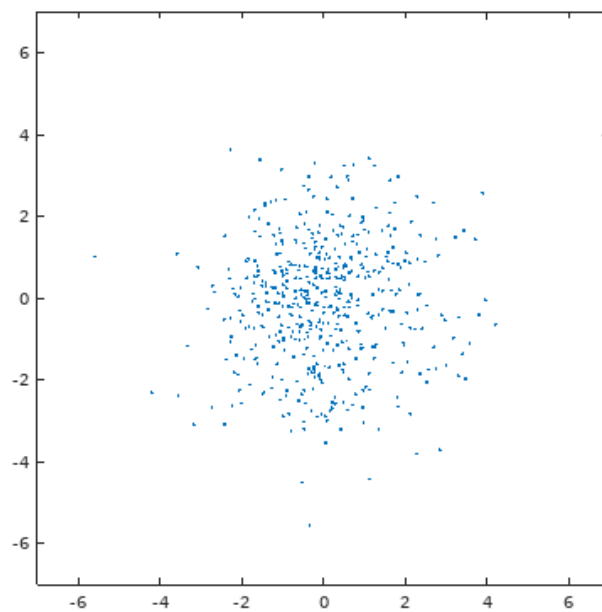
$$m = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad S = \begin{bmatrix} 0.2 & 0 \\ 0 & 0.2 \end{bmatrix}$$



Η περίπτωση αυτή είναι αντίστοιχη με τη προηγούμενη, καθώς όμως το $\alpha = \delta = 0.2 < 1$, τα δεδομένα βρίσκονται πιο κοντά στο (0,0).

γ)

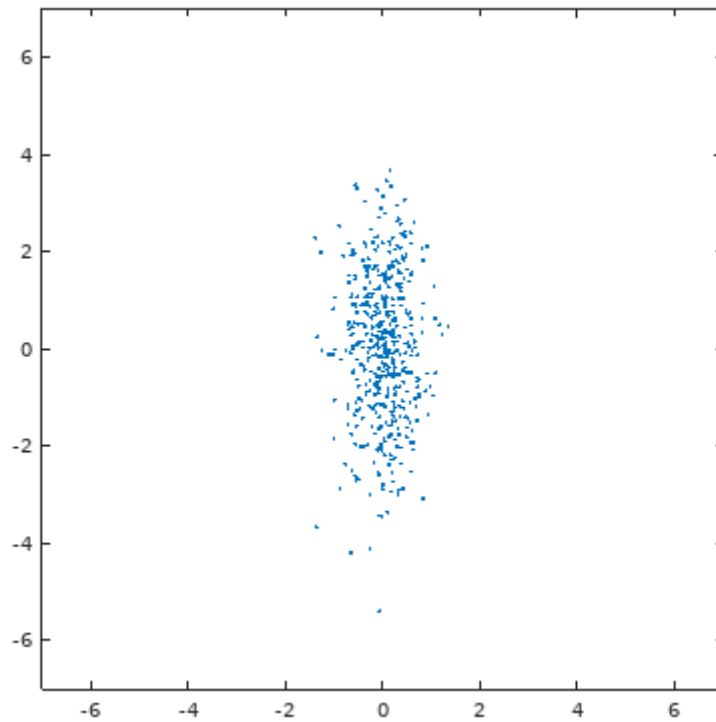
$$m = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad S = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



Ίδια περίπτωση με τις δύο προηγούμενες αλλά τα δεδομένα διασκορπίζονται στον χώρο περισσότερο καθώς $\alpha = \delta = 2$.

δ)

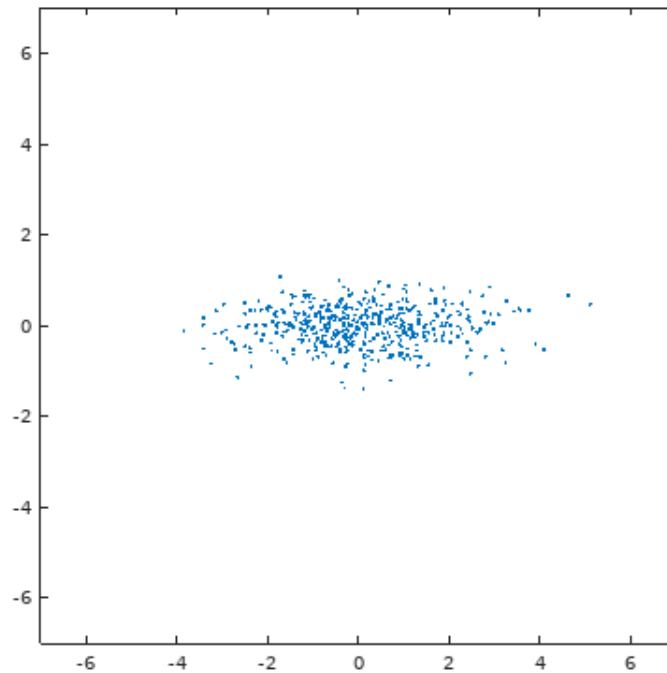
$$m = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad S = \begin{bmatrix} 0.2 & 0 \\ 0 & 0.2 \end{bmatrix}$$



Η περίπτωση αυτή είναι διαγώνιος πίνακας με $\alpha < \delta$ που σημαίνει πως η διασπορά των δεδομένων στον οριζόντιο άξονα είναι αρκετά μικρότερη από εκείνη στον κάθετο.

ε)

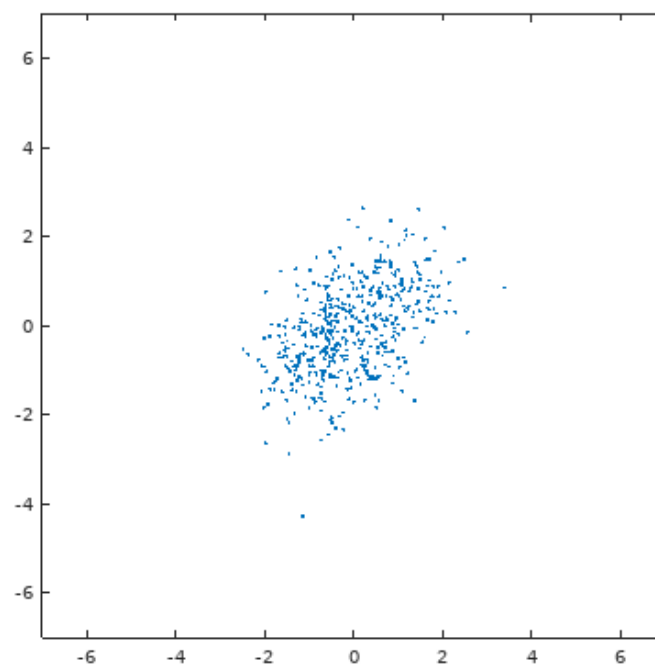
$$m = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad S = \begin{bmatrix} 2 & 0 \\ 0 & 0.2 \end{bmatrix}$$



Παρόμοια περίπτωση με την προηγούμενη αλλά τώρα ισχύει $\alpha \gg \delta$. Συνεπώς οι διασπορές των δεδομένων στον χώρο είναι αντεστραμμένες.

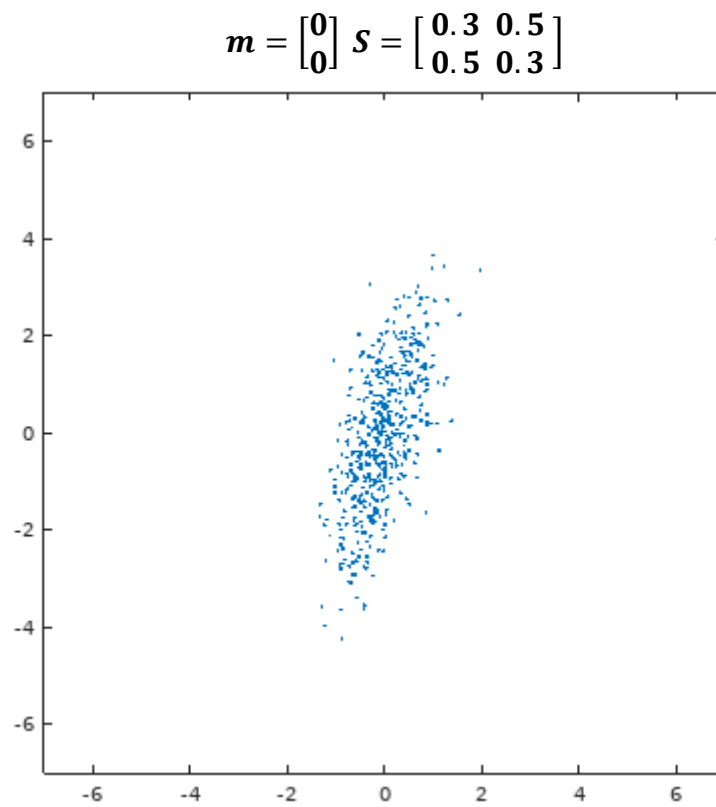
στ)

$$m = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad S = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



Καθώς τώρα το β και το γ έχουν τιμή, τα δεδομένα έχουν κλίση στον χώρο.

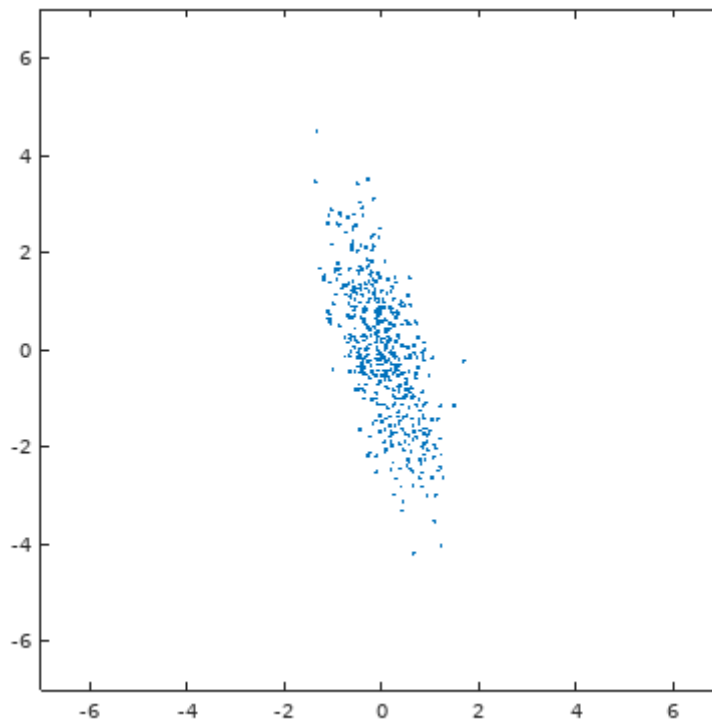
ζ)



Τώρα η κλίση διατηρείται, αλλά αλλάζει η πυκνότητα στον κάθετο και οριζόντιο άξονα.

η)

$$m = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad S = \begin{bmatrix} 0.3 & -0.5 \\ -0.5 & 0.3 \end{bmatrix}$$



Η αντίστροφη περίπτωση της προηγούμενης. (Η κλίση είναι από την άλλη μεριά).

2.

Έχουμε $m_1 = 1$, $m_2 = 4$, $s = 1$.

α)

Ισχύει ότι η a priori πιθανότητες των κλάσεων είναι $p_1 = p_2 = 0.5$. Θα υπολογίσουμε τις a-posteriori πιθανότητες P_1 , P_2 .

$$P_1 = \frac{1}{\sqrt{2\pi}} * \exp\left(-\frac{(1.7 - 1)^2}{2}\right), \quad P_2 = \frac{1}{\sqrt{2\pi}} * \exp\left(-\frac{(1.7 - 4)^2}{2}\right)$$

Καθώς $p_1 = p_2$, η κλάση που θα επιλεγθεί έχει να κάνει μονάχα μόνο με τις a-posteriori πιθανότητες.

$$\text{Έστω, } P_1 > P_2 \rightarrow \ln(P_1) > \ln(P_2) \rightarrow \left(-\frac{0.49}{2}\right) > \left(-\frac{5.29}{2}\right) \rightarrow 5.29 > 0.49$$

Η υπόθεση ισχύει, συνεπώς το $x = 1.7$ θα κατηγοριοποιηθεί στην πρώτη κλάση.

β)

Στο ερώτημα αυτό, χρησιμοποιήθηκε ο δοθέντας κώδικας με την προσθήκη του παρακάτω κομματιού, για τον υπολογισμό του λάθους.

```
...  
  
error = nnz(bayes_res); %Υπολογισμος των μη-μηδενικών αριθμών στον πίνακα  
fprintf("Error is %d percent\n", (error*100)/N);  
  
...
```

Ο κώδικας αυτός υπολογίζει τον αριθμό των μη μηδενικών αριθμών στον πίνακα σύγκρισης και βάση αυτού δίνει το ποσοστό του σφάλματος, που για $N = 100$ ήταν αρκετά καλό (3%).

γ)

Για $N=1000$ και $N=10000$ το σφάλμα κυμαίνεται στις ίδιες τιμές (2.9 και 3.19 αντίστοιχα)

δ)

Η πιθανότητα λάθους για τη κάθε μία πιθανότητα είναι το ποσοστό των 1 στα δύο μισά του πίνακα σύγκρισης. Πιο συγκεκριμένα, για την πρώτη κλάση και για $N=100$ είναι 6% και 8% αντίστοιχα. Ενώ για $N=1000$ και $N = 10000$, τα ποσοστά αυτά δεν αλλάζουν σημαντικά. ($N = 1000$: $e1 = 5.8\%$, $e2 = 8.2\%$. $N=10000$: $e1 = 6.38\%$, $e2 = 7.06\%$.)

3.

Για τους υπολογισμούς χρησιμοποιήθηκε ο δοθέντας κώδικας, σε συνδυασμό με το παρακάτω κομμάτι κώδικα.

```
t=[ones(1,N/2) 2*ones(1,N/2)];  
out_eucl = euclidean(m,X);  
eucl_res = (t~=out_eucl);  
out_mahal = mahalnobis(m,S,X);  
mahal_res = (t~=out_mahal);  
error1 = nnz(eucl_res);  
error2 = nnz(mahal_res);  
fprintf("Error euclidean is %d percent\n", (error1*100)/N);  
fprintf("Error mahalanobis is %d percent\n", (error2*100)/N);
```

β)

Και πήραμε τα αποτελέσματα που περιγράφονται στον πίνακα:

	N =100	N = 1000	N = 10000
Euclidean	23%	25.1%	25.51%
Mahalanobis	23%	22.2%	21.04%

γ)

Όπως φαίνεται, για χαμηλό αριθμό δεδομένων, οι δύο αλγόριθμοι δίνουν παρόμοια αποτελέσματα. Όσο αυξάνει όμως το N, φαίνεται πιο αισθητά η υπεροχή του αλγορίθμου Mahalanobis.

δ)

Όσον αφορά την απόδοση του ταξινομητή Bayes, χρησιμοποιήθηκε το παρακάτω κομμάτι κώδικα για τον υπολογισμό της. Τα αποτελέσματα όπως ήταν αναμενόμενο, ήταν ίδια με του αλγορίθμου ελάχιστης απόστασης Mahalanobis. Αυτό γιατί, οι δύο αλγόριθμοι ταυτίζονται σε αυτό το παράδειγμα.

```
[l,c]=size(m);  
for i=1:N  
    variable = X(:,i);  
    for j=1:c  
        mean = m(:,j);  
        exp_component = (variable - mean)'*(pinv(S))*(variable - mean);  
        val(j) = (1/sqrt(2*pi))*exp(-1/2*exp_component);  
    end  
    [num,z(i)]=max(val);  
end  
  
out_bayes = z;  
bayes_res = (t~=out_bayes);  
error3 = nnz(bayes_res);  
fprintf("Error bayes is %d percent\n", (error3*100)/N);
```


4.

γ)

Αφού ακολουθήσουμε τα βήματα του α και του β, με την υλοποίηση του γ, παίρνουμε $m_1 = 0.83590$ και $m_2 = 2.6196$.

δ)

Οι a priori πιθανότητες των κλάσεων είναι, $p_1 = p_2 = 0.5$. Με αυτά τα στοιχεία ο ταξινομητής Bayes δίνει πιθανότητα λάθους 0.163, που είναι ίση με την πιθανότητα λάθους του αλγορίθμου μέγιστης πιθανοφάνειας.

ε)

Για $N_1 = N_2 = 500$, ο ταξινομητής Bayes δίνει ίδια πιθανότητα λάθους, ενώ η ML, δίνει $p = 0.164$.

ζ)

Δεν θα είχε νόημα ο υπολογισμός των μέσων τιμών με τη βοήθεια της μεθόδου MAP, καθώς ισχύει ότι $\mu_{MAP} \cong \mu_{ML}$

5.

Για μικρό αριθμό δεδομένων ($N=40$), ο ταξινομητής Bayes δίνει καλύτερα αποτελέσματα από τα παράθυρα Parzen ($p_{e1} = 0.163$ $p_{e2} = 0.227$ αντίστοιχα). Με την αύξηση των δεδομένων κάθε κλάσης στα 200, η πιθανότητα λάθους του Bayes δεν αλλάζει, ενώ των Parzen windows μειώνεται σημαντικά ($p_{e2} = 0.166$). Γενικότερα, ο δεύτερος από τους αλγορίθμους αυτούς δίνει καλύτερα αποτελέσματα εάν έχουμε μεγαλύτερο σύνολο δεδομένων, καθώς όσο περισσότερα τα δεδομένα, τόσο περισσότερα θα είναι και τα τεμνόμενα παράθυρα, αρά θα μεγαλώσει και η ακρίβεια της ταξινόμησης.