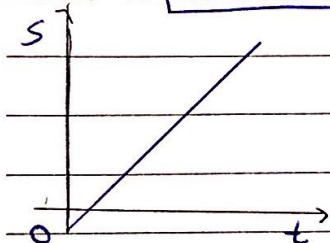


Δευτέρα 12/10/2020

~~Εισαγωγή γραμμική παλινδρόσηση~~

Παράδειγμα:

ΕΟΚ \rightarrow $s = u \cdot t$ \rightarrow Απόλυτο / αιτιοκρατικό / υπερμινιστικό μοντέλο



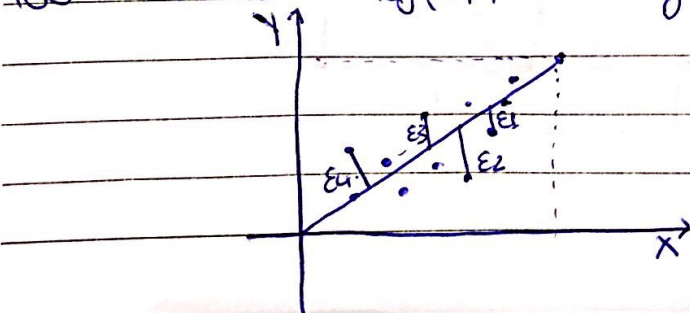
Όσα δεδομένα και να πάρω βρίσκονται απαραίτητα - απόλυτα πάνω στην γραμμή

Έστω ότι κάποιος σε μια μέση σχολείται με το βάρος $= X$ και το εύρος $= Y$ των παιδιών. Υπάρχει κάποια σχέση μεταξύ τους. Υπάρχει κάποιο f (μοντέλο) ώστε $Y = f(X)$?

Διαλέγω n μέλη του πληθυσμού και παράγω πειραματικά δεδομένα στις X, Y

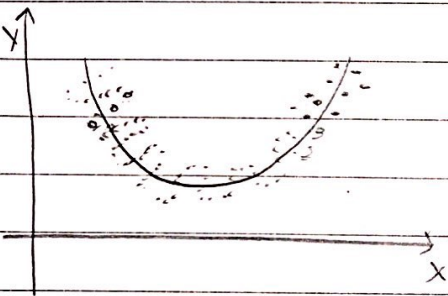
X	x_1	x_2	...	x_n
Y	y_1	y_2	...	y_n

Ερώτημα: Μπορούμε να κατασκευάσουμε ένα μαθηματικό μοντέλο στο οποίο να συνδέει τη X με την Y μεταβλητή. Δηλαδή έναν τύπο που να περιλαμβάνει τα σημεία (x_i, y_i) του ακόλατου διαγράμματος για $i=1, \dots, n$.



Στο 1^ο παράδειγμα είχαμε ένα αιτιοκρατικό πείραμα ενώ στο 2^ο παράδειγμα είχαμε ένα τυχαίο πείραμα με τυχαίες μεταβλητές. Επομένως, δεν αναφερόμαστε στο 2^ο παράδειγμα σε ένα μοντέλο αιτιώδες αλλά ~~είχαμε~~ έχουμε ένα μοντέλο τ.μ που διαφέρει από αυτό γιατί σε αυτό λόγω της διαφορετικότητας των συνθηκών. Για το σκοπό αυτό θεωρούμε ένα γραμμικό κατά προσέγγιση μοντέλο:

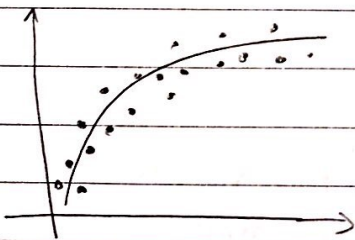
$$Y = \beta_0 + \beta_1 X + \varepsilon, \text{ με } \varepsilon: \begin{matrix} \text{οι} \\ \text{αποκλίσεις των σημείων } (x_i, y_i) \\ i=1, \dots, n \text{ από το γραμμικό μοντέλο} \end{matrix}$$



Όσοι εδώ υιοθετούμε το μοντέλο της γραμμής.

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \varepsilon$$

Η γραμμή είναι "επιλογή" προσέγγιση για τη συγκεκριμένη κατανομή των σημείων



$$Y = a \log x + \varepsilon, \text{ όπου με παραπάνω}$$

Σκοπός του μεθόδου είναι η κατασκευή ενός γραμμικού μοντέλου που περιγράφει τη σχέση μεταξύ δύο ή περισσότερων τ.μ.

ΕΝΟΤΗΤΑ 1 :

Απλή Γραμμική Παλινδρόμηση (α.γ.π)

Έστω δύο τ.μ X και Y . Ειδικότερα η X είναι η μεταβλητή την οποία χειρίζεται ο στατιστικός και Y είναι η τ.μ που επιθυμούμε να δούμε πως σχετίζεται πως συμπεριφέρεται σε μεταβολές της X . Θεωρώ τειχ. δείγματα από τις X και Y μεγέθους n .

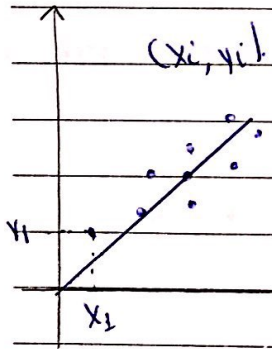
X	x_1	x_2	\dots	x_n
Y	y_1	y_2	\dots	y_n

Ερώτημα :

Υπάρχει κάποια σχέση (γραμμική) που συνδέει τα X και Y ?
Πόσο αξιόπιστη είναι αυτή?

Βήμα 1^ο : [Διόρθωση Διασποράς]

Απεικονίζουμε τα δεδομένα σε ένα σύστημα αξόνων

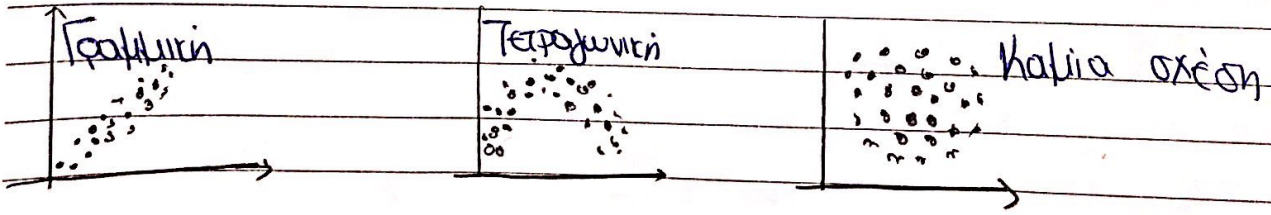


Βλέπουμε ότι δεν βρίσκονται όλα τα σημεία των πειραματικών δεδομένων πάνω στην ευθεία. Αντίθετα υπάρχει μια ευθεία $Y = \beta_0 + \beta_1 X$ που τείνει κατά προσέγγιση να συμπεριχθεί τα (x_i, y_i) . Τα σημεία (x_i, y_i) $i=1, 2, \dots, n$ παλινδρόμουν γύρω από το γραμμικό μοντέλο.

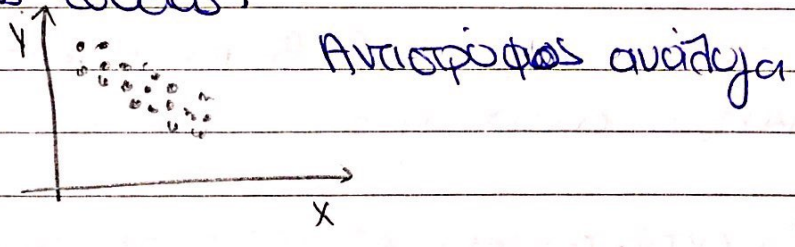
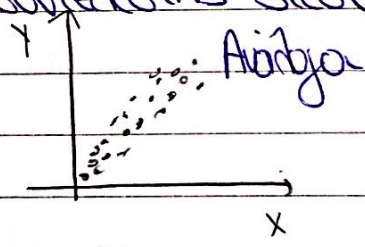
Ανταπόκριση υπάρχει η σχέση $Y = \beta_0 + \beta_1 X + \epsilon$, με ϵ : η απόκλιση των σημείων από την παραπάνω ευθεία.

Το Διόρθωση Διασποράς μας δείχνει:

a | Τη μορφή της σχέσης μεταξύ των X και Y .



81 Την κατεύθυνση της σχέσης (αν αυτή είναι γραμμική)
 ένδειξη αν το Y μεγαλώνει όταν μεγαλώνει το X ή αντιστοίχα
 αν μικραίνει. Λοιπόν η αντιστροφή ανάγλυφα, παίζει πολύ ο
 αντίστροφος δείκτης της σχέσης)



82 του βαθμού της σχέσης, ένδειξη ποσο κοντά είναι τα σημεία του
 διαγράμματος διασποράς στο μοντέλο που το περιγράφει.

Η αδυναμία του διαγράμματος διασποράς είναι ότι αναφέρεται
 σε δύο μόνο μεταβλητές.

Αν το διαγράμμα διασποράς υποθέτει γραμμικότητα των Y και X
 Στινάθε το μοντέλο της Άλβης Γραμμικής Πολυώνυμικής

Παράμετροι

$$Y = \beta_0 + \beta_1 X + \epsilon \Leftrightarrow Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, i=1, \dots, n$$

σφάλματα του μοντέλου

εξαρτημένη ή απόκριση

αυξήσιμη ή εξηγητική

Το ϵ περιέχει απροσδόκιστα γεγονότα μέχρι κυρίως την ύπαρξη
 άλλων παραγόντων που επηρεάζουν το X και Y και τους οποίους
 δεν λαμβάνω υπόψη στο μοντέλο

2^ο Βήμα: Κατασκευή του μοντέλου α.β.π.

Σκοπός είναι να προσδιορίσουν τα β_0 και β_1 .
 Εκτιμήτες ελάχιστων Τετραγώνων (Pearson ~ 1900)

$\left. \begin{array}{l} \text{Τρία εκτιμήσεις} \\ b_0, b_1 \end{array} \right\}$ Τα καλύτερα b_0, b_1 είναι εκείνα για τα οποία ελαχιστοποιούνται οι αποστάσεις (ϵ_i) των σημείων του διαγρ. διασκορπισμού από το μοντέλο της α.γ.π.

Οπότε θα βρούμε τα b_0, b_1 που ελαχιστοποιούν συνολικά τα ϵ_i ή κάποια συνάρτηση των ϵ_i ή την $S^2 = \sum_{i=1}^n \epsilon_i^2$

Οι εκτιμήσεις των ελαχίστων τετραγώνων των b_0, b_1 προκύπτουν από ελαχιστοποίηση ως προς b_0, b_1 της.

$$S^2 = S^2(b_0, b_1) = \sum_{i=1}^n \epsilon_i^2 \stackrel{\text{α.γ.π.}}{=} \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2$$

Για την εύρεση των b_0, b_1 που ελαχιστοποιούν την S^2 .

$$\left. \begin{array}{l} \frac{\partial S^2}{\partial b_0} = 0 \\ \frac{\partial S^2}{\partial b_1} = 0 \end{array} \right\} \Rightarrow \left. \begin{array}{l} n b_0 + b_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i \\ b_1 \sum_{i=1}^n x_i + b_0 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i \end{array} \right\} \begin{array}{l} \text{Κανονικές} \\ \text{Εξισώσεις} \end{array}$$

Καταλήγουμε μετά από πράξεις.

$$\hat{b}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

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

Η ~~ευθεία~~ $\hat{y} = \hat{b}_0 + \hat{b}_1 x$ είναι η πιο "καλή" ευθεία που παίρνεται με τέτοιο τρόπο ώστε από τα σημεία να πλησιάζουν όσο το δυνατόν πιο κοντά στο γραμμικό μοντέλο $y = b_0 + b_1 x + \epsilon$. Αυτό το μοντέλο λέγεται εκτιμώμενο μοντέλο α.γ.π. Χρησιμοποιείται για την πρόβλεψη της y για συγκεκριμένες τιμές της x .

Για το θέμα 3:

Υπόλοιπο: Η απόκλιση του μοντέλου από την πραγματικότητα που "επιτρέπει" το μοντέλο να περιγράψει
Συμβολίζουμε $e_i = y_i - \hat{y}_i$, $i=1, \dots, n$.

Ιδιότητα $\sum_{i=1}^n e_i = 0$

Αποδ:

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

$$\text{αφού } \sum_{i=1}^n (y_i - \bar{y}) = \sum y_i - \sum \bar{y} = n\bar{y}_e - n\bar{y} = 0$$

Ανάλυση της ολικής μεταβλητότητας στο μοντέλο α.γ.π.

Ολική μεταβλητότητα ανώτερες και σαν την συστηματική διακύμανση των y_1, \dots, y_n δίνεται την $S^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 \approx \sum_{i=1}^n (y_i - \bar{y})^2$

Θέλουμε να μετρήσουμε πως η ολική μεταβλητότητα δίνεται από το εκτιμημένο μοντέλο.

$\Rightarrow SS_{tot}$.

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

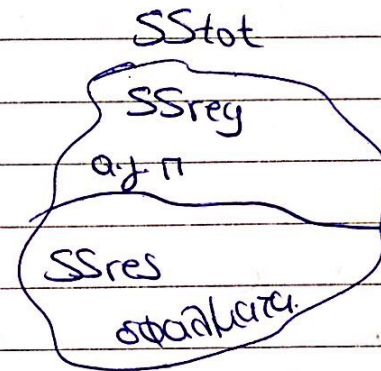
$$= \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \rightarrow \text{μεταβλητότητα που οφείλεται στα υψώματα (SSres)}$$

μεταβλητότητα που οφείλεται στο μοντέλο της α.φ.π (SSreg)

Αρα $SS_{tot} = \sum_{i=1}^n (Y_i - \bar{Y})^2$

$SS_{reg} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$

$SS_{res} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$



• Αν $SS_{reg} \gg SS_{res}$ το μοντέλο εγγίζει, εφικνότερα περνάει κέρως της ολικής μεταβλητότητας SS_{tot} .
Αρα είναι υψώσιμο μοντέλο α.φ.π.

• Αν $SS_{reg} \ll SS_{res}$. Τότε το μοντέλο της α.φ.π φαίνεται να μην είναι υψώσιμο