

► Απλή Γραμμική Πολυπρόσθετη

Δύο ή περισσότερες τυχαίες μεταβλητές π.χ. X, Y
 και έστω τυχαία δείγματα x_1, \dots, x_n από τη X
 και έστω τ.δ. y_1, \dots, y_n από τη Y

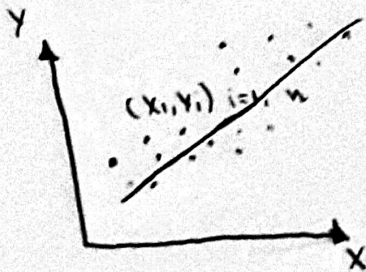
οπότε έχω δεδομένα

X	x_1	x_2	\dots	x_n
Y	y_1	y_2	\dots	y_n

Ερώτηση: Υπάρχει κάποια (σκέψη γραμμική?) που να συνδέει τις X, Y ?

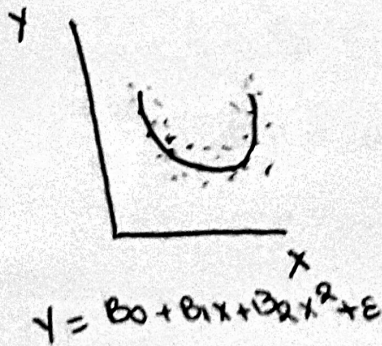
↳ Βίλια για να απαντήσω στο Ερώτημα

Διάγραμμα διασποράς

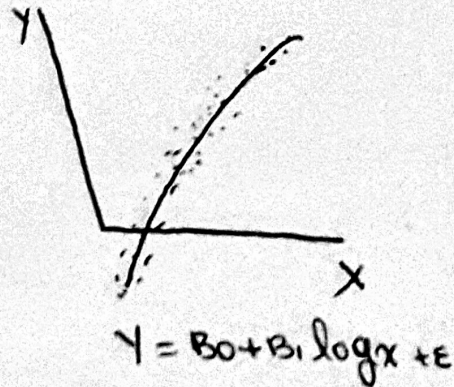


οι μεταβλητές X, Y τείνουν να
 προσεγγίζονται από μια ευθεία
 $Y = \beta_0 + \beta_1 X + \epsilon$

Στην περίπτωση που είχα τέτοιο διάγραμμα δεν θα μπορούσα να
 το προσεγγίσω από μια ευθεία



ή



Οι σχέσεις αυτές δεν είναι άπολετες, δηλ. περιμένω τα σφάλματα να παρατείνουν -
 παλινδρομούν γύρω από τη καμπύλη που τείνει να ελαττώσει

Αυτό αριθμώς περιγράφει το $-\epsilon-$
Περιοριζόμαστε σε γραμμικές σχέσεις

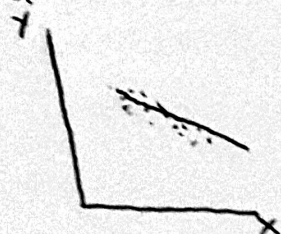
Το διάνυσμα διασποράς ενοφάει

1. Την μορφή της σχέσης X και Y ή X ερπυλιωτή σχέση $Y = \beta_0 + \beta_1 X + \epsilon$

2. Κατεύθυνση της ερπυλιωτής σχέσης



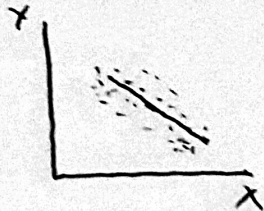
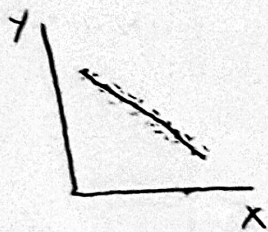
Μεγαλώνει η X και η Y
άρα θετική ματεύθ.



Μεγαλώνει η X λιμνώνει η Y
άρα αρνητική ματεύθ.

3. Πόσο έντονη είναι η ερπυλιωτή σχέση

Επιπλέον πόσο παλινδρόμηση από το ερπυλιωτό μοντέλο που τείνει να
τα περιγράψει



Δηλαδή η $Y = \beta_0 + \beta_1 X + \epsilon$ είναι το μοντέλο απλής ερπυλιωτής παλινδρόμησης
• Απλή = μεταξύ δύο μεταβλητών

ή αλλιώς $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, i=1, n$

Η μεταβλητή X ονομάζεται ανεξάρτητη ή επεξηγήσιμη

Η Y ονομάζεται εξαρτημένη ή απόκριση

Τα ϵ_i ονομάζονται εσφαλματα του μοντέλου, (η απόσταση του σημείου από το
μοντέλο)

Το β_0 ονομάζονται παράμετροι του μοντέλου

2^ο Βήμα Κατασκευή του μοντέλου α.δ.π

Εύρεση ευαθέτων των παραμέτρων β_0, β_1

Ευαθέτες. Ελαχίστων Τετραγώνων (Ε.Σ.Τ)
των β_0, β_1 (Pearson, 1910)

Η ιδέα της μεθόδου ευαθέτων: Θέλω να βρω την 'υαλύτερη' ευθεία που θα εσπεριλαμβάνει τα στοιχεία

Το κριτήριο είναι τα στοιχεία να εσπεριλαμβάνονται όσο πιο κοντά γίνεται ή στην ευθεία. ή δηλαδή έτσι που τα εσπεριτοποιούνται

Αρα λοιπόν αρμεί να βρω τα β_0, β_1 που ελαχιστοποιούν το

$$S = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

(το τετράγωνο είναι για να την ελαχιστοποιώ έχω αρνητικά)

• $\frac{\partial S}{\partial \beta_0} = 0$ και $\frac{\partial S}{\partial \beta_1} = 0$

δηλαδή
$$\begin{cases} \beta_0 + \beta_1 \sum x_i = \sum y_i \\ \beta_0 \sum x_i + \beta_1 \sum x_i^2 = \sum x_i y_i \end{cases}$$

μανουιές ερωτήσεις

Οπότε, η επίλυση των μανουιών ερωτήσεων οδηγεί στους Ε.Σ.Τ των

β_0 και β_1 που είναι
$$\hat{\beta}_1 = \frac{\sum x_i y_i - (\sum x_i)(\sum y_i)/n}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Τελικά η υαλύτερη ευθεία που θα υπορρώσω να έχω, δηλαδή με τα λιγρότερα σφάλματα προκύπτει από την αυτοαπόσπαση των παραμέτρων από τους ευαθέτες τους :

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

Ευτιώμενο μοντέλο α.δ.π

σπουηερμηνείας \blacktriangleright το $\hat{\beta}_0$ παριστά την τιμή της \hat{y} για $x=0$

\blacktriangleright το $\hat{\beta}_1$ παριστά την μεταβολή της \hat{y} σε μοναδιαία μεταβολή της x

Διόρθωση της ορθότητας του μοντέλου

► Υπόλοιπα: Ορίζονται ως οι αποκλίσεις του μοντέλου από την πραγματικότητα που το μοντέλο ελπίζει να προφέρει

Ορισμός: $e_i = y_i - \hat{y}_i \quad i=1, n$

Ιδιότητα των υπολοίπων $\sum_{i=1}^n e_i = 0$

$$\begin{aligned} \Phi \sum_{i=1}^n e_i &= \sum_{i=1}^n y_i - \hat{y}_i = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \\ &= \sum_{i=1}^n (y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) - \hat{\beta}_1 x_i) \\ &= \sum_{i=1}^n (y_i - \bar{y}) + \hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x}) = 0 \end{aligned}$$

αφού $\sum (x_i - \bar{x}) = \sum x_i - \sum \bar{x} = n\bar{x} - n\bar{x} = 0$

και $\sum (\bar{x} - x_i) = 0$ το ίδιο \uparrow

για πρώτη φορά είναι τα υπολείμματα να είναι κοντά στο 0 σημαίνει ότι υπάρχει διαφορά μεταξύ από την πραγματικότητα

► Ανάλυση Ολικής Μεταβλητότητας στο μοντέλο της α.δ.π

Ολική μεταβλητότητα \approx Μεταβλητή διασπορά

ο Μεταβλητή διασπορά των y_1, y_n

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

ο Ολική μεταβλητότητα $= \sum (y_i - \bar{y})^2$

Εδώ θέλω να δώ πως η ολική μεταβλητότητα δίνεται από το μοντέλο, άρα πρέπει να το ελαττώσω

Το πιο απλό είναι να το προσεγγίσω

3

Δηλαδή \hat{y}_i

οπότε μετά από πρώτα (άσκηση προηγ.)

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

ολική μεταβολή

μεταβλητότητα που οφείλεται στο μοντέλο

μεταβλητότητα που οφείλεται στα υπόλοιπα

Δηλαδή $SS_{tot} = SS_{reg} + SS_{res}$



Αν το μεγαλύτερο μέρος από την ολική μεταβλητότητα περιέχεται σε SS_{res} τότε αυτό είναι κακό.

Δηλαδή ένα δεύτερο κριτήριο είναι Αν $SS_{reg} > SS_{res}$ τότε το μοντέλο της α.δ.π φαίνεται να είναι υποκείμενο αλλιώς γυμνασίου το αντίθετο

• Άλλη κορφή για το άθροισμα τετραδ. που οφείλεται στην πολυωνομ SS_{reg}

$$SS_{reg} = \sum (\hat{y}_i - \bar{y})^2 = \sum (\hat{\beta}_0 + \hat{\beta}_1 x_i - \bar{y})^2$$

$$= \sum (\bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i - \bar{y})^2$$

$$= \hat{\beta}_1^2 \sum (x_i - \bar{x})^2$$

Άρα λοιπόν προσπαθούμε να αναδείξουμε αυτή τη σχέση κατασκευάζοντας έναν πίνακα:

Πίνακας Ανάλυσης Διακύμανσης του Μοντέλου της α.β.η

ΠΙΝΑΚΑ

ΑΝΑΔΙΑ

πηγή μεταβλητότητας	Αθρ. τετρα. SS	Β.ε	MS μέσο τετραδ.	F πηγή
Μοντέλο α.β.η	SS _{reg}	1	MS _{reg} = $\frac{SS_{reg}}{1}$	F = $\frac{MS_{reg}}{MS_{res}}$
Υπόλοιπο	SS _{res}	n-2	MS _{res} = $\frac{SS_{res}}{n-2}$	
ολική μεταβλητότητα	SS _{tot}	n-1		

Για να αξιοποιήσω αυτόν τον πίνακα χρησιμοποιώ το 2^ο κριτήριο που προαναφέρθηκε

► Βαθμός ελευθερίας ενός αθροίσματος τετραγώνου: είναι το πλήθος των ανεξάρτητων πληροφοριών στο Y_1, \dots, Y_n τις οποίες πρέπει να διαθέσουμε ώστε να μπορούμε να υπολογίσουμε το αντίστοιχο άθροισμα

~~Παράδειγμα~~

Επειδή

Ας πούμε για το SS_{tot} γέρω αρκεί να έχω n-1 πληροφορίες για να το υπολογίσω

$$SS_{tot} = \sum (Y_i - \bar{Y})^2 \text{ δηλαδή}$$

$$\begin{matrix} Y_1 - \bar{Y} \\ Y_2 - \bar{Y} \\ \vdots \\ Y_n - \bar{Y} \end{matrix}$$

όμως !! $\bar{Y} = \frac{1}{n} \sum Y_i$ δηλαδή

μπορώ να εκφράσω το Y_n συναρτήσει των άλλων άρα έχω n-1 β.ε

4

Απόδειξη ότι Β.Ε του SS_{tot} είναι $n-1$

ο $SS_{tot} = \sum (y_i - \bar{y})^2$

όπου $\bar{y} = \frac{1}{n} \sum y_i \Rightarrow \sum y_i = n\bar{y} \Rightarrow y_n + \sum_{i=1}^{n-1} y_i = n\bar{y}$

$\Rightarrow y_n = n\bar{y} - \sum_{i=1}^{n-1} y_i$

ο $y_n - \bar{y} = n\bar{y} - \sum_{i=1}^{n-1} y_i - \bar{y} = (n-1)\bar{y} - \sum_{i=1}^{n-1} y_i$

$\Rightarrow y_n - \bar{y} = - \sum_{i=1}^{n-1} (y_i - \bar{y})$

$\Rightarrow (y_n - \bar{y})^2 = \left[\sum_{i=1}^{n-1} (y_i - \bar{y}) \right]^2$

δηλαδή $SS_{tot} = \sum (y_i - \bar{y})^2 = \sum_{i=1}^{n-1} (y_i - \bar{y})^2 + (y_n - \bar{y})^2$

$= \sum_{i=1}^{n-1} (y_i - \bar{y})^2 + \left[\sum_{i=1}^{n-1} (y_i - \bar{y}) \right]^2$

Άρα χρειάζεται να ξέρω $y_1 - \bar{y} \dots y_{n-1} - \bar{y}$

για να υπολογίσω το SS_{tot}

Άρα οι Β.Ε είναι $n-1$

Επιπλέον κοινός για Β.Ε

Το SS_{tot} έχει Β.Ε: αριθμός παρατηρήσεων - 1
(κερδός δείγματος) - 1

Το SS_{reg} έχει Β.Ε: άσπληθος ανεξάρτητων μεταβλητών
(για την περίπτωση μας είναι μόνο η X)

Το SS_{res} έχει Β.Ε την διαφορά

Συντελεστής αφορισμού ή προσδιοριστικότητας

$$R^2 = \frac{SS_{\text{reg}}}{SS_{\text{tot}}} = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

από τη σχέση $SS_{\text{reg}} + SS_{\text{res}} = SS_{\text{tot}}$

και 2 ιδιότητες είναι οι εξής:

1. ο R^2 είναι αριθμός

2. $0 \leq R^2 \leq 1$ και άρα ευφραίνεται και ως ποσοστό!

Ενδιαφέρον παρουσιάζουν οι τιμές 0 και 1

- ▶ Τιμή του R^2 κοντά στο 1: δηλαδή θα πρέπει $SS_{\text{reg}} \approx SS_{\text{tot}}$
δηλ $SS_{\text{reg}} \gg SS_{\text{res}}$
υποκείμενο μοντέλο
- ▶ Τιμή του R^2 κοντά στο 0: δηλαδή θα πρέπει $SS_{\text{reg}} \ll SS_{\text{res}}$
 ~~$SS_{\text{reg}} \approx SS_{\text{res}}$~~ λη υποκείμενο μοντέλο

Μια ερμηνεία διορθωμένο R^2 : Ευφραίνει το ποσοστό της ολικής μεταβλητότητας των Y_i , Y_n που ερμηνεύεται από το μοντέλο