

# Προβλεπτική Αναλυτική

Μεταπτυχιακή απαλλακτική εργασία ΙΟΥΝΙΟΥ 2024

Καθηγητής:  
Μ. Φιλιππάκης

Μαρκουλής Δημήτριος  
Καμπάσης Παναγιώτης  
Πανεπιστήμιο Πειραιώς  
Τμήμα Ψηφιακών Συστημάτων  
Μεγάλα Δεδομένα και Αναλυτική  
25 Ιουνίου 2024

# Προβλεπτική Αναλυτική

Πάνος Καμπάσης  
27/05/2024

## Περιεχόμενα

<b>I</b>	<i>Linear Regression</i>	2
I-A'	<i>Simple linear Regression</i> . . . . .	2
I-B'	<i>Linear Regression με Stochastic Gradient Descent (SGD)</i> . . . . .	3
<b>II</b>	<i>Multivariate Linear Regression</i>	4
II-A'	<i>Multivariate Linear Regression με Stochastic Gradient Descent</i> . . . . .	4
<b>III</b>	<i>Classification</i>	5
III-A'	<i>Logistic Regression</i> . . . . .	5
III-B'	<i>Logistic Regression με την μέθοδο Newton</i> . . . . .	6
<b>IV</b>	<i>Unsupervised Learning</i>	7
IV-A'	<i>K – meansClustering</i> . . . . .	7

## I. Linear Regression

### Task 1

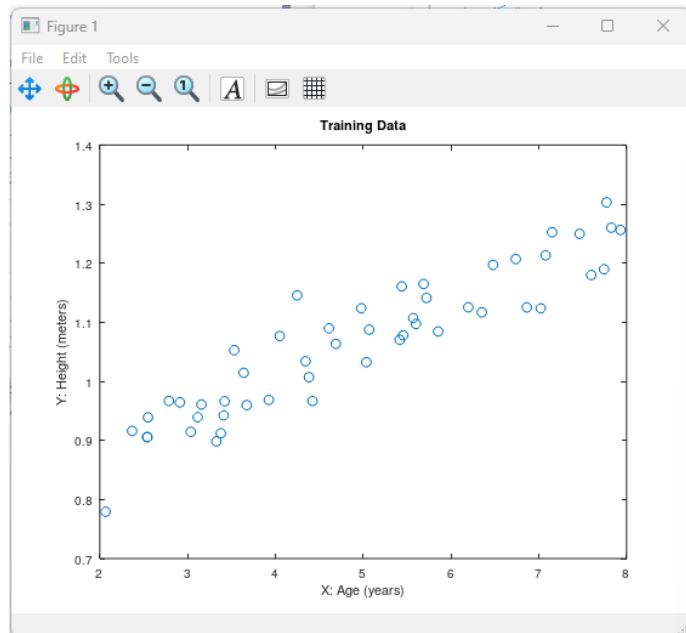
#### A'. Simple linear Regression

Χρησιμοποιώντας το κατάλληλο στατιστικό πακέτο της *Octave* υλοποιούμε το μοντέλο της απλής γραμμικής παλινδρόμησης για περιβάλλοντα τύπου *Matlab* και εξετάζουμε τα αποτελέσματα, τους χρόνους σύγκλισης και τα ποιοτικά χαρακτηριστικά της εκάστοτε μεθόδου σε συνδυασμό με το ποιόν των αρχικών μας δεδομένων. Οι μέθοδοι αυτοί έχουν ενίοτε την δυνατότητα να προσδώσουν προστιθέμενη πληροφορία στα δεδομένα μας και να εκμαιεύσουν χαρακτηριστικά που πρότινος δεν ήταν διαυγή.

Φορτώνουμε το στατιστικό πακέτο της *Octave* με: *pkgload statistics* (που είναι η αντίστοιχη βιβλιοθήκη της *MATLAB's Statistics and Machine Learning Toolbox*).

Για δεδομένα  $x, y$ :

```
9 clear all; close all; clc;
10 x = load('linear_regressionx.dat');
11 y = load('linear_regressiony.dat');
12
13 m = length(y); % number of data points = 50
14
15 % Add a column of ones to x for the intercept term
16 X = [ones(m, 1), x];
17
18 % Fit the linear regression model using regress
19 theta = regress(y, X);
20
21 % Predict y values using the model
22 y_pred = X * theta;
23
24 % Visualize the training data and the regression line
25 figure;
26 plot(x, y, 'o'); % Plot the data points
27 hold on;
28 plot(x, y_pred, '-r'); % Plot the regression line
29 ylabel('Y: Height (meters)')
30 xlabel('X: Age (years)')
31 title('Simple Linear Regression using regress(x, y)')
32 legend('Training data', 'Linear regression')
33 hold off;
34
35 % Display the theta values
36 disp('Using regress function:');
37 disp('Coefficients (Intercept and Slope):');
38 disp(theta);
```



where :

$y$  : Dependent variable (response)

$x$  : Independent variable (predictor)

$\theta_1$  : Intercept

$\theta_2$  : Slope of the regression line

$\epsilon$  : Error term

*Estimated Regression Equation*

The estimated regression equation is given by:

$$\hat{y} = \hat{\theta}_1 + \hat{\theta}_2 x \quad (2)$$

where:

$\hat{y}$  : Predicted value of  $y$

$\hat{\theta}_0$  : Estimated intercept

$\hat{\theta}_1$  : Estimated slope

Και ο αντίστοιχος γράφος:

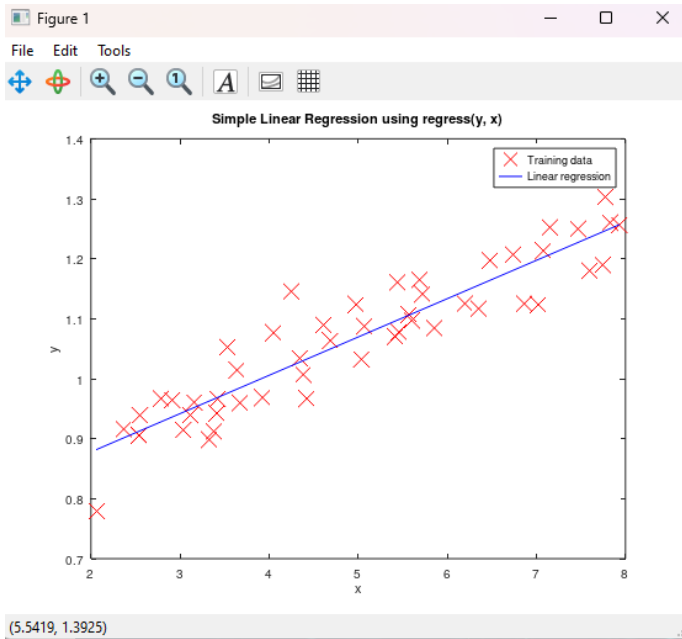
Χρησιμοποιώντας τον παρακάτω κώδικα και εκτέλωντας τον στην *Octave* παίρνουμε τα εξής αποτελέσματα. Επομένως το αποτέλεσμα μας είναι: *Theta values*:

$$\hat{\theta}_1 = 0.750163$$

$$\hat{\theta}_2 = 0.063881$$

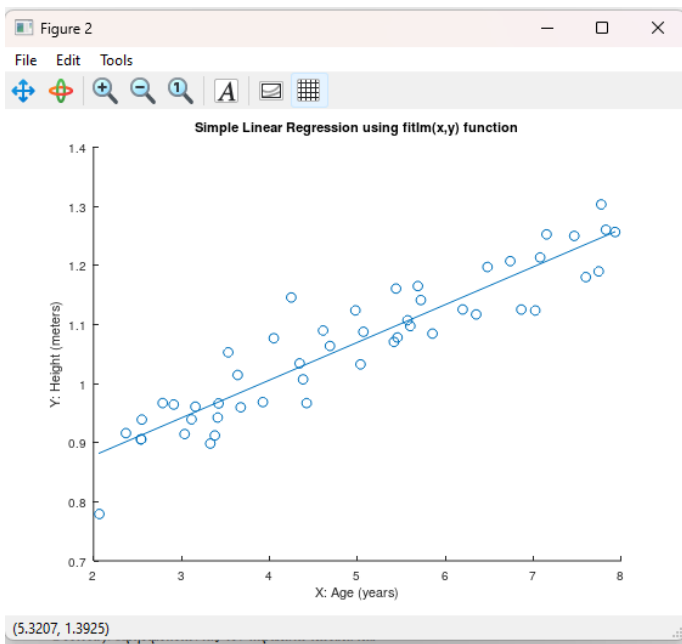
Επομένως η ευθεία μας προσαρμόζεται ως:

$$y = \theta_1 + \theta_2 x + \epsilon \quad (1)$$



Σχήμα 1: Γράφημα ΑΓΠ.

Συγκρίνουμε το αποτέλεσμα μας με εκείνο του αρχείου *simple\_lineaerregression.m* και παρατηρούμε:



Σχήμα 2: Γράφημα ΑΓΠ.

$$\hat{\beta}_0 = 0.750163$$

$$\hat{\beta}_1 = 0.063881$$

Βλέπουμε ότι τα αποτελέσματα είναι ίδια καθώς έχουμε ακριβώς τις ίδιες εκτιμήσεις  $\theta$ . Η μια προέκυψε από την *regress(y, X)* ενώ η άλλη από *fitlm(x, y)*

## Task 2

### B'. Linear Regression με Stochastic Gradient Descent (SGD)

Ομοίως ανακαλούμε την παραπάνω υλοποίηση προσθέτοντας μια άλλη μέθοδο σύγκλισης για την εύρεση της καλύτερης δυνατής κλίσης. Η συνάρτησή μας είναι:

$$h_{\theta}(x) = \theta^T x = \sum_{i=0}^n \theta_i x_i$$

Και με κανόνα ανενέωσης για να βρούμε το κατάλληλο  $\theta$ :

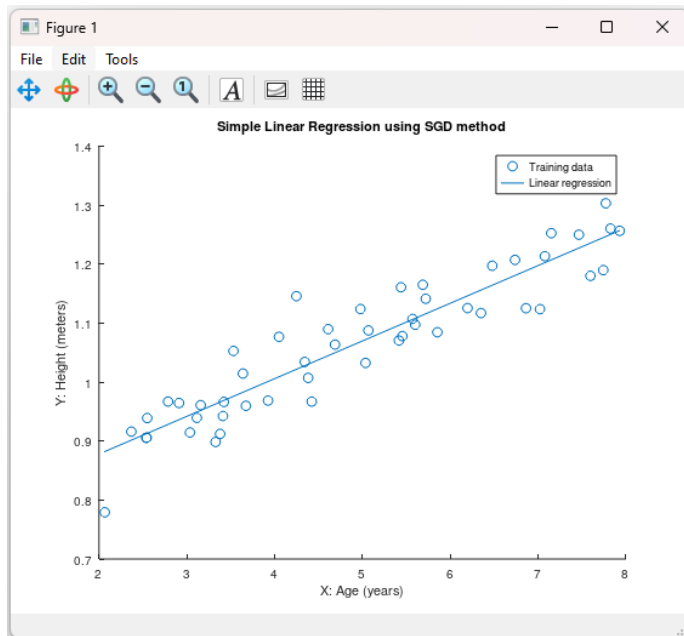
$$\theta_i := \theta_i - \alpha \frac{1}{m} \sum_{t=1}^m (h_{\theta}(x^{(t)}) - y^{(t)}) x_j^{(t)}$$

Χρησιμοποιούμε τον κώδικα *Task1.m*:

```
58 alpha = 0.07; % Learning rate
59 iterations = 1500; % Number of iterations
60 theta = zeros(2, 1); % Initialize fitting parameters
61
62 % Add a column of ones to x for the intercept term
63 X = [ones(m, 1), x];
64
65 for iter = 1:iterations
66     gradient = (1/m) * X' * (X * theta - y); % Compute the gradient
67     theta = theta - alpha * gradient; % Update theta
68 end
69
70 % Visualize results
71 figure;
72 hold on;
73 scatter(x, y);
74 plot(x, X * theta, '-');
75 ylabel('Y: Height (meters)');
76 xlabel('X: Age (years)');
77 title('Simple Linear Regression using SGD method');
78 legend('Training data', 'Linear regression');
79 hold off;
80
81 % Display theta values
82 disp('Using Stochastic Gradient Descent:');
83 disp('Coefficients (Intercept and Slope):');
84 disp(theta);
```

Σχήμα 3: Simple linear Regression with Stochastic Gradient Descent

Έχουμε το παρακάτω γράφημα που είναι σχεδόν πανομοιότυπο με τις προηγούμενες μεθόδους: με  $\theta$ :



Σχήμα 4: Simple linear Regression with Stochastic Gradient Descent

$\begin{bmatrix} 0.750150 \\ 0.063883 \end{bmatrix}$

## II. Multivariate Linear Regression

### Task3

#### A'. Multivariate Linear Regression με Stochastic Gradient Descent

Χρησιμοποιώντας νέα δεδομένα μεγαλύτερων διαστάσεων και τον κώδικα *Matlab*

*multivariate\_linear\_regression.m* Υλοποιούμε την παραπάνω διαδικασία για μια ακόμα φορά σε μεγαλύτερες διαστάσεις (γεγονός που βοηθά την εκμάθηση περισσότερης πληροφορίας καθώς στις πολλαπλές διαστάσεις εύκολα χάνεται η αντιληπτική ικανότητα) και ρυθμίζουμε τις διάφορες τιμές του  $\alpha$ .

Θα χρησιμοποιήσουμε το παρακάτω εύρος τιμών για να ερευνήσουμε την συμπεριφορά της σύγκλισης

$$\alpha = [0.01, 0.03, 0.1, 0.3, 1, 1.3]$$

Τότε ο κώδικας *Task3.m* μας δίνει:

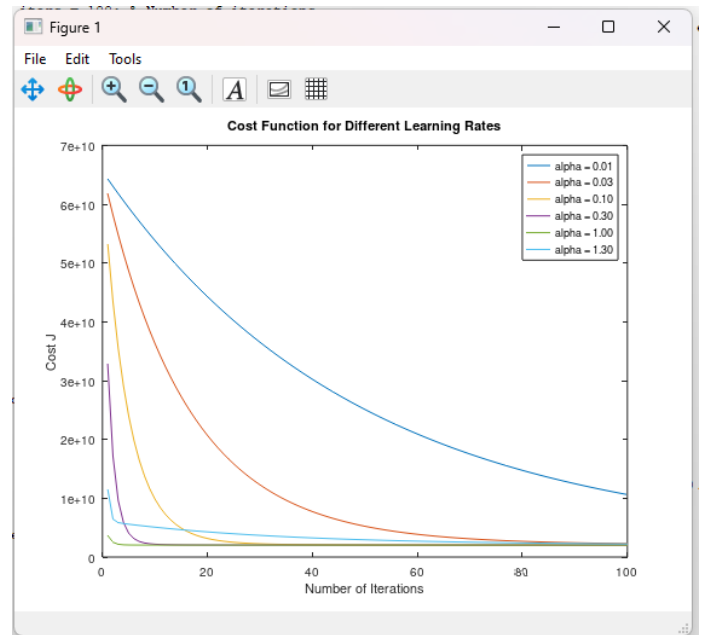
Καταγράφουμε επίσης και τις παρακάτω τιμές για τις διάφορες περιπτώσεις:

Σύγκριση Ρυθμού Μάθησης και Ανάλυση Σύγκλισης

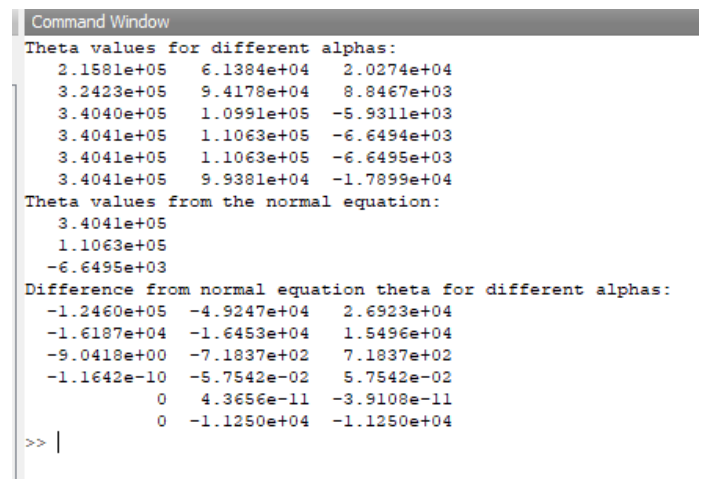
Άλφα = 0.01

- Διαφορά:  $[-1.2460 \times 10^5, -4.9247 \times 10^4, 2.6923 \times 10^4]$

Οι διαφορές είναι αρκετά μεγάλες, υποδηλώνοντας ότι ο ρυθμός μάθησης είναι πολύ μικρός και ο αλγόριθμος δεν έχει συγκλίνει καλά στις βέλτιστες τιμές  $\theta$  ήτα.



Σχήμα 5: Γράφημα σύγκλισης το  $\sigma$



Σχήμα 6: Simple linear Regression with Stochastic Gradient Descent

Άλφα = 0.03

- Διαφορά:  $[-1.6187 \times 10^4, -1.6453 \times 10^4, 1.5496 \times 10^4]$

Οι διαφορές είναι μικρότερες από ότι για  $\alpha = 0.01$ , αλλά ακόμα σημαντικές, υποδηλώνοντας καλύτερη σύγκλιση αλλά ακόμα όχι βέλτιστη.

Άλφα = 0.1

- Διαφορά:  $[-9.0418, -7.1837 \times 10^2, 7.1837 \times 10^2]$

Οι διαφορές είναι πολύ μικρότερες, υποδηλώνοντας ότι ο ρυθμός μάθησης είναι πιο κοντά στο βέλτιστο και ο αλγόριθμος έχει σχεδόν συγκλίνει στις σωστές τιμές  $\theta$  ήτα.

Άλφα = 0.3

- Διαφορά:  $[-1.1642 \times 10^{-10}, -5.7542 \times 10^{-2}, 5.7542 \times 10^{-2}]$

Οι διαφορές είναι εξαιρετικά μικρές, υποδηλώνοντας εξαιρετική σύγκλιση και σχεδόν τέλεια συμφωνία με τις τιμές θήτα της κανονικής εξίσωσης.

Άλφα = 1.0

- Διαφορά:  $[0, 4.3656 \times 10^{-11}, -3.9108 \times 10^{-11}]$

Οι διαφορές είναι ουσιαστικά μηδενικές, υποδηλώνοντας τέλεια σύγκλιση στις τιμές θήτα της κανονικής εξίσωσης.

Άλφα = 1.3

- Διαφορά:  $[0, -1.1250 \times 10^4, -1.1250 \times 10^4]$

Οι διαφορές είναι πάλι μεγαλύτερες, υποδηλώνοντας ότι ο ρυθμός μάθησης είναι πολύ υψηλός και προκαλεί απόκλιση ή υπέρβαση.

Από την παραπάνω ανάλυση και το γράφημα που μας υποδεικνύει τα 2 καλύτερα άλφα με γαλάζιο και πράσινο καταλήγουμε στο

Συμπέρασμα

- Καλύτερο Άλφα:  $\alpha = 1.0$
- Λόγος: Για  $\alpha = 1.0$ , οι διαφορές μεταξύ των τιμών θήτα του SGD και των τιμών θήτα της κανονικής εξίσωσης είναι ουσιαστικά μηδενικές. Αυτό υποδηλώνει τέλεια σύγκλιση και ότι ο ρυθμός μάθησης είναι βέλτιστος για αυτό το συγκεκριμένο πρόβλημα.

### III. Classification

#### Task 4

##### A'. Logistic Regression

Χρησιμοποιώντας νέα αριθμητικά δεδομένα με σκοπό την μελέτη σύγκλισης με την σιγμοειδή συνάρτηση έχουμε τις δύο κλάσεις μαθητών που πέτυχαν την είσοδο τους στο πανεπιστήμιο και εκείνων που δεν τα κατάφεραν εκείνη την χρονιά. Κατασκευάζουμε το αρχείο *logistic\_regression.m* και χρησιμοποιώντας την παρακάτω μορφή και SGD όπως παραπάνω, έχουμε:

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} = P(y = 1|x; \theta)$$

,και:

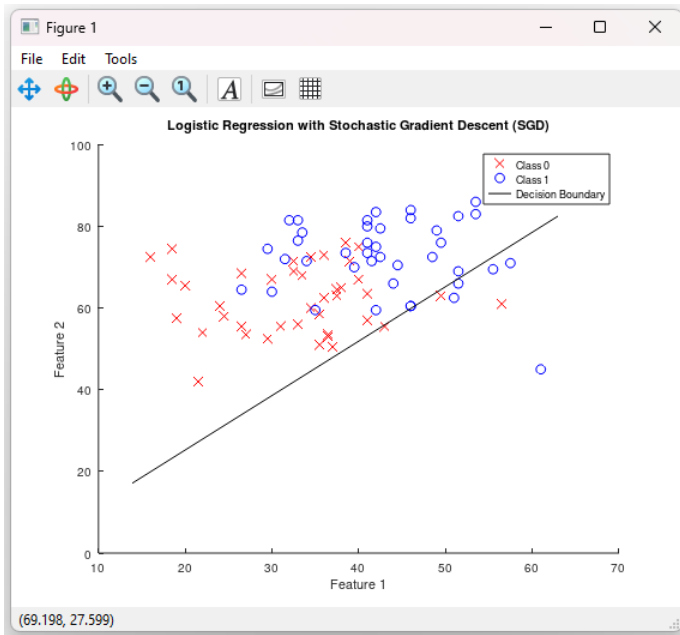
```

1 % Logistic Regression with Stochastic Gradient Descent (SGD)
2 % We will use log_regressionx.dat and log_regressiony.dat
3 % x = features
4 % y = labels (0 or 1)
5 pkg load statistics
6 clear all; close all; clc;
7
8 % Load data
9 x = load('log_regressionx.dat');
10 y = load('log_regressiony.dat');
11
12 % Initialize parameters
13 [m, n] = size(x);
14 x = [ones(m, 1), x]; % Add a column of ones to x for the intercept term
15 theta = zeros(n + 1, 1); % Initialize fitting parameters
16
17 alpha = 0.01; % Learning rate
18 iterations = 1500; % Number of iterations
19
20 % Sigmoid function
21 sigmoid = @(z) 1 ./ (1 + exp(-z));
22
23 % Gradient Descent
24 for iter = 1:iterations
25     % Compute hypothesis
26     h = sigmoid(x * theta);
27     % Compute gradient
28     gradient = (1/m) * x' * (h - y);
29     % Update theta
30     theta = theta - alpha * gradient;
31 end
32
33 % Visualize results
34 figure;
35 gscatter(x(:, 2), x(:, 3), y, 'rb', 'xo');
36 hold on;
37 x_values = [min(x(:, 2)) - 2, max(x(:, 2)) + 2];
38 y_values = -(theta(1) + theta(2) * x_values) / theta(3);
39 plot(x_values, y_values, 'k-');
40 xlabel('Feature 1');
41 ylabel('Feature 2');
42 title('Logistic Regression with Stochastic Gradient Descent (SGD)');
43 legend('Class 0', 'Class 1', 'Decision Boundary');
44 hold off;
45
46 % Display theta values
47 disp('Using Stochastic Gradient Descent:');
48 disp('Coefficients (including intercept):');
49 disp(theta);
50
51 % Prediction function
52 predict = @(x) round(sigmoid([ones(size(x, 1), 1) x] * theta));
53
54 % Compute accuracy
55 predictions = predict(x(:, 2:end));
56 accuracy = mean(double(predictions == y)) * 100;
57 disp(['Training Accuracy: ', num2str(accuracy), '%']);
58

```

Σχήμα 7: Logistic Regression with Stochastic Gradient Descent

Έχουμε τα εξής αποτελέσματα με την εκτέλεση του παραπάνω κώδικα:



Σχήμα 8: Graph Logistic Regression with Stochastic Gradient Descent

Το οποίο δίνει ακρίβεια 52.5%. Το ποσοστό αυτό είναι χαμηλό καθώς μας δείχνει ότι το μοντέλο είναι ελαφρώς καλύτερο στο να κατηγοριοποιεί από την κατηγοριοποίηση που θα μας έδινε ένα τίμιο κέρμα. Επομένως δεν θα μέναμε στην παρούσα υλοποίηση εάν θέλαμε να πετύχουμε αυστηρές κατηγοριοποιήσεις.

#### Task5

Β'. Logistic Regression με την μέθοδο Newton

Χρησιμοποιώντας την μέθοδο Newton με την συνάρτηση:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left( -y^{(i)} \log(h_{\theta}(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right) \quad (3)$$

$$\theta^{(t+1)} = \theta^{(t)} - H^{-1} \nabla_{\theta} J \quad (4)$$

όπου:

$$\nabla_{\theta} J = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \quad (5)$$

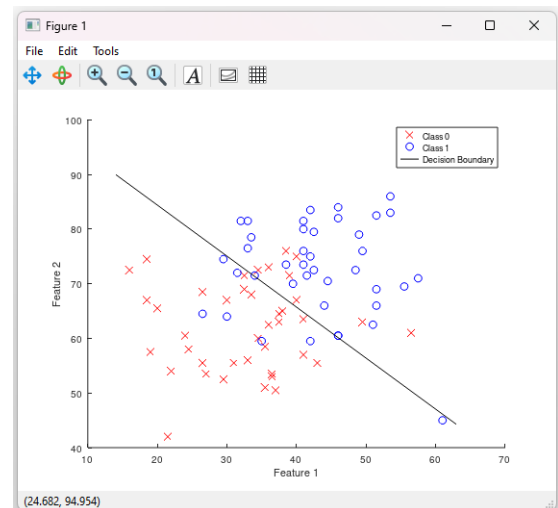
και

$$H = \frac{1}{m} \sum_{i=1}^m \left[ h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) x^{(i)} (x^{(i)})^T \right] \quad (6)$$

Ο κώδικας μας *Task5.m* Μας δίνει τα εξής αποτελέσματα:

```
1 % Logistic Regression with Newton's Method
2 % We will use log_regressionn.dat and log_regressionny.dat
3 % x = features
4 % y = labels (0 or 1)
5
6 clear all; close all; clc;
7
8 % Load data
9 x = load('log_regressionn.dat');
10 y = load('log_regressionny.dat');
11
12 % Initialize parameters
13 [m, n] = size(x);
14 x = [ones(m, 1), x]; % Add a column of ones to x for the intercept term
15 theta = zeros(n + 1, 1); % Initialize fitting parameters
16
17 % Sigmoid function
18 sigmoid = @(z) 1 ./ (1 + exp(-z));
19
20 % Newton's method
21 iterations = 15; % Number of iterations, usually converges within 5-15 iterations
22
23 for iter = 1:iterations
24     % Compute hypothesis
25     h = sigmoid(x * theta);
26
27     % Compute gradient
28     gradient = (1/m) * x' * (h - y);
29
30     % Compute Hessian
31     H = (1/m) * x' * diag(h .* (1 - h)) * x;
32
33     % Update theta
34     theta = theta - H \ gradient;
35 end
36
37 % Visualize results
38 figure;
39 scatter(x(:, 2), x(:, 3), y, 'xb', 'xo');
40 hold on;
41 x_values = [min(x(:, 2)) - 2, max(x(:, 2)) + 2];
42 y_values = -(theta(1) + theta(2) * x_values) / theta(3);
43 plot(x_values, y_values, 'k-');
44 xlabel('Feature 1');
45 ylabel('Feature 2');
46 title('Logistic Regression with Newton's Method');
47 legend('Class 0', 'Class 1', 'Decision Boundary');
48 hold off;
49
50 % Display theta values
51 disp('Using Newton's Method:');
52 disp('Coefficients (including intercept):');
53 disp(theta);
54
55 % Prediction function
56 predict = @(x) round(sigmoid([ones(size(x, 1), 1) x] * theta));
57
58 % Compute accuracy
59 predictions = predict(x(:, 2:end));
60 accuracy = mean(double(predictions == y)) * 100;
61 disp(['Training Accuracy: ', num2str(accuracy), '%']);
```

Σχήμα 9: Logistic Regression with Newton's Method



Σχήμα 10: Logistic Regression with Newton's Method Graph

- Ακρίβεια Εκπαίδευσης: 81.25
- Έλεγχος Σύγκλισης : Σύγκλιση σε 7 επαναλήψεις.

- Τιμές Θήτα  
-16.3787  
0.1483  
0.1589

Πράγματι επιτεύχθηκε η σύγκλιση μέσα στο προτεινόμενο όριο των 5 έως 15 επαναλήψεων και μάλιστα πολύ γρήγορα με μόλις 7. Επίσης πολύ καλύτερο ποσοστό ακρίβειας και ακόμα και από την εικόνα βλέπουμε ότι η ευθεία διαχωρίζει πολύ καλύτερα τις δύο κλάσεις, ασφαλώς με ένα ποσοστό σφάλματος.

#### IV. Unsupervised Learning

##### Task 6

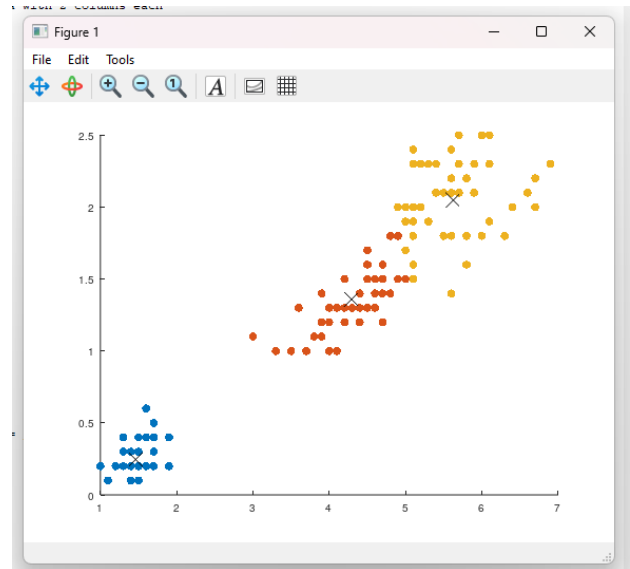
##### A'. *K – meansClustering*

Με την χρήση των αριθμητικών δεδομένων στα οποία θα έχουμε 150 δείγματα και 2 χαρακτηριστικά το καθένα θα δημιουργήσουμε ένα μοντέλο *K-means* για διαφορετικό αριθμό πυρήνων έτσι με την κατασκευή του παρακάτω αλγορίθμου έχουμε:

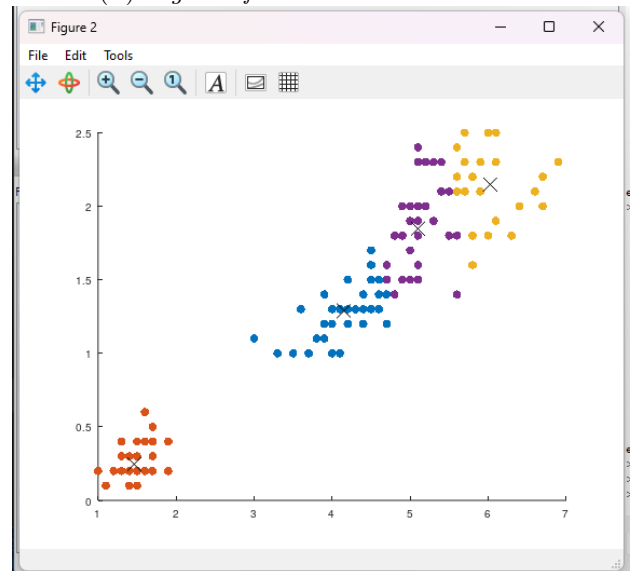
```
1 function k_means_clustering(k)
2   % Load the data
3   data = load('kmeans.dat'); % Assuming kmeans.dat contains 150 rows of data with 2 columns each
4
5   % Number of data points
6   num_points = size(data, 1);
7
8   % Randomly initialize the cluster centers
9   centers = data(randperm(num_points, k), :);
10
11  % To store the cluster assignments
12  cluster_assignments = zeros(num_points, 1);
13
14  % To store the previous cluster centers for convergence check
15  previous_centers = zeros(size(centers));
16
17  % Iterate until convergence
18  while true
19    % Assign each point to the nearest center
20    for i = 1:num_points
21      distance = sum((centers - data(i, :)).^2, 2);
22      [~, cluster_assignments(i)] = min(distance);
23    end
24
25    % Recompute the cluster centers
26    for j = 1:k
27      cluster_assignments == j;
28      centers(j, :) = mean(data(cluster_assignments == j, :), 2);
29    end
30
31    % Check for convergence (if centers do not change)
32    if isequal(centers, previous_centers)
33      break;
34    end
35
36    % Update previous centers
37    previous_centers = centers;
38  end
39
40  % Plot the data points and the cluster centers
41  figure;
42  hold on;
43  colors = lines(k);
44  scatter(data(cluster_assignments == 1, :), data(cluster_assignments == 1, :), colors(1, :), 'filled');
45  end
46  scatter(centers(1, :), centers(1, :), 100, 'k', 'x');
47  hold off;
48
49  % Save the cluster centers
50  save('k_means_centers.mat', 'centers');
51 end
```

Σχήμα 11: *Logistic Regression with Newton's Method Graph*

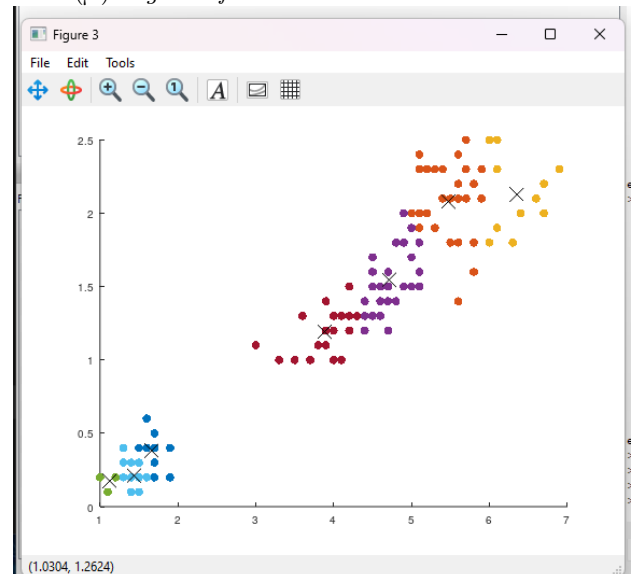
Και αφού τον εκτελέσουμε μπορούμε να εκτελέσουμε *k\_means\_clustering(k)*, όπου  $k = 2, 3, 4, \dots$  και να λάβουμε το αντίστοιχο γράφημα. Δοκιμάζουμε για  $k = 3, 4$  και 7.



(α') Figure of K-means with 3 clusters



(β') Figure of K-means with 4 clusters

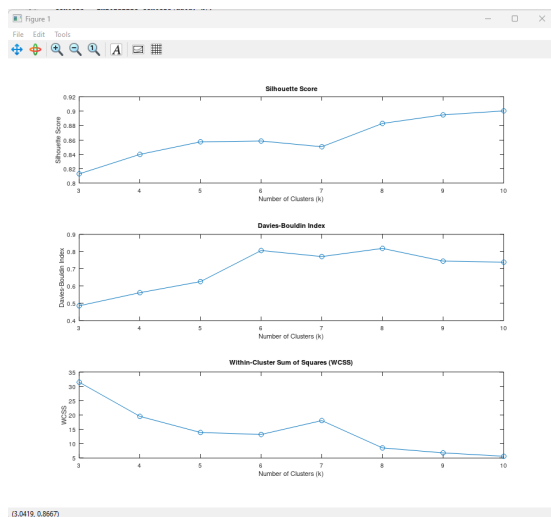


(γ') Figure of K-means with 7 clusters

Σχήμα 12: *K\_means*



Συγκρίνοντας οπτικά προτιμούμε τους λιγότερους πυρήνες καθώς βλέπουμε ότι στους 7 πυρήνες μια περιοχή που φαίνεται ξεκάθαρα αποκομμένη από τα υπόλοιπα δεδομένα μπορεί εύκολα να αναπαρασταθεί από ένα κέντρο όπως στις πρώτες 2 εικόνες. Για να έχουμε μια καλύτερη εικόνα αναβαθμίζουμε τον κώδικα μας ώστε να βρούμε τον καλύτερο αριθμό πυρήνων. Υλοποιούμε τον *KmeansUltraPlus.m* και εκτελούμε *evaluate\_k\_means()*. Μπορούμε ευκολά να δώσουμε στον κώδικα να ελέγξει περισσότερους πυρήνες αλλά



Σχήμα 13: *Silhouette's*, *Davies – Bouldin Index*, *WSS* for  $K$  clusters where  $\kappa = [1,10]$

Σύμφωνα με τα γραφήματα που φαίνονται στην εικόνα, ο καλύτερος αριθμός των κλάσεων  $k$  για τον αλγόριθμο  $k$  – means μπορεί να εκτιμηθεί ως εξής:

- *SilhouetteScore*: Ο δείκτης *Silhouette* αυξάνεται σταθερά μέχρι το  $k = 6$  και μετά παραμένει περίπου σταθερός με μικρές διακυμάνσεις. Υψηλότερη τιμή υποδεικνύει καλύτερο *clustering*.
- *Davies – BouldinIndex*: Ο δείκτης *Davies – Bouldin* είναι χαμηλότερος για  $k = 5$ , το οποίο υποδεικνύει καλύτερη κατανομή των κλάσεων καθώς χαμηλότερη τιμή υποδεικνύει καλύτερο *clustering*.
- *Within – Cluster Sum of Squares (WS)*: Η καμπύλη του *WSS* παρουσιάζει μεγάλη μείωση μέχρι το  $k = 5$  και μετά από εκεί η μείωση είναι μικρότερη. Το σημείο όπου η καμπύλη αρχίζει να γίνεται πιο επίπεδη, γνωστό και ως "*Elbow*", μπορεί να υποδεικνύει τον καλύτερο αριθμό των κλάσεων.

Συμπερασματικά, ο καλύτερος αριθμός των κλάσεων  $k$  φαίνεται να είναι είτε  $k = 5$  είτε  $k = 6$ . Ωστόσο, λαμβάνοντας υπόψη τον *Davies – Bouldin Index*, το  $k = 5$  φαίνεται να είναι η καλύτερη επιλογή.