

Επεξεργασία Δεδομένων και Αλγόριθμοι Μάθησης

Αναφορά Σετ Ασκήσεων 2

Όνομα: Μακάριος Χρηστάκης

Περιεχόμενα

Περιεχόμενα	1
Πρόβλημα 1	1
Πρόβλημα 2	4

Πρόβλημα 1

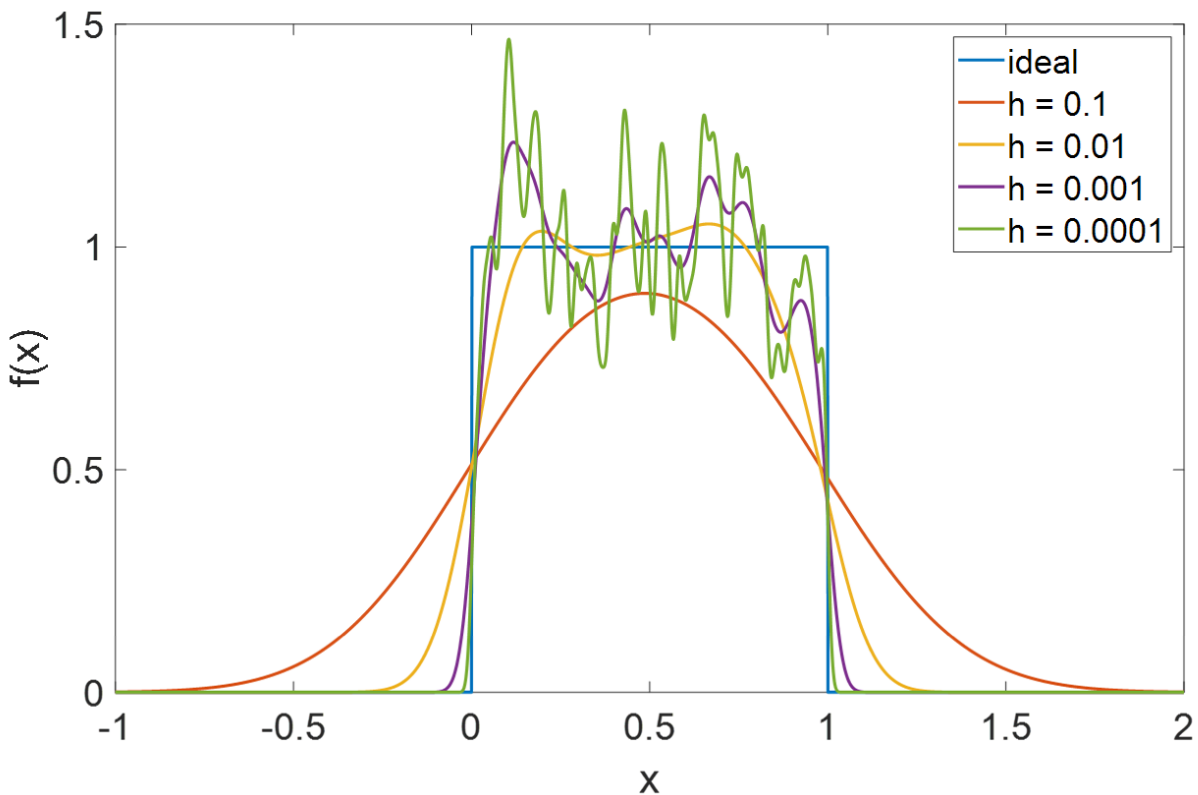
Όπως μάθαμε στην θεωρία για να προσεγγίσουμε την συνάρτηση πυκνότητας πιθανότητας $f(x)$ μιας τυχαίας μεταβλητής στο σημείο x , έχοντας στη διάθεση μας N υλοποιήσεις της τυχαίας μεταβλητής αυτής $x_i, i = 1 \dots N$ μπορούμε να χρησιμοποιήσουμε κάποια από τις συναρτήσεις kernel που προσεγγίζουν την κρουστική συνάρτηση και να υπολογίσουμε την παράσταση:

$$\hat{f}(x) = \frac{\sum_{i=1}^N K(x - x_i, h)}{N} \quad (1)$$

Για το σκοπό αυτό δημιουργήθηκαν 1000 υλοποιήσεις από ομοιόμορφη κατανομή στο διάστημα $[0,1]$. Παρακάτω θα μελετηθεί η επίδραση των διαφόρων kernel function, όπως και της παραμέτρου h , στην προσέγγιση αυτής της συνάρτησης πυκνότητας πιθανότητας.

Ερώτημα (α) Gaussian Kernel

Χρησιμοποιώντας τον Gaussian Kernel υπολογίζουμε την προσέγγιση της PDF μέσω της σχέσης (1) σε 2000 σημεία στο διάστημα $[-1,2]$ και μεταβάλλουμε την τιμή της παραμέτρου h . Τα αποτελέσματα φαίνονται παρακάτω:

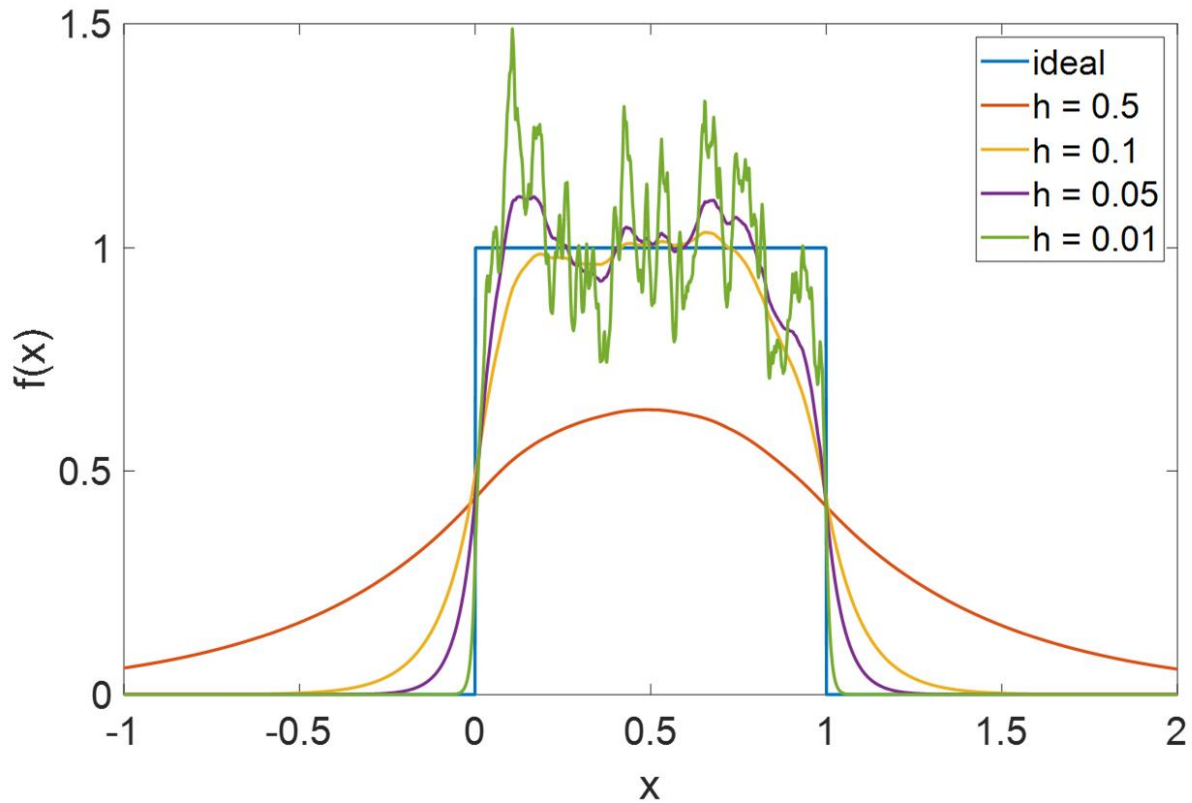


Σχήμα 1: Προσέγγιση με Gaussian Kernel

Όπως φαίνεται στο Σχήμα 1, όσο μειώνουμε την τιμή του h η προσέγγιση πετυχαίνει με όλο και μεγαλύτερη ακρίβεια τα όρια της ιδανικής καμπύλης όπου υπάρχει η ασυνέχεια, ωστόσο εισάγεται θόρυβος στην προσέγγιση της τιμής 1 της ιδανικής καμπύλης. Ο θόρυβος αυτός μειώνεται όσο αυξάνεται το h , αλλά φυσικά χάνονται τα όρια της ιδανικής μεταβολής.

Ερώτημα (β) **Laplacian Kernel**

Για τον Laplacian Kernel υπολογίζοντας στα ίδια 2000 σημεία την προσέγγιση της PDF για διάφορες τιμές του h παίρνουμε:



Σχήμα 2: Προσέγγιση με Laplacian Kernel

Όπως φαίνεται στο Σχήμα 2 και με τον Laplacian Kernel έχουμε το tradeoff μεταξύ ακρίβειας προσέγγισης των ορίων της καμπύλης. Ο Laplacian Kernel ωστόσο αρχίζει να εμφανίζει αυτόν τον θόρυβο σε μεγαλύτερες κατά περίπου μια τάξη μεγέθους τιμές του h απ' ότι ο Gaussian Kernel.

Πρόβλημα 2

Ερώτημα (α)

Η συνάρτηση $\Phi(\mathbf{X})$ που θέλουμε να προσεγγίσουμε έγκειται στο χώρο V που ορίζει ο Gaussian Kernel:

$$K(\mathbf{X}, \mathbf{Y}) = e^{\frac{-\|\mathbf{X}-\mathbf{Y}\|^2}{h}}$$

Η συνάρτηση $\hat{\Phi}(\mathbf{X})$ ανήκει στον γραμμικό υποχώρο Ω του V που ορίζεται από τις συναρτήσεις $K(\mathbf{X}, \mathbf{X}_i)$ όπου \mathbf{X}_i τα γνωστά μας παραδείγματα. Συνεπώς από το representer theorem θα έχουμε:

$$\langle \Phi(\mathbf{X}) - \hat{\Phi}(\mathbf{X}), K(\mathbf{X}, \mathbf{X}_i) \rangle = 0 \Rightarrow \langle \Phi(\mathbf{X}), K(\mathbf{X}, \mathbf{X}_i) \rangle = \langle \hat{\Phi}(\mathbf{X}), K(\mathbf{X}, \mathbf{X}_i) \rangle \Rightarrow \hat{\Phi}(\mathbf{X}_i) = \Phi(\mathbf{X}_i)$$

Συνεπώς οι τιμές της ιδανικής, άγνωστης συνάρτησης $\Phi(\mathbf{X})$ στα σημεία \mathbf{X}_i που γνωρίζουμε θα ταυτίζονται με τις τιμές της $\hat{\Phi}(\mathbf{X})$, στα σημεία αυτά. Άρα μπορούμε στο κριτήριο κόστους να αντικαταστήσουμε στα πρώτα 2 αθροίσματα την συνάρτηση $\Phi(\mathbf{X})$ με $\hat{\Phi}(\mathbf{X})$.

Ερώτημα (β)

Αρχικά μπορούμε να εκφράσουμε την τελευταία νόρμα του κριτηρίου ως:

$$\|\Phi(\mathbf{X})\|^2 = \|\Phi(\mathbf{X}) - \hat{\Phi}(\mathbf{X}) + \hat{\Phi}(\mathbf{X})\|^2 \Rightarrow$$

$$\|\Phi(\mathbf{X})\|^2 = \|\Phi(\mathbf{X}) - \hat{\Phi}(\mathbf{X})\|^2 + \|\hat{\Phi}(\mathbf{X})\|^2 + 2 \langle \Phi(\mathbf{X}) - \hat{\Phi}(\mathbf{X}), \hat{\Phi}(\mathbf{X}) \rangle \quad (2)$$

Όμως αφού η $\hat{\Phi}(\mathbf{X})$ είναι κάθετη προβολή της $\Phi(\mathbf{X})$ στον υποχώρο Ω , το σφάλμα $\Phi(\mathbf{X}) - \hat{\Phi}(\mathbf{X})$, θα είναι κάθετο στον υποχώρο Ω που βρίσκεται η $\hat{\Phi}(\mathbf{X})$, οπότε το εσωτερικό γινόμενο των δύο διανυσμάτων αυτών θα είναι μηδενικό. Συνεπώς από τη σχέση (2) θα έχουμε:

$$\|\Phi(\mathbf{X})\|^2 = \|\hat{\Phi}(\mathbf{X})\|^2 + \|\Phi(\mathbf{X}) - \hat{\Phi}(\mathbf{X})\|^2 \Rightarrow \|\hat{\Phi}(\mathbf{X})\|^2 = \|\Phi(\mathbf{X})\|^2 - \|\Phi(\mathbf{X}) - \hat{\Phi}(\mathbf{X})\|^2$$

Κι αφού ο τελευταίος όρος είναι μη αρνητικός καταλήγουμε στο:

$$\|\hat{\Phi}(\mathbf{X})\|^2 \leq \|\Phi(\mathbf{X})\|^2$$

Αφού λοιπόν αποδείξαμε πως οι τιμές της συνάρτησης $\hat{\Phi}(\mathbf{X})$ στα σημεία \mathbf{X}_i είναι ίδιες με αυτές της συνάρτησης $\Phi(\mathbf{X})$ στα σημεία αυτά και η νόρμα της $\hat{\Phi}(\mathbf{X})$ στον τελευταίο όρο του κριτηρίου θα είναι το πολύ ίση με την νόρμα της $\Phi(\mathbf{X})$, αντικαθιστώντας την $\Phi(\mathbf{X})$ με $\hat{\Phi}(\mathbf{X})$ στο κριτήριο και ελαχιστοποιώντας το ως προς την $\hat{\Phi}(\mathbf{X})$, το Loss που θα προκύψει θα είναι μικρότερο ή ίσο του βέλτιστου που θα πρόκυπτε εάν ελαχιστοποιούσαμε την αρχική παράσταση! Συνεπώς, μπορούμε να αντικαταστήσουμε την $\Phi(\mathbf{X})$ με $\hat{\Phi}(\mathbf{X})$ στο κριτήριο αυτό.

Ερώτημα (γ)

Στο κριτήριο προς ελαχιστοποίηση είναι

$$\min_{\hat{\Phi} \in \Omega} \sum_{X_i \in stars} (1 - \hat{\Phi}(X_i))^2 + \sum_{X_j \in circles} (1 + \hat{\Phi}(X_j))^2 + \lambda \|\hat{\Phi}(X)\|^2 \quad (3)$$

Ορίζοντας $y_i = \pm 1$ η επιθυμητή τιμή της $\Phi(X)$ για το παράδειγμα X_i και

$$\hat{\Phi}(X) = \mathbf{c}^T \mathbf{K}(X)$$

όπου

$$\mathbf{c} = [a_1 \dots a_N \ \beta_1 \dots \beta_M]^T$$

$$\mathbf{K}(X) = [K(X, X_1) \dots K(X, X_N) \ K(X, X_{N+1}) \dots K(X, X_{N+M})]^T$$

$$\text{με τα γνωστά παραδείγματα } X_i \in \begin{cases} stars, & i = 1 \dots N \\ circles, & i = N + 1 \dots N + M \end{cases}$$

Η εξίσωση (3) μπορεί να γραφεί ισοδύναμα ως:

$$\min_{\mathbf{c}} \sum_{i=1}^{N+M} (1 - y_i \mathbf{c}^T \mathbf{K}(X_i))^2 + \lambda \|\hat{\Phi}(X)\|^2$$

Η νόρμα στον τελευταίο όρο μπορεί να γραφεί ως:

$$\|\hat{\Phi}(X)\|^2 = \langle \hat{\Phi}(X), \hat{\Phi}(X) \rangle = \sum_{i=1}^{N+M} \sum_{j=1}^{N+M} c_i c_j K(X_i, X_j) = \mathbf{c}^T \mathbf{A} \mathbf{c}$$

Όπου \mathbf{A} ο $(N+M) \times (N+M)$ συμμετρικός πίνακας με στοιχεία $A_{i,j} = K(X_i, X_j)$

Οπότε το κριτήριο προς ελαχιστοποίηση είναι:

$$\min_{\mathbf{c}} L(\mathbf{c}) = \min_{\mathbf{c}} \sum_{i=1}^{N+M} (1 - y_i \mathbf{c}^T \mathbf{K}(X_i))^2 + \lambda \mathbf{c}^T \mathbf{A} \mathbf{c}$$

Οπότε για να βρούμε το βέλτιστο διάνυσμα συντελεστών \mathbf{c} , αφού το κριτήριο $L(\mathbf{c})$ είναι τετραγωνικό ως προς το \mathbf{c} αρκεί να λύσουμε την εξίσωση

$$\frac{\partial L(\mathbf{c})}{\partial \mathbf{c}} = \mathbf{0}$$

Οπότε έχουμε:

$$\begin{aligned} \frac{\partial L(\mathbf{c})}{\partial \mathbf{c}} = \mathbf{0} &\Rightarrow -2 \sum_{i=1}^{N+M} (1 - y_i \mathbf{c}^T \mathbf{K}(X_i)) y_i \mathbf{K}(X_i)^T + 2\lambda \mathbf{c}^T \mathbf{A} = \mathbf{0} \\ &\Rightarrow - \sum_{i=1}^{N+M} (y_i \mathbf{K}(X_i)^T) + \sum_{i=1}^{N+M} (y_i^2 \mathbf{c}^T \mathbf{K}(X_i) \mathbf{K}(X_i)^T) + \lambda \mathbf{c}^T \mathbf{A} = \mathbf{0} \end{aligned}$$

$$\begin{aligned}
&\Rightarrow \mathbf{c}^T \sum_{i=1}^{N+M} (\mathbf{K}(\mathbf{X}_i) \mathbf{K}(\mathbf{X}_i)^T) + \lambda \mathbf{c}^T \mathbf{\Lambda} = \sum_{i=1}^{N+M} (y_i \mathbf{K}(\mathbf{X}_i)^T) \\
&\Rightarrow \mathbf{c}^T (\mathbf{\Lambda}^2 + \lambda \mathbf{\Lambda}) = \sum_{i=1}^{N+M} (y_i \mathbf{K}(\mathbf{X}_i)^T) \\
&\Rightarrow \mathbf{c}^T = \sum_{i=1}^{N+M} (y_i \mathbf{K}(\mathbf{X}_i)^T) (\mathbf{\Lambda}^2 + \lambda \mathbf{\Lambda})^{-1} \quad (4)
\end{aligned}$$

Ερώτημα (δ)

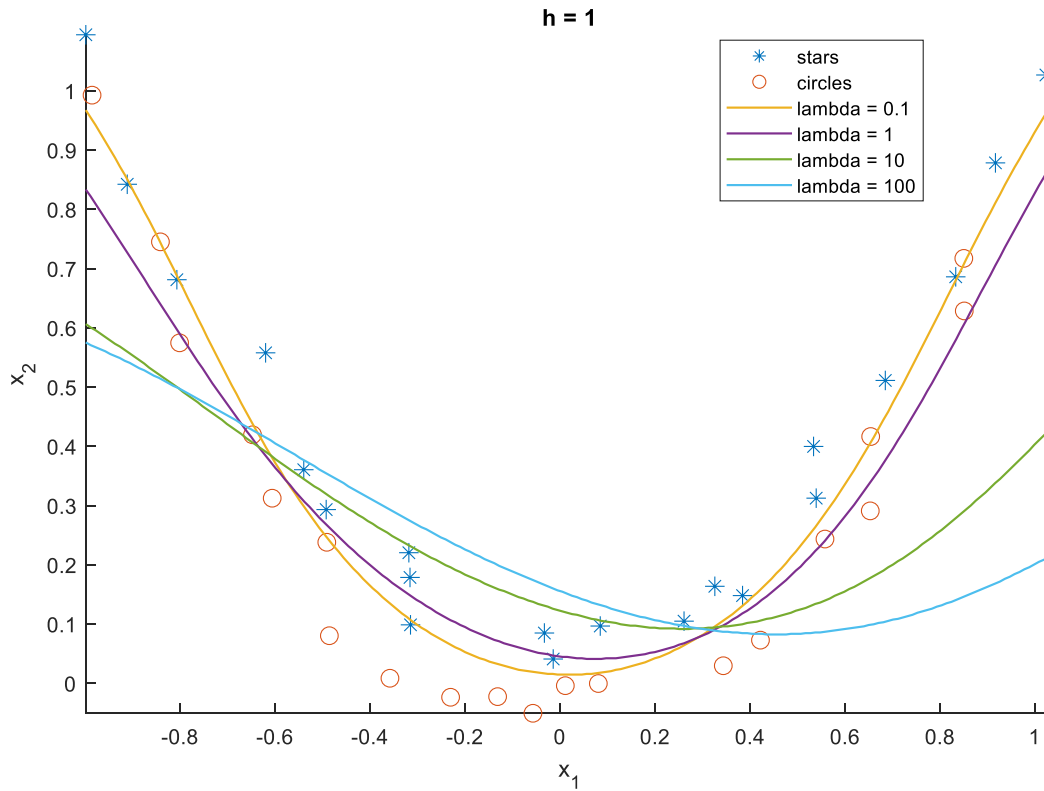
Αφού προσδιορίσαμε μέσω της εξίσωσης (4) τους βέλτιστους συντελεστές, έχουμε ορίσει πλήρως την συνάρτηση $\hat{\Phi}(\mathbf{X})$ η οποία μας δίνει μια προσέγγιση του ιδανικού κατηγοριοποιητή $\Phi(\mathbf{X})$. Συνεπώς για να κατηγοριοποιήσουμε ένα νέο διάνυσμα δεδομένων \mathbf{X}_{new} αρκεί να κοιτάξουμε την τιμή της $\hat{\Phi}(\mathbf{X}_{new})$ και να δούμε εάν είναι κοντινότερα στο 1 ή στο -1, δηλαδή ο κανόνας κατηγοριοποίησης θα είναι:

$$class\{\mathbf{X}_{new}\} = sign(\hat{\Phi}(\mathbf{X}_{new}))$$

Και το διαχωριστικό σύνορο φυσικά θα δίνεται από την εξίσωση $\hat{\Phi}(\mathbf{X}) = 0$.

Ερώτημα (ε)

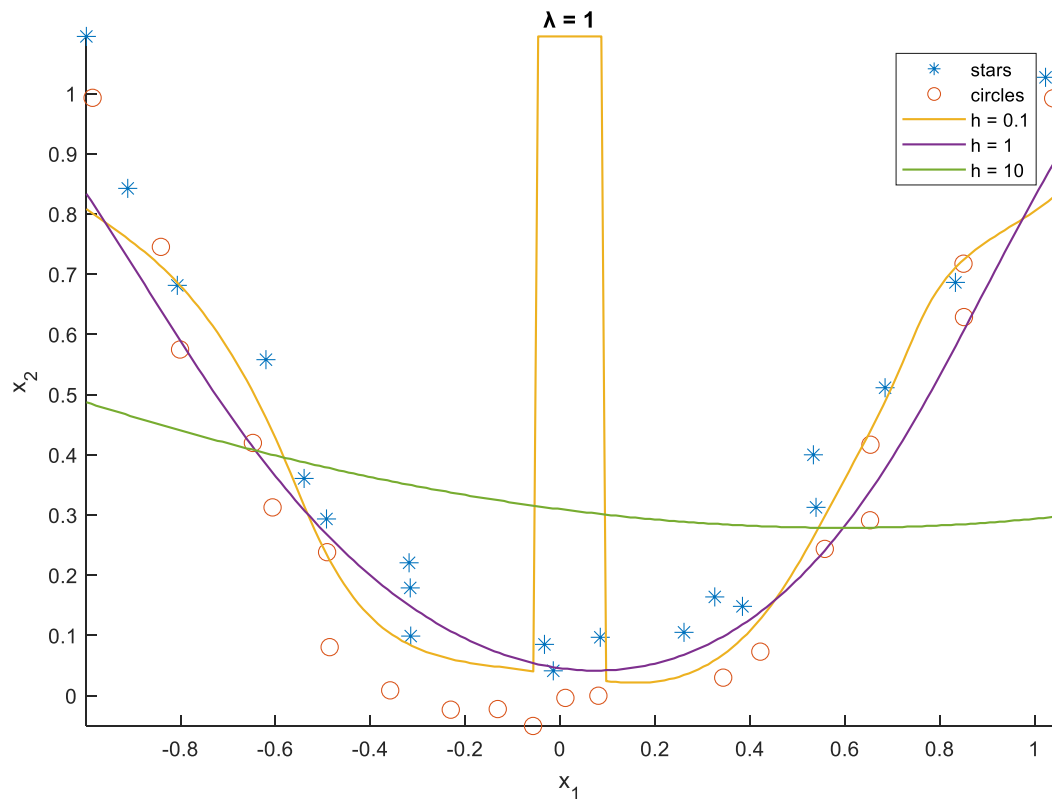
Υπολογίζοντας αριθμητικά το όριο διαχωρισμού για $h=1$ Και διάφορες τιμές της παραμέτρου λ έχουμε:



Σχήμα 3: Διαχωριστικό σύνορο μεταβάλλοντας την παράμετρο λ

Μπορούμε να παρατηρήσουμε πως όσο μικρότερο είναι το λ τόσο πιο κοντά στα δεδομένα εκπαίδευσης βρίσκεται το διαχωριστικό σύνορο, και εάν το μειώσουμε επαρκώς μπορεί να διαχωρίσει πλήρως τα δεδομένα αυτά. Ωστόσο αυτό δεν είναι επιθυμητό διότι θέλουμε ο κατηγοριοποιητής να μπορεί να διακρίνει και δεδομένα πέραν αυτών που χρησιμοποιήθηκαν για την εκπαίδευση. Το φαινόμενο αυτό ονομάζεται overfitting και πρέπει να αποφεύγεται όταν σχεδιάζουμε κατηγοριοποιητές.

Κρατώντας τώρα σταθερή την παράμετρο $\lambda=1$ και μεταβάλλοντας το h παίρνουμε τις παρακάτω καμπύλες:



Σχήμα 4: Διαχωριστικό σύνορο για μεταβλητό h

Μικραίνοντας το h φαίνεται η καμπύλη να μετατοπίζεται όλο και πιο κοντά στα δεδομένα εκπαίδευσης, αλλά εμφανίζονται περίεργες συμπεριφορές για μικρά h , όπως φαίνεται στην κίτρινη καμπύλη γύρω από το 0. Αυτό μπορεί ενδεχομένως να οφείλεται σε ανακρίβειες κατά την αντιστροφή του πίνακα για τον υπολογισμό των βέλτιστων συντελεστών.