



Data Preprocessing στην R

1 Εισαγωγή

1.1 Εισαγωγή στο Data Preprocessing

Η διαδικασία της προεπεξεργασίας δεδομένων αποτελεί ίσως το πιο σημαντικό βήμα στην εξόρυξη δεδομένων. Τα σύνολα δεδομένων που χρησιμοποιούνται μπορεί συχνά να περιέχουν θόρυβο, περιττές πληροφορίες κ.α. ενώ πολλές φορές και η μορφή τους δεν είναι κατάλληλη για την εφαρμογή αλγορίθμων. Η προεπεξεργασία που γίνεται στα δεδομένα αφορά σε πλήθος μεθόδων που αφορούν τον καθαρισμό των δεδομένων (Data Cleansing) (π.χ. Fill in missing values, smooth noisy data, identify or remove outliers and noisy data, and resolve inconsistencies), τη δειγματοληψία τους (Data Sampling), το μετασχηματισμό τους π.χ. κανονικοποίηση (Data Normalization), την επιλογή και την εξαγωγή χαρακτηριστικών (Feature Selection and Extraction) κ.α. Τέλος, όσον αφορά την εξαγωγή χαρακτηριστικών, διακρίνουμε τις μεθόδους σε γραμμικές και μη γραμμικές, όπως PCA και ISOMAP αντίστοιχα.

1.2 Data Preprocessing στην R

Για αν αφαιρέσουμε τις διπλοεγγραφές από τα δεδομένα εκτελούμε την εντολή:

```
> data = unique(data)
```

Για απλούς μετασχηματισμούς χρησιμοποιούμε τη συνάρτηση scale. Η scale έχει τα ορίσματα center και scale. Με το center = TRUE αφαιρεί από τα δεδομένα το μέσο όρο τους (δηλαδή τα κεντράρει στην αρχή των αξόνων). Με το παράμετρο scale = TRUE διαιρεί τα δεδομένα με την τυπική τους απόκλιση, όπως παρακάτω:

```
> transformed <- scale(data, center = TRUE, scale = TRUE)
```

Νέα δεδομένα μπορούν και αυτά να γίνουν scale, αρκεί ως ορίσματα να δοθούν οι υπάρχοντες μέσοι όροι και τυπικές αποκλίσεις:

```
> scale(newdata, center = attr(transformed, "scaled:center"),  
      scale = attr(transformed, "scaled:scale"))
```

Για να κάνουμε normalize ένα διάνυσμα x στο εύρος [0, 1] εκτελούμε την εντολή:

```
> normalized_x = (x - min(x)) / (max(x) - min(x))
```

Για να εφαρμόσουμε το normalization σε ένα data frame εκτελούμε την εντολή:

```
> normalized = as.data.frame(lapply(pdata, function(x)  
      (x - min(x)) / (max(x) - min(x))))
```

Για να διακριτοποιήσουμε τα δεδομένα χρησιμοποιούμε την `cut`, η οποία λαμβάνει ως όρισμα μια ακολουθία από τα σημεία στα οποία θα γίνει διακριτοποίηση, π.χ. η εντολή:

```
> cut(x, seq(0,100,10))
```

διακριτοποιεί τα δεδομένα στα διαστήματα (0, 10], (10, 20], ..., (90, 100].

Για να εφαρμόσουμε δειγματοληψία στα δεδομένα χρησιμοποιούμε την εντολή `sample`. Για να λάβουμε π.χ. 100 εγγραφές εκτελούμε την εντολή:

```
> data_sample = data[sample(nrow(data), 100, replace = TRUE),]
```

όπου με την παράμετρο `replace` ελέγχουμε αν η δειγματοληψία θα γίνει με αντικατάσταση.

Τέλος, για να υπολογίσουμε το `correlation` των δεδομένων μπορούμε να χρησιμοποιήσουμε την `cor`, ενώ για το `covariance` μπορούμε να χρησιμοποιήσουμε την `cov`.

1.2 PCA στην R

Για να υπολογίζουμε το PCA εκτελούμε τη συνάρτηση `prcomp`:

```
> pca_model <- prcomp(data, center = TRUE, scale = TRUE)
```

Όπου μπορούμε να επιλέξουμε να κάνουμε `center` και `scale` τα δεδομένα με τις αντίστοιχες παραμέτρους. Μπορούμε να λάβουμε τις ιδιοτιμές και τα ιδιοδιανύσματα με τις εντολές:

```
> eigenvalues = pca_model$sdev^2
```

```
> eigenvectors = pca_model$rotation
```

1.3 ISOMAP στην R

Για να υπολογίζουμε το ISOMAP χρησιμοποιούμε τη συνάρτηση `isomap` που βρίσκεται στη βιβλιοθήκη `vegan`. Εφόσον έχουμε φορτώσει τη βιβλιοθήκη (`library(vegan)`), αρχικά κατασκευάζουμε ένα πίνακα αποστάσεων για τα δεδομένα:

```
> data_dist <- dist(data)
```

και στη συνέχεια εφαρμόζουμε τον αλγόριθμο:

```
> isom <- isomap(data_dist, ndim = 2, k = 4)
```

όπου με την παράμετρο `ndim` επιλέγουμε τον αριθμό των διαστάσεων που θα μετασχηματίσουμε τα δεδομένα και με το `k` επιλέγουμε το πλήθος των κοντινότερων αποστάσεων για κάθε σημείο. Μπορούμε να πάρουμε τις τιμές στις νέες διαστάσεις ως εξής:

```
> data_2d <- isom$points
```

2 Προεπεξεργασία Δεδομένων – Χρήση Μεθόδων για Καθαρισμό, Κανονικοποίηση, Διακριτοποίηση και Δειγματοληψία

Για την προεπεξεργασία δεδομένων θα χρησιμοποιήσουμε τις δύο πρώτες στήλες των δεδομένων του αρχείου που δίνεται (engdata.txt).

Ένα summary των δεδομένων με την R είναι το παρακάτω:

Age	Salary
Min. :20.00	Min. : 475
1st Qu.:44.00	1st Qu.:1156
Median :51.00	Median :1517
Mean :49.94	Mean :1592
3rd Qu.:57.00	3rd Qu.:1969
Max. :70.00	Max. :3900

Εισάγουμε τα δεδομένα με τις παρακάτω εντολές:

```
> engdata = read.csv("engdata.txt")  
> pdata = engdata[, 1:2]
```

Θα απαντήσουμε στα παρακάτω ερωτήματα:

α) Διαγράψτε τις διπλοεγγραφές από τα δεδομένα.

β) Εφαρμόστε τους μετασχηματισμούς center και scale στα δεδομένα και σχεδιάστε τα δεδομένα πριν και μετά από αυτούς τους μετασχηματισμούς.

γ) Εφαρμόστε δειγματοληψία με αντικατάσταση στα δεδομένα επιλέγοντας 150 τιμές. Σχεδιάστε τα δεδομένα πριν και μετά τη δειγματοληψία. Διατηρείται η δομή των δεδομένων;

δ) Εφαρμόστε διακριτοποίηση στα δεδομένα και κάντε τα plot σε bar charts.

2.1 Καθαρισμός και Μετασχηματισμοί Δεδομένων

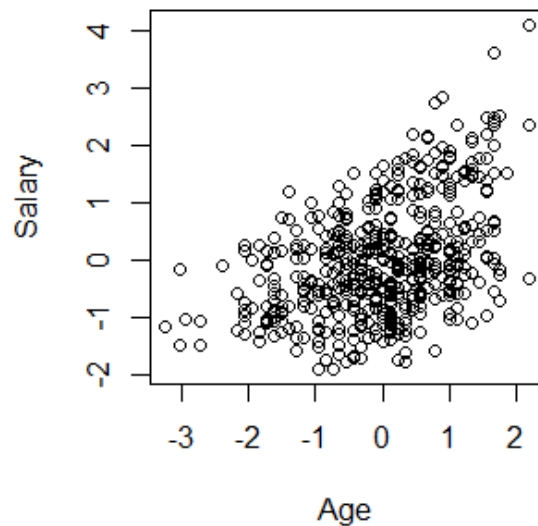
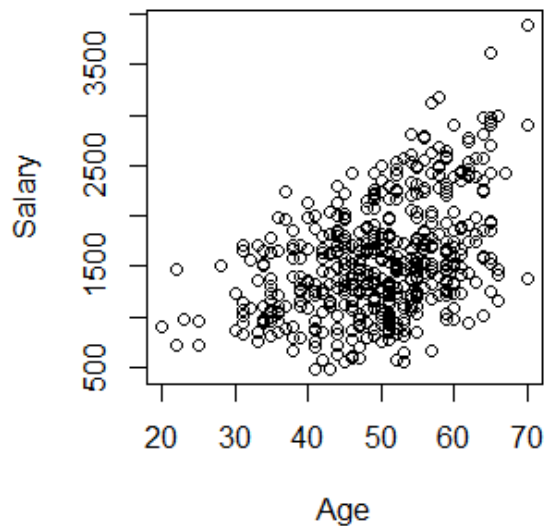
Αρχικά αφαιρούμε τα διπλότυπα από τα δεδομένα (ερώτημα (α)) με την παρακάτω εντολή:

```
> pdata = unique(pdata)
```

Στη συνέχεια, αφαιρούμε το μέσο όρο και διαιρούμε με την τυπική απόκλιση (ερώτημα (β)):

```
> transformed <- scale(pdata, center = TRUE, scale = TRUE)
```

Μπορούμε να σχεδιάσουμε τα δεδομένα πριν και μετά τους μετασχηματισμούς με τις εντολές `plot(pdata)` και `plot(transformed)`.

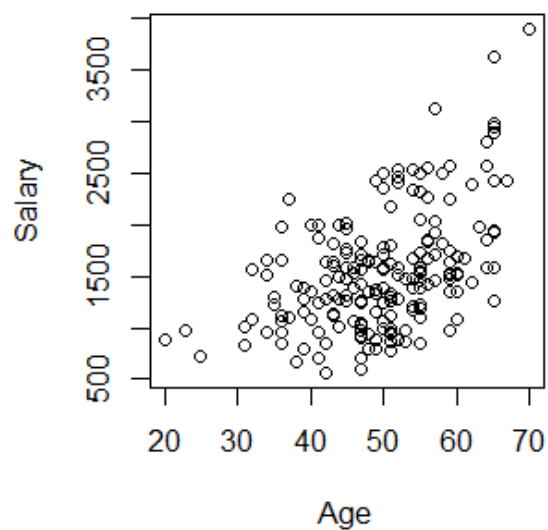
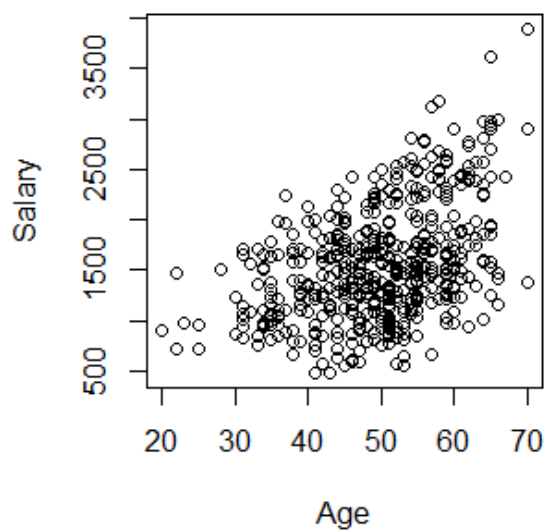


2.2 Δειγματοληψία Δεδομένων

Η δειγματοληψία (ερώτημα (γ)) γίνεται με την εντολή `sample`:

```
> sampdata = pdata[sample(nrow(pdata), 250, replace = TRUE),]
```

Μπορούμε να σχεδιάσουμε τα δεδομένα πριν και μετά τους μετασχηματισμούς με τις εντολές `plot(pdata)` και `plot(sampdata)`.



2.3 Διακριτοποίηση Δεδομένων

Η διακριτοποίηση (ερώτημα (γ)) γίνεται με την εντολή `cut`. Για τις μεταβλητές `Age` και `Salary` επιλέγουμε κατάλληλα `bins` με τις παρακάτω εντολές:

```
> discAge = cut(pdata$Age, seq(0,80,10))
```

```
> discSalary = cut(pdata$Salary, seq(0,4000,400), dig.lab = 4)
```

Μπορούμε να σχεδιάσουμε τα δεδομένα σε bar charts με τις εντολές `plot(discAge)` και `plot(discSalary, las=2)`.

3 Εξαγωγή Χαρακτηριστικών με PCA

Για το μετασχηματισμό ενός συνόλου δεδομένων με PCA θα χρησιμοποιήσουμε ως εφαρμογή τα δεδομένα του παρακάτω πίνακα.

	X	Y	Z
X1	1	0	-1
X2	0	1	-1
X3	-1	1	0
X4	0	-1	1
X5	-1	0	1
X6	1	-1	0

Θα απαντήσουμε στα παρακάτω ερωτήματα:

α) Σχεδιάστε τα δεδομένα σε 3 διαστάσεις.

β) Να εφαρμοστεί ο αλγόριθμος PCA για τα διανύσματα $x_1 - x_6$ και να αντικατασταθούν από διανύσματα 2 διαστάσεων κατά τέτοιον τρόπο ώστε η απώλεια πληροφορίας να είναι ελάχιστη.

γ) Επαναλάβετε το ερώτημα (α) χρησιμοποιώντας την R.

δ) Σχεδιάστε τα νέα δεδομένα μετά το μετασχηματισμό.

3.1 Κατασκευή Δεδομένων και Εισαγωγή Βιβλιοθηκών

Αρχικά κατασκευάζουμε τα δεδομένα και φορτώνουμε τις απαραίτητες βιβλιοθήκες:

```
> pdata = data.frame(X = c(1,0,-1,0,-1,1), Y = c(0,1,1,-1,0,-1),  
                     Z = c(-1,-1,0,1,1,0))  
  
> row.names(pdata) <- c("x1", "x2", "x3", "x4", "x5", "x6")  
  
> library(scatterplot3d)
```

3.2 Εφαρμογή Αλγορίθμου PCA

Για το ερώτημα (α) εκτελούμε:

```
> s3d = scatterplot3d(pdata, color = "blue", pch = 19, scale.y = 1.5)  
> coords <- s3d$xyz.convert(pdata)  
> text(coords$x, coords$y, labels=row.names(pdata), pos=2)
```

Για το ερώτημα (β), αρχικά ορίζουμε τους πίνακες X και $X^T X$:

$$X = [x_1 \ x_2 \ \dots \ x_n]^T = \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & -1 \\ -1 & 1 & 0 \\ 0 & -1 & 1 \\ -1 & 0 & 1 \\ 1 & -1 & 0 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 1 & 0 & -1 & 0 & -1 & 1 \\ 0 & 1 & 1 & -1 & 0 & -1 \\ -1 & -1 & 0 & 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & -1 \\ -1 & 1 & 0 \\ 0 & -1 & 1 \\ -1 & 0 & 1 \\ 1 & -1 & 0 \end{bmatrix} = \begin{bmatrix} 4 & -2 & -2 \\ -2 & 4 & -2 \\ -2 & -2 & 4 \end{bmatrix}$$

Στην συνέχεια απαιτείται ο υπολογισμός των ιδιοτιμών και των ιδιοδιανυσμάτων του πίνακα $X^T X$. Οι ιδιοτιμές του πίνακα $X^T X$ είναι οι ρίζες του χαρακτηριστικού του πολυνύμου. Ισχύει:

$$|X^T X - \lambda I_3| = 0 \Rightarrow \begin{vmatrix} 4-\lambda & -2 & -2 \\ -2 & 4-\lambda & -2 \\ -2 & -2 & 4-\lambda \end{vmatrix} = 0 \Rightarrow (\lambda - 6) \cdot \lambda \cdot (6 - \lambda) = 0$$

Άρα οι λύσεις του χαρακτηριστικού πολυνύμου είναι $\lambda_{1,2} = 6, \lambda_3 = 0$

Επόμενο βήμα αποτελεί η εύρεση των ιδιοδιανυσμάτων που αντιστοιχούν στις 2 μεγαλύτερες ιδιοτιμές αφού ζητείται η αναπαράσταση να γίνει στις 2 διαστάσεις. Για το 1^ο ιδιοδιάνυσμα \underline{u}_1 ισχύει:

$$(X^T X - \lambda_1 I_3) \underline{u}_1 = 0 \Rightarrow \begin{bmatrix} -2 & -2 & -2 \\ -2 & -2 & -2 \\ -2 & -2 & -2 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = 0 \Rightarrow x + y + z = 0 \Rightarrow z = -x - y \quad (1)$$

Η παραπάνω σχέση αποτελεί συνθήκη που πρέπει να ισχύει λόγω της 1^{ης} ιδιοτιμής λ_1 , για το 1^ο ιδιοδιάνυσμα. Οπότε προκειμένου να υπολογίσουμε το ιδιοδιάνυσμα που αντιστοιχεί στην ιδιοτιμή λ_1 έχουμε:

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = x \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + y \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} + z \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \stackrel{(4)}{\Rightarrow} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = x \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + y \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} - (x + y) \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} x \\ y \\ z \end{bmatrix} = x \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} + y \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix}$$

Η παραπάνω σχέση μας δίνει μια οικογένεια ιδιοδιανυσμάτων για τις διάφορες τιμές των x, y που όμως δεν είναι όλα πάντα γραμμικώς ανεξάρτητα και δεν έχουν μέτρο όσο με 1. Έστω επιλέγω ως \underline{u}_1 το ακόλουθο διάνυσμα:

$$\underline{u}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \quad (2)$$

Για το 2^ο ιδιοδιάνυσμα που επίσης αντιστοιχεί στην ιδιοτιμή $\lambda_1 = \lambda_2 = 6$ πρέπει να ισχύει:

$$\underline{u}_1 \perp \underline{u}_2 \Rightarrow \underline{u}_2^T \underline{u}_1 = 0 \Rightarrow \frac{1}{\sqrt{2}} \begin{bmatrix} x & y & z \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} = 0 \Rightarrow x - z = 0 \Rightarrow x = z \quad (3)$$

Με συνδυασμό των σχέσεων (1) και (3) προκύπτει ότι για το \underline{u}_2 πρέπει να ισχύει:

$$\xrightarrow{\text{από (1) και (3)}} y = -2x$$

Οπότε τελικά έχουμε:

$$\underline{u}_2 = \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} x \\ -2x \\ x \end{bmatrix} \xrightarrow{x=1} \underline{u}_2 = \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} \xrightarrow{\text{μοναδιαίο διάνυσμα}} \underline{u}_2 = \frac{1}{\sqrt{6}} \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} \quad (4)$$

Έχοντας υπολογίσει τα 2 ιδιοδιανύσματα που αντιστοιχούν στις 2 μεγαλύτερες ιδιοτιμές (σχέσεις (2) και (4)), η αναγωγή των 3D διανυσμάτων x_i σε αντίστοιχα 2D με την ελάχιστη απώλεια πληροφορίας πραγματοποιείται βάση της σχέσης:

$$x'_i = \begin{bmatrix} \underline{u}_1^T x_i \\ \underline{u}_2^T x_i \end{bmatrix}$$

Έτσι προκύπτουν τα νέα δισδιάστατα διανύσματα.

$$x'_1 = \begin{bmatrix} 2 \\ \frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}, x'_2 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ 3 \\ -\frac{1}{\sqrt{6}} \end{bmatrix}, x'_3 = \begin{bmatrix} -1 \\ \frac{1}{\sqrt{2}} \\ 3 \end{bmatrix}, x'_4 = \begin{bmatrix} -1 \\ \frac{1}{\sqrt{2}} \\ 3 \end{bmatrix}, x'_5 = \begin{bmatrix} -2\sqrt{2} \\ 0 \end{bmatrix}, x'_6 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ 3 \end{bmatrix}$$

Η ποσότητα της πληροφορίας που χάνεται με τη νέα αυτή αναπαράσταση σε σχέση με αυτήν της εκφώνησης είναι βάση της θεωρίας:

$$\text{απώλεια πληροφορίας} \Rightarrow \frac{\sum(\text{ιδιοτιμών που δε χρησιμοποιούνται})}{\sum(\text{όλων των ιδιοτιμών})} = \frac{\lambda_3}{\lambda_1 + \lambda_2 + \lambda_3} = 0$$

Επομένως δεν υπάρχει καμιά απώλεια πληροφορίας από την μετάβαση από τις τρεις διαστάσεις στις 2.

Επαναλαμβάνουμε τη λύση με την R (ερώτημα (γ)).

Σημείωση: Η λύση παρακάτω πραγματοποιείται με την εντολή `prcomp`. Για μια λύση σε αντιστοιχία με την παραπάνω χρησιμοποιήστε τις παρακάτω εντολές:

```
> covmat = cov(pdata)
> eigenvalues = eigen(covmat)$values
```

```
> eigenvectors = eigen(covmat)$vectors
> pdata_pc = as.matrix(pdata) %*% eigenvectors[, 1:2]
> info_loss = eigenvalues[3] / sum(eigenvalues)
```

Για να εφαρμόσουμε PCA στην R εκτελούμε την παρακάτω εντολή:

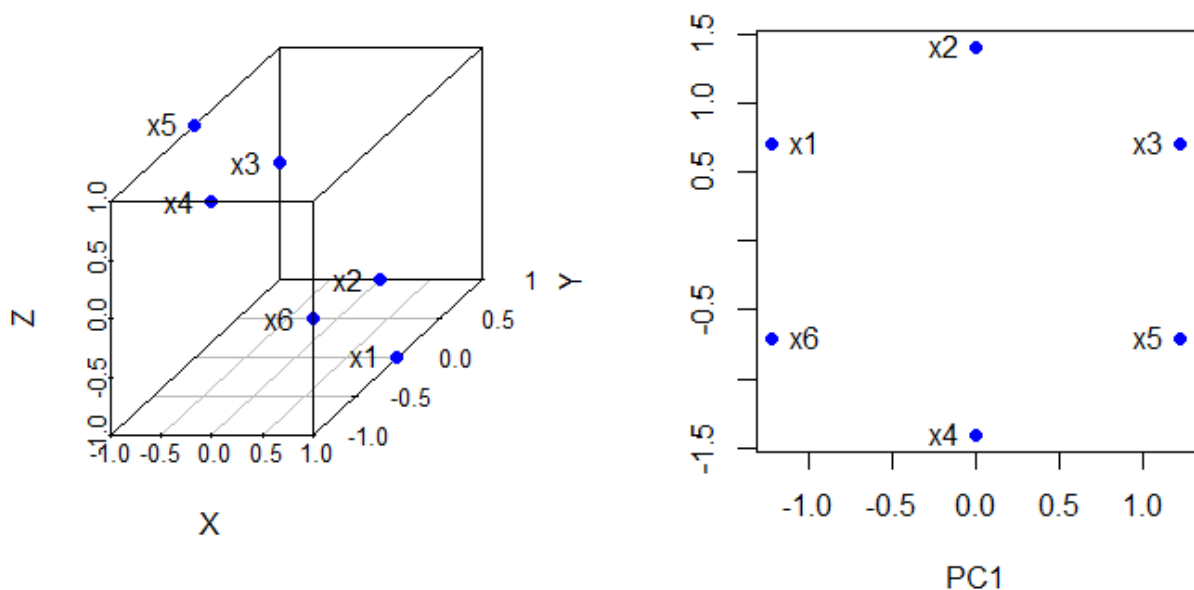
```
> pca_model <- prcomp(pdata)
```

Μπορούμε επιπλέον να δούμε τις ιδιοτιμές και τα ιδιοδιανύσματα που προκύπτουν και να εμφανίσουμε τη μεταβλητότητα για κάθε principal component εκτελώντας τις παρακάτω εντολές:

```
> eigenvalues = pca_model$sdev^2
> eigenvectors = pca_model$rotation
> barplot(pca_model$sdev ^ 2 / sum(pca_model$sdev ^ 2))
> info_loss = eigenvalues[3] / sum(eigenvalues)
```

Μπορούμε τέλος να εφαρμόσουμε το μετασχηματισμό και να σχεδιάσουμε τα νέα δεδομένα εκτελώντας τις παρακάτω εντολές:

```
> pc <- predict(pca_model, pdata)
> plot(pc, col = "blue", pch = 19)
> text(pc, labels=row.names(pdata), pos=2)
```



4 Επιλογή και Εξαγωγή Χαρακτηριστικών (Correlation, PCA)

Για την κατασκευή ενός μοντέλου PCA θα χρησιμοποιήσουμε ως εφαρμογή τα δεδομένα εκπαίδευσης του αρχείου που δίνεται (engdata.txt) για ένα πρόβλημα δυαδικής ταξινόμησης.

Ένα summary των δεδομένων με την R είναι το παρακάτω:

Age	Salary	YearsOfStudy	WorkExp	Location
Min. :20.00	Min. : 475	Min. : 3.000	Min. :10.0	EU:345
1st Qu.:44.00	1st Qu.:1156	1st Qu.: 5.000	1st Qu.:39.0	US:305
Median :51.00	Median :1517	Median : 6.000	Median :46.0	
Mean :49.94	Mean :1592	Mean : 6.032	Mean :44.9	
3rd Qu.:57.00	3rd Qu.:1969	3rd Qu.: 7.000	3rd Qu.:52.0	
Max. :70.00	Max. :3900	Max. :10.000	Max. :68.0	

Αρχικά, εισάγουμε τα δεδομένα και τα διαχωρίζουμε σε χαρακτηριστικά και χαρακτηριστικό κλάσης:

```
> engdata = read.csv("engdata.txt")  
> Location = engdata[, 5]  
> engdata = engdata[, 1:4]
```

Στη συνέχεια, θα απαντήσουμε στα παρακάτω ερωτήματα:

α) Σχεδιάστε τις τιμές του συνόλου δεδομένων ανά δύο χαρακτηριστικά με διαφορετικά χρώματα/σύμβολα για κάθε κλάση.

β) Υπολογίστε το correlation matrix για τα χαρακτηριστικά. Ποια χαρακτηριστικά θα μπορούσαν να παραληφθούν;

γ) Κατασκευάστε ένα μοντέλο PCA από τα δεδομένα και βρείτε τις ιδιοτιμές και τα ιδιοδιανύσματα. Σχεδιάστε το ποσοστό πληροφορίας για τα principal components.

δ) Εφαρμόστε το μοντέλο στα δεδομένα και διατηρήστε τις δύο πρώτες διαστάσεις. Σχεδιάστε τα νέα δεδομένα σε δισδιάστατο άξονα.

ε) Ανακατασκευάστε τα δεδομένα και σχεδιάστε τα εκ νέου (όπως στο (α)). Υπολογίστε την απώλεια της πληροφορίας.

4.1 Σχεδίαση Δεδομένων και Υπολογισμός Correlation

Για τα ερωτήματα (α) και (β) εκτελούμε τις εντολές:

```
> plot(engdata, col = Location, pch = c("o", "+")[Location])  
> cor(engdata)
```

4.2 Εφαρμογή Αλγορίθμου PCA

Για το (γ) υπολογίζουμε αρχικά το PCA και βρίσκουμε ιδιοτιμές και ιδιοδιανύσματα:

```
> pca_model <- prcomp(engdata, center = TRUE, scale = TRUE)  
> eigenvalues = pca_model$sdev^2  
> eigenvectors = pca_model$rotation
```

Για να σχεδιάσουμε το ποσοστό πληροφορίας για τα principal components εκτελούμε την παρακάτω εντολή:

```
> barplot(eigenvalues / sum(eigenvalues))
```

Εφαρμόζουμε το PCA στα δεδομένα ώστε να μειώσουμε τις διαστάσεις τους σε δύο και στη συνέχεια τα σχεδιάζουμε (ερώτημα (δ)):

```
> engdata_pc <- as.data.frame(predict(pca_model, engdata)[, 1:2])  
> plot(engdata_pc, col = Location, pch = c("o", "+")[Location])
```

Τέλος, για να ανακατασκευάσουμε τα δεδομένα (ερώτημα (ε)) και να τα σχεδιάσουμε εκτελούμε τις εντολές:

```
> engdata_pc[, 3:4] <- 0  
  
> engdata_rec = data.frame(t(t(as.matrix(engdata_pc) %*%  
    t(pca_model$rotation)) * pca_model$scale  
    + pca_model$center))  
  
> plot(engdata_rec, col = Location, pch = c("o", "+")[Location])
```

Η απώλεια της πληροφορίας υπολογίζεται με την εντολή:

```
> info_loss = (eigenvalues[3] + eigenvalues[4]) / sum(eigenvalues)
```

5 Μη Γραμμική Εξαγωγή Χαρακτηριστικών (ISOMAP)

Για την κατασκευή ενός μοντέλου ISOMAP θα χρησιμοποιήσουμε ως εφαρμογή τα δεδομένα εκπαίδευσης του αρχείου που δίνεται (srdata.txt).

Ένα summary των δεδομένων με την R είναι το παρακάτω:

V1	V2	V3
Min. : -0.4738629	Min. : -0.55193	Min. : 0.0006573
1st Qu.: -0.1502609	1st Qu.: -0.21918	1st Qu.: 0.1290643
Median : -0.0006371	Median : -0.02540	Median : 0.2416298
Mean : 0.0058976	Mean : -0.05905	Mean : 0.2498458
3rd Qu.: 0.1935721	3rd Qu.: 0.08422	3rd Qu.: 0.3734662
Max. : 0.6282747	Max. : 0.39584	Max. : 0.4999653

Αρχικά, εισάγουμε τα δεδομένα:

```
> srdata = read.csv("srdata.txt")
```

Στη συνέχεια, θα απαντήσουμε στα παρακάτω ερωτήματα:

α) Σχεδιάστε τα δεδομένα σε τρισδιάστατο γράφημα.

β) Εφαρμόστε τον αλγόριθμο ISOMAP με $k = 4$ ώστε να μετασχηματίσετε τα δεδομένα σε δύο διαστάσεις.

γ) Σχεδιάστε εκ νέου τα δεδομένα σε τρισδιάστατο γράφημα με διαφορετικά χρώματα ανάλογα με το αποτέλεσμα του ISOMAP (την πρώτη στήλη).

δ) Σχεδιάστε τα νέα δεδομένα σε ένα νέο δισδιάστατο γράφημα, με διαφορετικά χρώματα ανάλογα με το αποτέλεσμα του ISOMAP (την πρώτη στήλη).

5.1 Σχεδίαση Δεδομένων

Για το ερώτημα (α) εκτελούμε την εντολή:

```
> scatterplot3d(srdata, angle = 88, scale.y = 5)
```

οπότε προκύπτει το σχήμα που ζητείται.

5.2 Εφαρμογή Αλγορίθμου ISOMAP

Για το (β) υπολογίζουμε το ISOMAP και λαμβάνουμε τα νέα δεδομένα:

```
> srdata_dist <- dist(srdata)
> isom <- isomap(srdata_dist, ndim=2, k = 4)
> srdata_2d <- isom$points
```

Για να χρησιμοποιήσουμε το αποτέλεσμα του ISOMAP ως το χρώμα των δεδομένων (ερωτήματα (γ), (δ)), αρχικά κατασκευάζουμε την παρακάτω μεταβλητή:

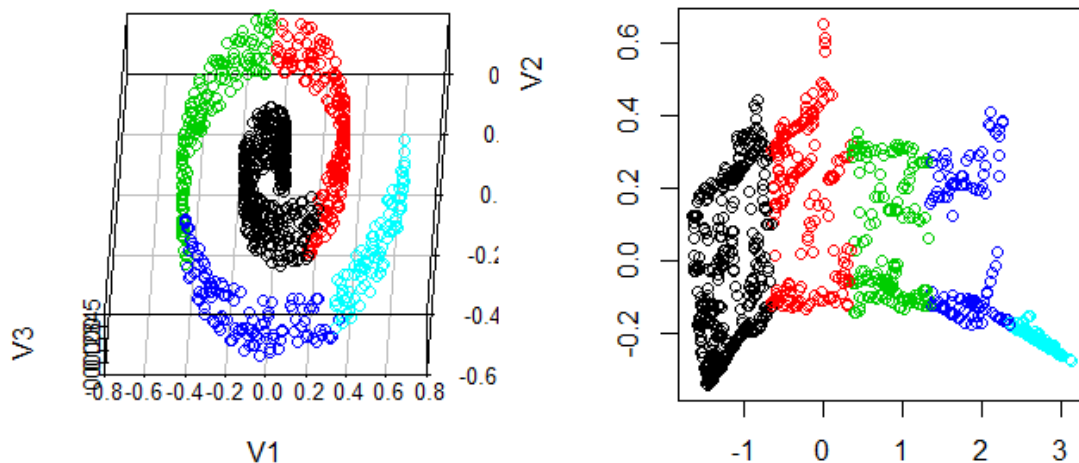
```
> colors = srdata_2d[,1] - min(srdata_2d[,1]) + 1
```

Στη συνέχεια εκτελούμε τις παρακάτω εντολές ώστε να σχεδιάσουμε με το χρωματισμό του ISOMAP τα αρχικά δεδομένα (ερώτημα (γ)):

```
> scatterplot3d(srdata, angle = 88, scale.y = 5, color = colors)
```

και τα μετασχηματισμένα δεδομένα σε νέο γράφημα (ερώτημα (δ)):

```
> x11(); plot(srdata_2d, col = colors)
```



6 Πρόβλημα για Εξάσκηση

Εισάγετε τα δεδομένα Glass του πακέτου mlbench και διαχωρίστε τα σε δεδομένα εκπαίδευσης και δεδομένα ελέγχου καθώς και σε χαρακτηριστικά και χαρακτηριστικό κλάσης με τις παρακάτω εντολές:

```
> data(Glass, package = "mlbench")  
> training = Glass[c(1:50, 91:146), -10]  
> trainingType = factor(Glass[c(1:50, 91:146), 10])  
> testing = Glass[51:90, -10]  
> testingType = factor(Glass[51:90, 10])
```

Απαντήστε στα παρακάτω ερωτήματα:

α) Εφαρμόστε PCA στα δεδομένα εκπαίδευσης.

β) Μετασχηματίστε τα δεδομένα εκπαίδευσης και τα δεδομένα ελέγχου και διατηρήστε τα 2 πρώτα principal components.

γ) Εφαρμόστε τον kNN με $k = 3$ ώστε να κατατάξετε τα δεδομένα ελέγχου και υπολογίστε το accuracy για τα δεδομένα ελέγχου.

δ) Επαναλάβετε τα ερωτήματα (β) και (γ) διατηρώντας κατά σειρά τα 3, 4 και 5 πρώτα principal components. Τι παρατηρείτε;