

Άσκηση 2

Στατιστική Μοντελοποίηση

Αναστάσιος Πατερίτσας , ΑΜ: 2016030065

Θέμα 1.

Ο σκοπός του συγκεκριμένου θέματος είναι να κάνουμε μια εκτίμηση μέσω λογιστικής παλινδρόμησης, όποτε ο στόχος μας είναι να βρούμε τις κατάλληλες παραμέτρους για να κάνουμε καλύτερα την πρόβλεψη. Από την εκφώνηση της άσκησης γνωρίζουμε ότι η συνάρτηση λογιστικής παλινδρόμησης είναι $h_{\theta}(x) = f(\theta^T x)$ όπου υλοποιείται στην συναρτηση sigmoid. Έπειτα υπολογίσαμε την cross entropy (loss function) με βάση τον τύπο :

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left(-y^{(i)} \ln(h_{\theta}(x^{(i)})) - (1 - y^{(i)}) \ln(1 - h_{\theta}(x^{(i)})) \right)$$

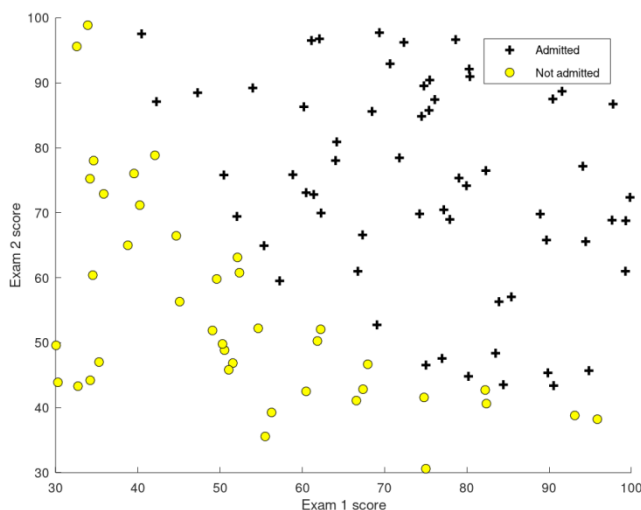
Και υπολογίζουμε την κλήση του σφάλματος με τον τύπο :

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Απόδειξη τέλος pdf*.

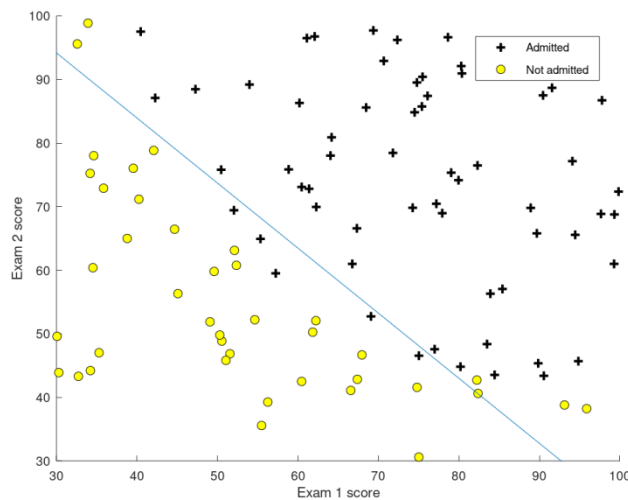
Και οι δυο τύποι υπολογίζονται στην συνάρτηση costFunction.

Αρχικά απεικονίσαμε τα δεδομένα με την plotData



Ενώ στην συνέχεια εκτελέσαμε τα βήματα που προαναφέρθηκαν για τον υπολογισμό του σφάλματος και της κλήσης. Τα αποτελέσματα της συνάρτησης του σφάλματος είναι πολύ λογικά και είναι φυσιολογικό η τιμή του J να είναι πολύ υψηλή (0.69) ενώ το gradient διάνυσμα είναι $[-0.1, -12, -11.26]$.

Εφόσον εφαρμόσαμε την συνάρτηση `fminunc` η οποία βρίσκει το τοπικό ελάχιστο μιας συνάρτησης πολλών παραμέτρων. Τα αποτελέσματα που βρεθήκαν ήταν εμφανώς καλύτερα από τα προηγούμενα (Κόστος: 0.2, gradient: $[-25.16, 0.2, 0.2]$).

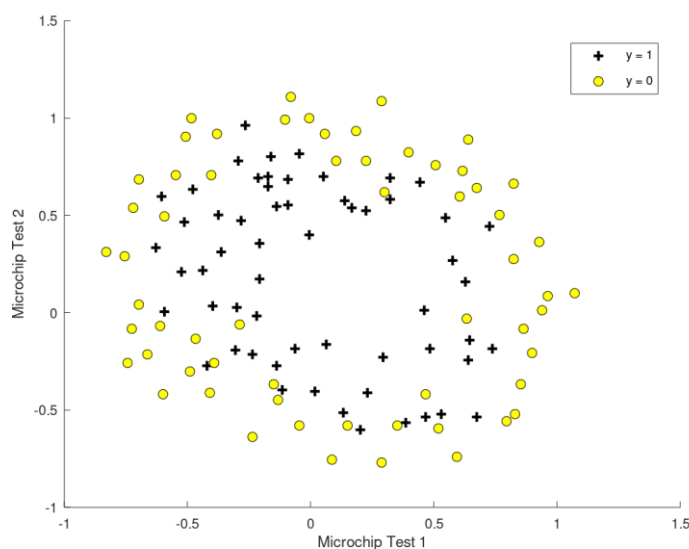


Έχοντας στην διάθεση μας τις συνιστώσες του decision boundary μπορούμε να κάνουμε πρόβλεψη για τον μαθητή που έχει γράψει στο πρώτο διαγώνισμα 45 και στο δεύτερο 85. Η πιθανότητα να περάσει είναι 0.77 με ακρίβεια 0.89. Αν δούμε και στην παραπάνω εικόνα, θα παρατηρήσουμε ότι 11 από τα 100 δείγματα είναι λάθος ταξινομημένα.

Για να πάρουμε την βέλτιστη λύση θα πρέπει η σιγμοειδούς συνάρτησης σε συνδυασμό με την cross entropy θα πρέπει να συνδυαστεί με την εύρεση κατάλληλου τοπικού ελάχιστου.

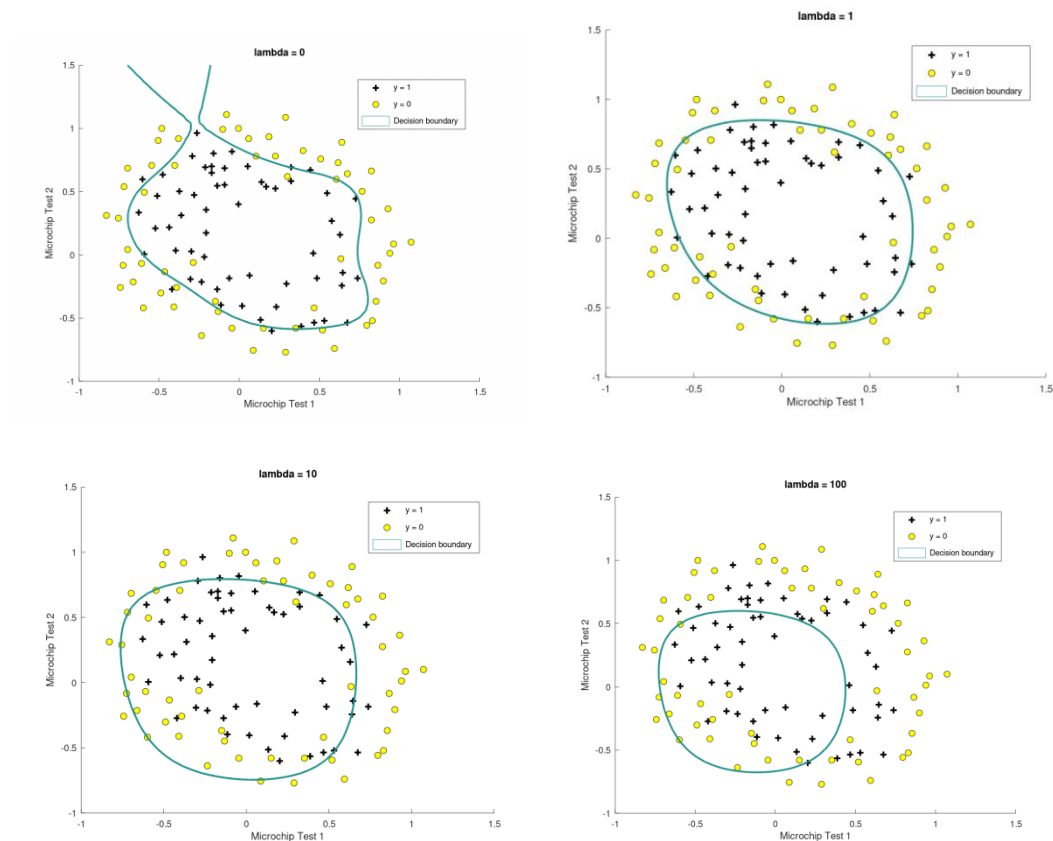
Θέμα 2.

Για το δεύτερο θέμα χρειάστηκε να εφαρμόσουμε ομαλοποιημένη λογιστική παλινδρόμηση για να προβλέψουμε αν περνούν τον έλεγχο ποιότητας χρησιμοποιώντας δυο από τα αποτελέσματα δοκιμών. Τα δεδομένα ήταν τα εξής:



Στην συνέχεια χρειάστηκε να απεικονίσουμε τα δεδομένα σε μεγαλύτερη διάσταση διότι με αυτό τον τρόπο θα μπορούμε να τα διαχωρίσουμε ευκολότερα με την λογιστική παλινδρόμηση. Ο τύπος που μετατρέπει τα δεδομένα σε μεγαλύτερη διάσταση είναι $\sum_{i=0}^6 \sum_{j=0}^i x1^{i-j} x2^j$. Έπειτα χρησιμοποιήσαμε αυτά τα δεδομένα και με $\lambda = 1$ βρήκαμε το κόστος με την τιμή του J να είναι υψηλή (0.69 το initial cost).

Εφόσον εφαρμόσαμε την συνάρτηση `fminunc` η οποία βρίσκει το τοπικό ελάχιστο μιας συνάρτησης πολλών παραμέτρων. Τα αποτελέσματα που βρεθήκαν ήταν εμφανώς καλύτερα από τα προηγούμενα (Κόστος: 0.28, $\lambda = 0$).



Όπως μπορούμε να παρατηρήσουμε , το πόσο καλά αποτελέσματα θα πάρουμε όσον αφορά τα σύνορα εξαρτάται κατά πολύ από το λ . Είναι εύκολο να παρατηρηθεί ότι για $\lambda = 0$ έχουμε τα καλύτερα αποτελέσματα με train accuracy = 86.44 και κόστος 0.28(overfitting λογο του κακο σχημα). Για $\lambda = 1$ έχουμε train accuracy 83.1 , για $\lambda = 10$ train accuracy 74.57 ενώ για $\lambda = 100$ train accuracy 61.01(underfitting) με το κόστος να γίνεται 0.68 δηλαδή λίγο καλύτερο από το initial cost.

Θέμα 3.

Για το τρίτο θέμα έπρεπε να βρούμε τον εκτιμητή μεγίστης πιθανοφανείας της παραμέτρου λ για n δείγματα που παράγονται ανεξάρτητα από μια κατανομή Poisson με παράμετρο λ .

Likelihood function:

$$L_D(\lambda) = \prod_{i=1}^n \lambda^{x_i} \frac{e^{-\lambda}}{x_i!}$$

Η log-likelihood είναι :

$$\begin{aligned} l_D(\lambda) &= \log(L_D(\lambda)) = \sum_{i=1}^n \log\left(\lambda^{x_i} \frac{e^{-\lambda}}{x_i!}\right) = \\ &= \sum_{i=1}^n (-\lambda + x_i \log(\lambda) - \log x_i!) = \\ &= -n\lambda - \sum_{i=1}^n \log x_i! + \log \sum_{i=1}^n x_i \end{aligned}$$

Η Maximum Likelihood Estimation μπορεί να βρεθεί με την επίλυση του κυρτού προβλήματος βελτιστοποίησης $\operatorname{argmax}_{\lambda}(l_D(\lambda))$

$$\nabla_{\lambda_{MLE}} l_D(\lambda) = 0 \leftrightarrow -n + \frac{1}{\lambda_{MLE}} \sum_{i=0}^n x_i = 0$$

$$\lambda_{MLE} = \frac{1}{n} \sum_{i=0}^n x_i$$

Ωστόσο το λ_{MLE} είναι ο μέσος Όρος των n παρατηρήσεων στο δείγμα.

Θέμα 4.

Σε αυτό το θέμα κληθήκαμε να φτιάξουμε ένα Naïve Bayes classifier για αναγνώριση ψηφίων από το 0 έως το 9. Στο αρχείο digits.mat υπάρχουν 10 κλάσεις όπου σε κάθε μια κλάση 500 στοιχεία τα οποία είναι μια ασπρόμαυρη εικόνα με διαστάσεις 28x28 οι οποίες αναπαριστώνται σαν πίνακας 784x1 με τιμές 0 η 1. Για τον Naïve Bayes ,υποθέσαμε ότι κάθε pixel είναι ανεξάρτητο από τα άλλα , και χρησιμοποιήσαμε την κατανομή Bernouli για να μοντελοποιήσουμε κάθε πιξελ.

Likelihood function:

$$L_{D_{x^{yi}}} = P(D_{x^{yi}} | p^{yi}) = \prod_{i=1}^n P(X = x_j^{yi}) = (p^{yi})^k (1 - p^{yi})^{k'}$$

Το $k = \sum_{j=1}^n x_j^{yi}$ και $k' = 1 - k$.

Log likelihood

$$l_{D_{x^{yi}}}(p^{yi}) = \log \left((p^{yi})^k (1 - p^{yi})^{k'} \right) = k \log(p^{yi}) + k' \log(1 - p^{yi})$$

$$\nabla_{p^{yi}_{MLE}} l_{D_{x^{yi}}}(p^{yi}) = 0 \leftrightarrow \frac{k}{p^{yi}_{MLE}} - \frac{k'}{1 - p^{yi}_{MLE}} = 0$$

$$p^{yi}_{MLE} = \frac{k}{k + k'}$$

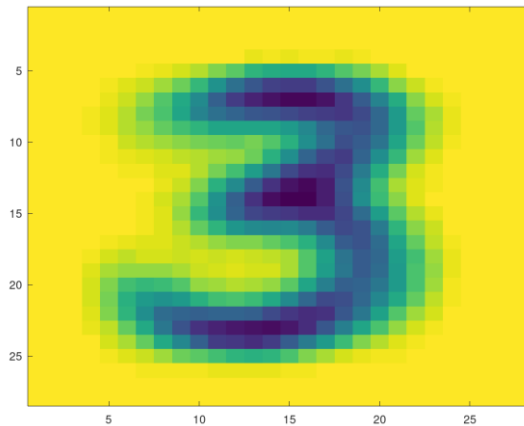
$$p^{yi}_{MLE} = \frac{1}{n} \sum_{j=1}^n x_j^{yi}$$

Όποτε η p^{yi}_{MLE} δίνεται από τον μέσο όρο του πιξελ της κλασης.

Χρησιμοποιώντας τον παραπάνω τύπο εκπαιδεύσαμε τα μοντέλα για κάθε ψηφίο αφού πρώτα υπολογίσαμε τους εκτιμητές μεγίστης πιθανοφανείας (άθροισμα όλων των γραμμών της κάθε κλάσης διαιρεμένο με τον αριθμό των εικόνων που έχει δηλαδή 500) Στην συνέχεια αφού

εκπαιδεύσαμε τα μοντέλα για κάθε αριθμό οπτικοποιήσαμε το
κάθε μοντέλο για το κάθε ψηφίο.

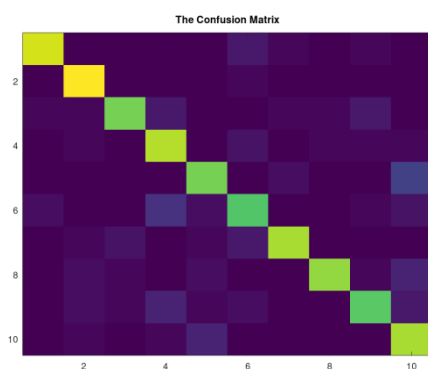
Παράδειγμα του αριθμού 3



Έπειτα εφαρμόσαμε την αντίστοιχη διαδικασία με τα
εκπαιδευμένα μοντέλα για την πρόβλεψη των αριθμών στο
τεστ σετ .Για κάθε εικόνα του σετ επιλέχτηκε το στοιχείο που
είχε την μεγαλύτερη λογαριθμημένη πιθανότητα. Με τα εξής
αποτελέσματα :

```
Accuracy for testing digits 0 : 0.882
Accuracy for testing digits 1 : 0.946
Accuracy for testing digits 2 : 0.748
Accuracy for testing digits 3 : 0.83
Accuracy for testing digits 4 : 0.752
Accuracy for testing digits 5 : 0.69
Accuracy for testing digits 6 : 0.824
Accuracy for testing digits 7 : 0.788
Accuracy for testing digits 8 : 0.696
Accuracy for testing digits 9 : 0.818
Accuracy: 0.7974
```

Ενώ ο confusion matrix είναι ο εξής :



Θέμα 5.

Στο συγκεκριμένο θέμα έπρεπε να πειραματιστούμε με n μετρήσεις τις χρονικές στιγμές $k=1,2,\dots,n$. Ενώ υποθέτουμε ότι η συνάρτηση πυκνότητας πιθανότητας των μετρήσεων είναι Gauss.

Για το α ερώτημα έπρεπε να ζωγραφίσουμε στο ίδιο σχήμα την δεσμευμένη πιθανότητα $p(\mu | H(n))$ καθώς το n μεταβάλλεται από 1 έως 25. Από θεωρία γνωρίζουμε ότι:

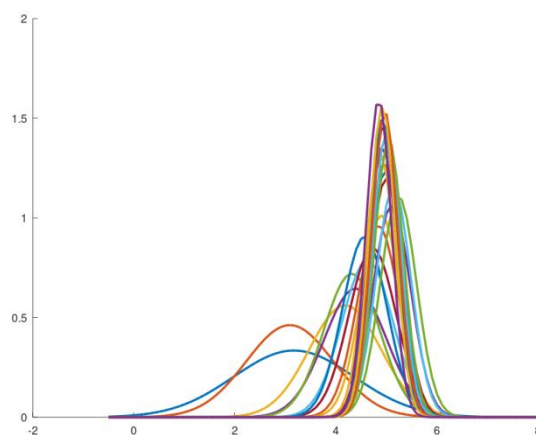
$$p(\mu | H(n)) = \frac{1}{\sqrt{2\pi} \sigma_n} \exp\left[-\frac{1}{2} \left(\frac{\mu - \mu_n}{\sigma_n}\right)^2\right]$$

Με

$$\mu_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \bar{x}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0, \quad \sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}.$$

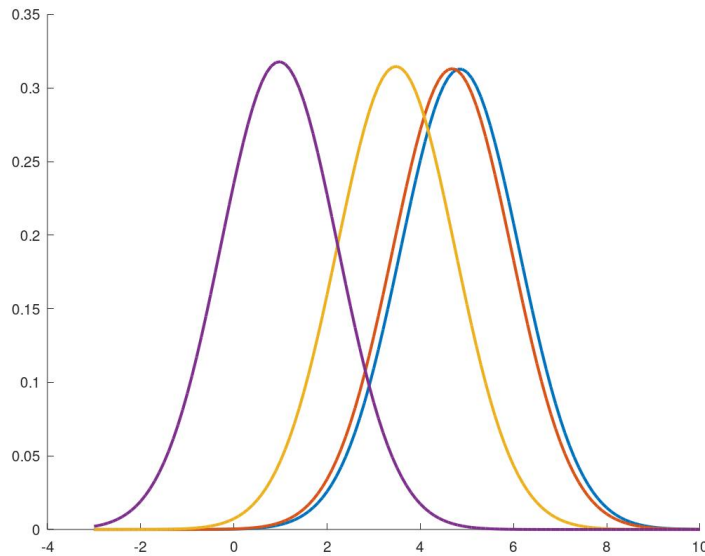
$$\bar{x}_n = \frac{1}{n} \sum_{k=1}^n x_k.$$

Όμως από την εκφώνηση γνωρίζουμε το $\mu_0=0$ όποτε ο δεύτερος Όρος της πρόσθεσης στο μ_n μηδενίζεται.



Για το δεύτερο ερώτημα κληθήκαμε να σχεδιάσουμε $p(\chi|H(n))$ που προκύπτουν από την Bayesian εκτίμηση όταν το $\sigma_0^2 = 10\sigma^2$, $\sigma_0^2 = \sigma^2$, $\sigma_0^2 = 0.1\sigma^2$ και $\sigma_0^2 = 0.01\sigma^2$

Από την θεωρία προκύπτει ότι η $p(\chi|H(n))$ είναι ομοιομορφα κατανοημένη με mean μ_n και variance $\sigma^2 + \sigma_n^2$. Όποτε με τον ίδιο τρόπο όπως στο ερώτημα α αλλά για $n=25$ έχουμε :



Σε αυτό το σημείο παρατηρήσαμε ότι η πιθανότητες δε αλλάζουν για τις διαφορές τιμές του σ^2 (πλάτος σταθερό) αλλά μόνο μετατοπίζονται στον άξονα του χ (ολισθαίνει αριστερά).

Θέμα 7.

Για το μέρος πρώτο έπρεπε να εφαρμόσουμε SVM για γραμμικά διαχωρίσιμα δείγματα. Για να το κάνουμε αυτό έπρεπε να υπολογίσουμε την Λαγκαρσιανή του δυικού προβλήματος. Για να γίνει αυτό χρησιμοποιούμε γραμμικό kernel που δίνεται από τον τύπο :

$$K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \mathbf{x}^{(i)} \cdot (\mathbf{x}^{(j)})^T$$

Ενώ η Λαγκαρσιανή

$$L(\lambda) = \sum_{i=1}^n \lambda^{(i)} - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda^{(i)} \lambda^{(j)} y^{(i)} y^{(j)} K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

Την όποια την υπολογίσαμε με την συνάρτηση της matlab quadprog. Για τον υπολογισμό του διανύσματος \mathbf{w} χρησιμοποιήσαμε τον τύπο :

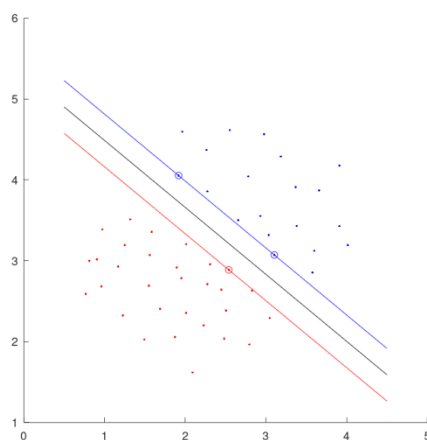
$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha} \mathbf{y}$$

Όπου α είναι τα λ . Ενώ υπολογίσαμε το bias με τον τύπο :

$$w_0 = (\mathbf{y} - \mathbf{X}\mathbf{w})_{\epsilon S} \text{ που } S \text{ τα Support vector subset}$$

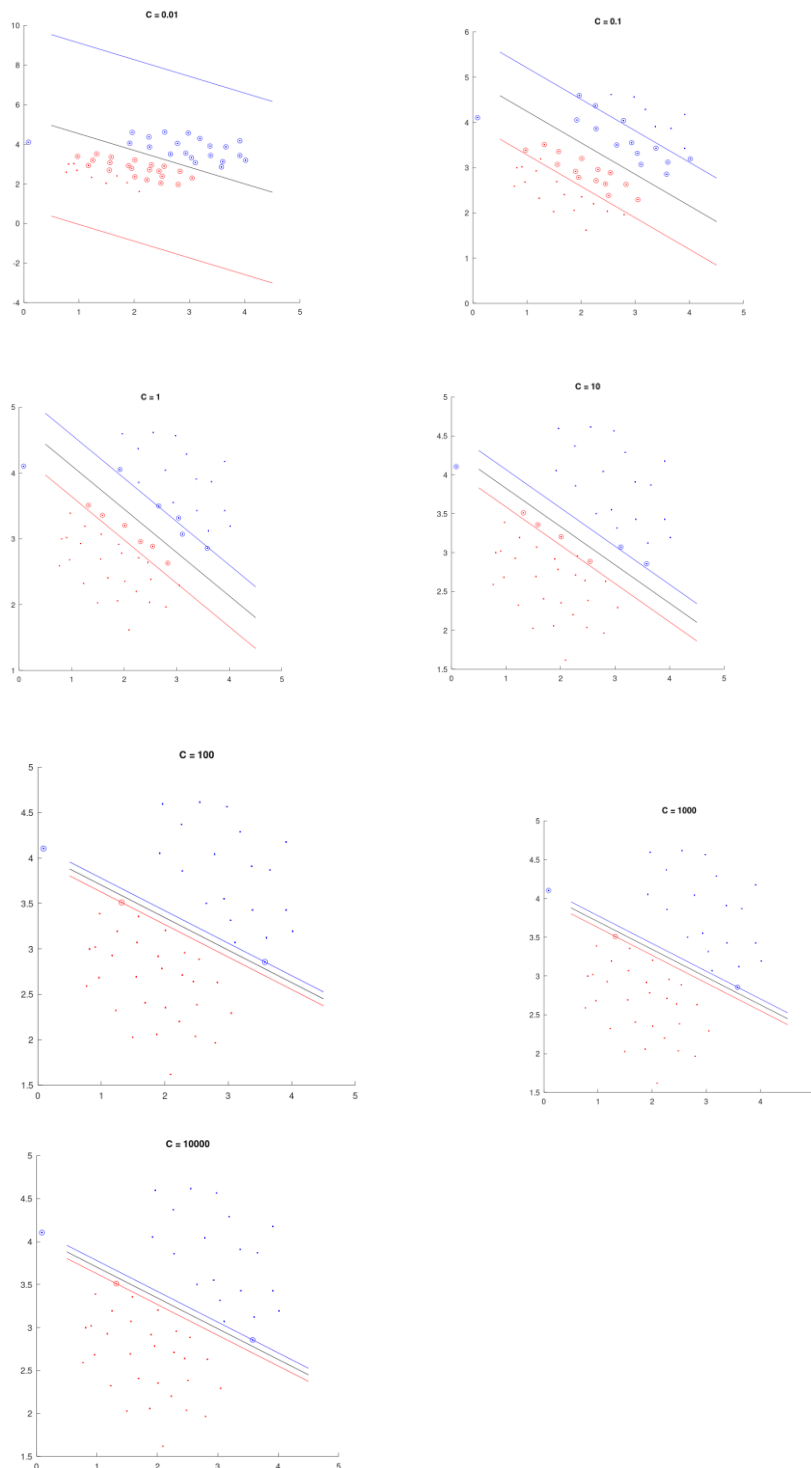
Ενώ το width = $-w_0 / \max(\mathbf{w})$

Η γραφική παράσταση διαχωρισμού είναι



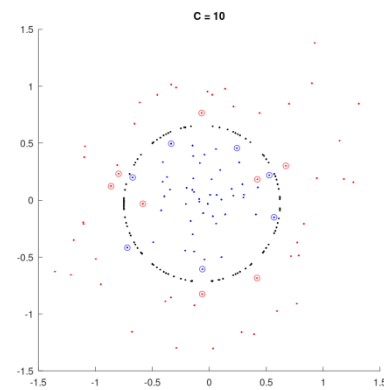
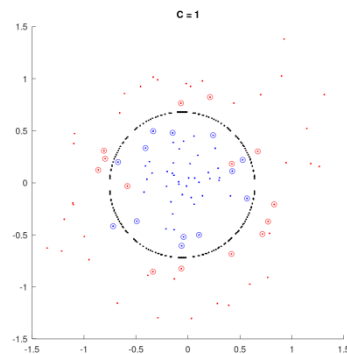
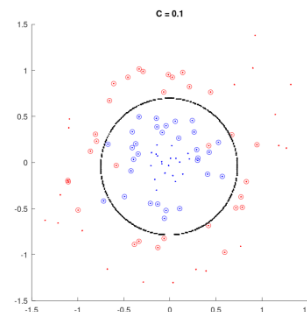
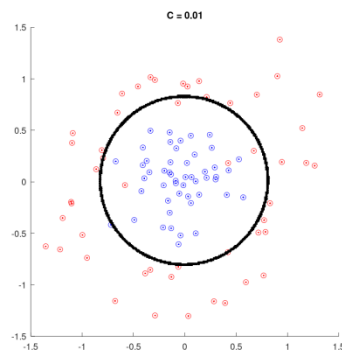
Εύκολα παρατηρούνται τα 3 Support Vector Machines 1 τα αρνητικά και 2 για τα θετικά. Με κόκκινο είναι η παράλληλη στην γραμμή διαχωρισμού που παρανοεί από το αρνητικό SVM ενώ η μπλε αντίστοιχα για τα θετικά SVM.

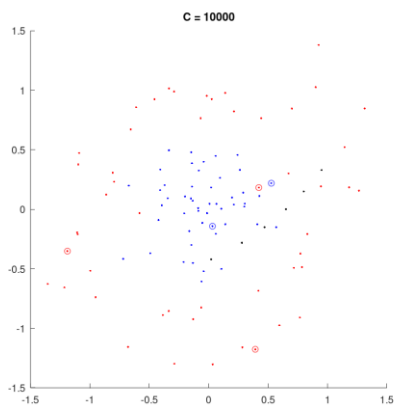
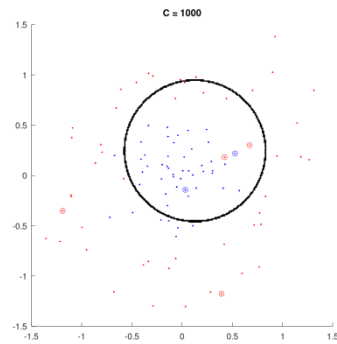
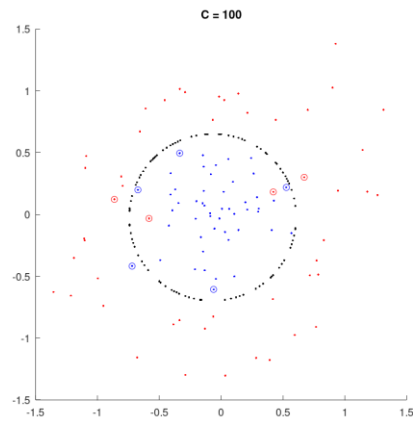
Για το β μέρος χρησιμοποιήσαμε τον ίδιο κώδικα αλλά για τα διάφορα C τα όποια περιορίζουν τα λ :



Όπως μπορεί να παρατηρηθεί όσο μικρότερο είναι το C τόσο λιγότερα λ έχουμε άρα τα αποτελέσματα είναι χειρότερα και έχουμε παρά πολλά SVM με αποτέλεσμα να έχουμε misclassifications.

Για το γ ερώτημα τα δεδομένα μας δε είναι γραμμικά όποτε δε θα μπορούσαμε να χρησιμοποιήσουμε τον συγκεκριμένο kernel. Για τον λόγο αυτό προσθέτουμε ένα τρίτο χαρακτηριστικό στο διάνυσμα , και με αυτό το trick τα δεδομένα γίνονται γραμμικά χωρισμένα στον τρισδιάστατο χώρο.





Όπως και πριν για μικρά C βρίσκουμε πολλά SVM όμως για πολύ μεγάλο C βρίσκουμε πολύ λίγα και δε μπορεί να βρει decision Boundary.

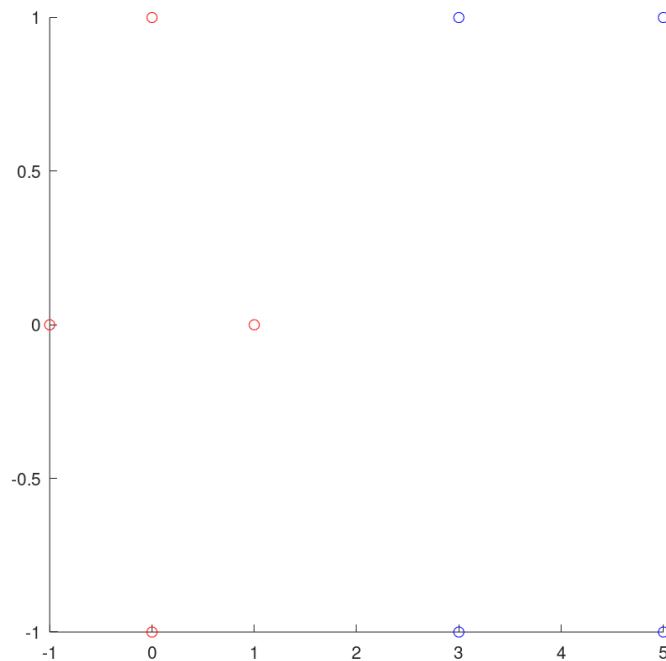
Συμπέρασμα για β, γ :

Το C πρέπει να επιλεγχεται στο χερι για την καθε διανομη.

Θεμα 6.

A)

Τα δειγματα στον χώρο είναι τα εξής:



Τα 3 svm θα είναι το σημείο (1,0) για τα αρνητικά και (3,-1) , (3,1) για τα θετικά . Η γραμμή διαχωρισμού θα είναι η κάθετη ευθεία που διέρχεται από $x = 2$.

Λόγος 6

Τα 3 SVM

$$S_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, S_2 = \begin{pmatrix} 3 \\ 1 \end{pmatrix}, S_3 = \begin{pmatrix} 3 \\ -1 \end{pmatrix}$$

$$0 \leq \text{αριθμοί που ισχύουν } \begin{matrix} \omega^T x + \omega_0 \geq 1 \\ \omega^T x + \omega_0 \leq 1 \end{matrix}$$

για να επιβληθούν τις αυξήσεις θα χρησιμοποιήσουμε $\omega_0 = -2$

$$\begin{aligned} [3, 1] & \quad 3\omega_1 + \omega_2 - 3 \geq 0 \\ [3, -1] & \quad 3\omega_1 - \omega_2 - 3 \geq 0 \\ [1, 0] & \quad -\omega_1 + 1 \geq 0 \end{aligned}$$

$$L(\omega, d_1, d_2, d_3) = \frac{\omega_1^2 + \omega_2^2}{2} - d_1(3\omega_1 + \omega_2 - 3) - d_2(3\omega_1 - \omega_2 - 3) - d_3(-\omega_1 + 1)$$

$$\frac{dL(\omega, d_1, d_2, d_3)}{d\omega_1} = 0 \Rightarrow \omega_1 = 3d_1 + 3d_2 - d_3$$

$$\frac{dL(\omega, d_1, d_2, d_3)}{d\omega_2} = 0 \Rightarrow \omega_2 = d_1 - d_2$$

$$d_1(3\omega_1 + \omega_2 - 3) = 0$$

$$d_2(3\omega_1 - \omega_2 - 3) = 0$$

$$d_3(-\omega_1 + 1) = 0$$

για $d_1, d_2, d_3 \neq 0$

$$3\omega_1 + \omega_2 - 3 = 0$$

$$3\omega_1 - \omega_2 - 3 = 0$$

$$-\omega_1 + 1 = 0$$

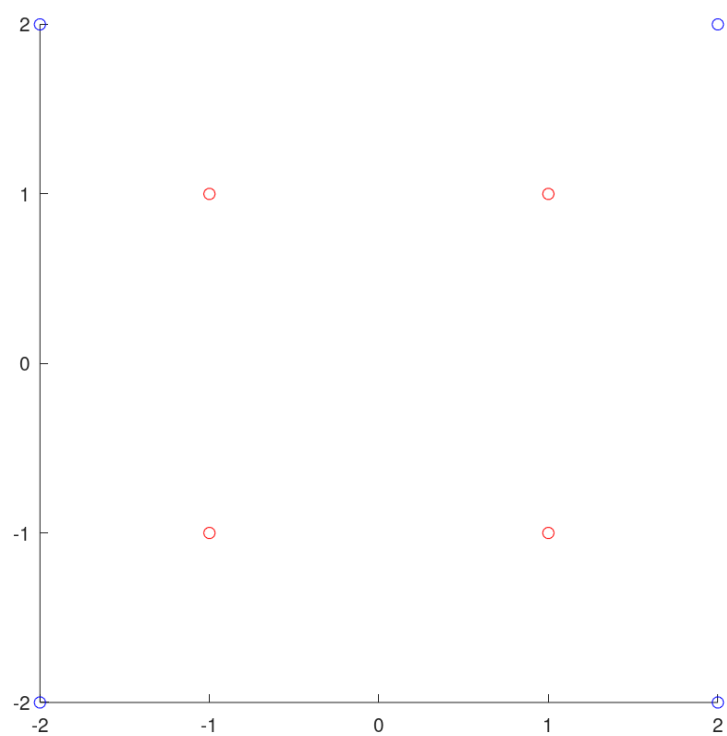
$$\omega_1 = 1 \quad \text{και} \quad \omega_2 = 0$$

$$\text{Άρα τότε έχουμε } y = \omega x + b \Rightarrow y = \begin{pmatrix} 1 \\ 0 \end{pmatrix} x + 2$$

Άρα είναι η ευθεία $x = 2$

B)

Τα δείγματα στο B ερώτημα δε είναι γραμμικά χωρισμένα οπότε θα πρέπει να προσθέσουμε ένα τρίτο χαρακτηριστικό για να τα κάνουμε γραμμικά διαχωρίσιμα.



Λόγιστος 1

$$J(\theta) = \frac{1}{n} \sum_{i=1}^m (-y^i \log(h\theta(x^i)) - (1-y^i) \log(1-h\theta(x^i)))$$

Η απάντηση κολλάει είναι η εξής

$$\begin{aligned} \frac{dJ(\theta)}{d\theta_j} &= \frac{d}{ds} \left[\frac{1}{n} \sum_{i=1}^m [y^i \log(h\theta(x^i)) + (1-y^i) \log(1-h\theta(x^i))] \right] \stackrel{\text{για κάθε } i \text{ και } \theta_j}{=} h\theta(x^i) = \sigma(\theta^T x^i) \\ &= -\frac{1}{n} \sum_{i=1}^m \left[y^i \frac{\frac{d}{d\theta_j} \sigma(\theta^T x^i)}{h\theta(x^i)} + (1-y^i) \frac{\frac{d}{d\theta_j} (1-h\theta(x^i))}{1-h\theta(x^i)} \right] = \\ &= -\frac{1}{n} \sum_{i=1}^m \left[y^i \frac{\sigma(\theta^T x^i)(1-\sigma(\theta^T x^i)) \frac{d}{d\theta_j} \theta^T x^i}{h\theta(x^i)} - (1-y^i) \frac{\sigma(\theta^T x^i)(1-\sigma(\theta^T x^i)) \frac{d}{d\theta_j} \theta^T x^i}{1-h\theta(x^i)} \right] \\ \frac{d}{d\theta_j} \theta^T x^i &= x_j^i \\ &= -\frac{1}{n} \sum_{i=1}^m [y^i (1-h\theta(x^i)) x_j^i - (1-y^i) h\theta(x^i) x_j^i] = \\ &= -\frac{1}{n} \sum_{i=1}^m [x_j^i - y^i h\theta(x^i) - h\theta(x^i) + y^i h\theta(x^i)] x_j^i = \\ &= -\frac{1}{n} \sum_{i=1}^m [y^i - h\theta(x^i)] x_j^i = \frac{1}{n} \sum_{i=1}^m h\theta(x^i - y^i) x_j^i \end{aligned}$$

Λόγιστος 2

Ομοια με πριν με την μικρή διαφορά

$$\frac{d}{d\theta} \left[\frac{1}{n} \sum_{j=1}^n \theta_j^2 \right] = \frac{2}{n} \sum_{j=1}^n \theta_j = \sum_{j=1}^n \frac{2}{n} \theta_j$$