

Ανάλυση Δεδομένων με χρήση του Στατιστικού Πακέτου R



Δημήτρης Φουσκάκης,
Καθηγητής,
Τομέας Μαθηματικών,
Σχολή Εφαρμοσμένων Μαθηματικών και Φυσικών Εφαρμογών,
Εθνικό Μετσόβιο Πολυτεχνείο.

Περιεχόμενα

- Εισαγωγή στη Στατιστική
- Εισαγωγή στο Στατιστικό Πακέτο R
- Περιγραφική Στατιστική
- Προσομοίωση
- Στατιστική Συμπερασματολογία
 - Ένα Δείγμα
 - Δύο Ανεξάρτητα Δείγματα
 - Δείγματα κατά Ζεύγη
 - Ποσοστά
 - Έλεγχος καλής προσαρμογής
 - Πίνακες Συνάφειας 2×2 .
- **Ανάλυση Παλινδρόμησης**
- Ανάλυση Διασποράς

Εισαγωγή

- Αρκετές φορές σε μια στατιστική μελέτη ερχόμαστε αντιμέτωποι με το πρόβλημα της *πρόβλεψης* μιας μεταβλητής (*μεταβλητή απόκρισης*) όταν γνωρίζουμε τις τιμές κάποιας ή κάποιων άλλων μεταβλητών (*επεξηγηματικές μεταβλητές*).
- Ας θεωρήσουμε καταρχήν ότι έχουμε μία μόνο επεξηγηματική μεταβλητή.

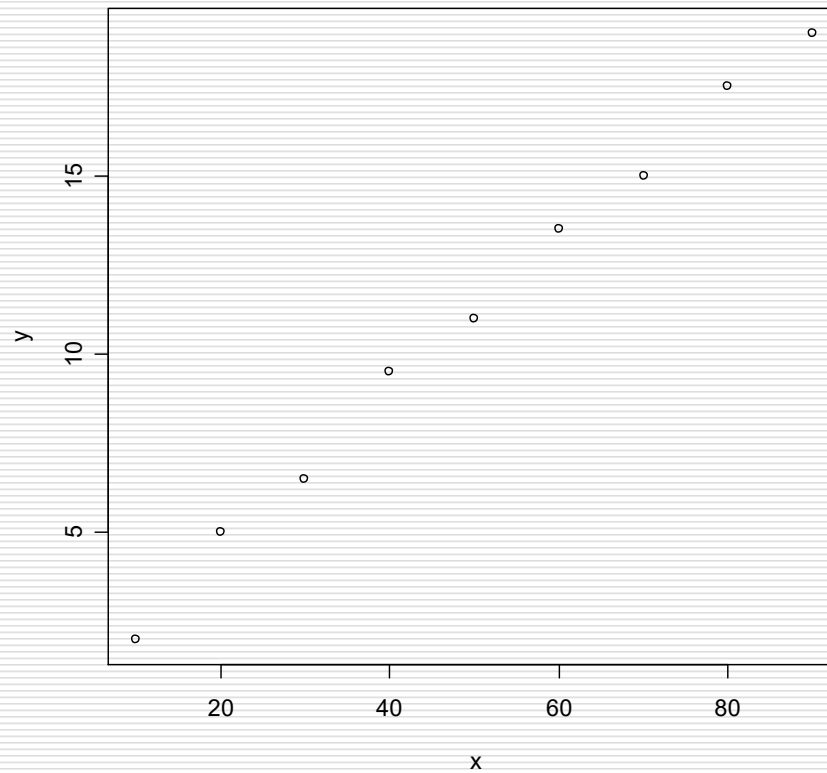
Απλό Γραμμικό Μοντέλο

- Πιο συγκεκριμένα ας υποθέσουμε ότι Y είναι η μεταβλητή απόκρισης και X η επεξηγηματική μεταβλητή και ας υποθέσουμε ότι και οι δύο μεταβλητές είναι **ποσοτικές**. Σκοπός μας είναι να δημιουργήσουμε ένα **μοντέλο**, έστω $Y=g(X)$, έτσι ώστε στο μέλλον να μπορούμε να προσδιορίσουμε την τιμή του Y με βάση την τιμή του X .

Απλό Γραμμικό Μοντέλο

- Πώς όμως επιλέγουμε την συναρτησιακή μορφή της $g(X)$; Η επιλογή της κατάλληλης συνάρτησης $g(X)$ μπορεί γίνει με την βοήθεια ενός τυχαίου δείγματος $(Y_1, X_1), \dots, (Y_n, X_n)$. Πιο συγκεκριμένα αν $(y_1, x_1), \dots, (y_n, x_n)$ οι παρατηρήσεις μας, τότε μπορούμε να σχηματίσουμε το γράφημα των σημείων (y_i, x_i) , γνωστό ως **διάγραμμα διασποράς** των σημείων, και να εκτιμήσουμε την συναρτησιακή μορφή της g .
- Στο επόμενο γράφημα, π.χ., το διάγραμμα διασποράς υποδεικνύει ότι η σχέση της X με την Y είναι γραμμική και άρα μπορούμε να θεωρήσουμε ότι $g(X) = a + bX$. Μένει τότε να εκτιμήσουμε τις τιμές των a και b , με την βοήθεια και πάλι της πληροφορίας που έχουμε από το τυχαίο δείγμα.

Απλό Γραμμικό Μοντέλο



Απλό Γραμμικό Μοντέλο

- Η ανάλυση μας λοιπόν είναι **εμπειρική** (βασίζεται στην υπάρχουσα εμπειρία μας με βάση το τυχαίο δείγμα) και άρα το μοντέλο μας είναι **στοχαστικό**. Αντιθέτως αν η ανάλυσή μας ήταν **θεωρητική**, γνωρίζαμε δηλαδή όλον τον πληθυσμό, το μοντέλο θα ήταν **προσδιοριστικό**.
- Στα στοχαστικά μοντέλα προφανώς έχουμε ελλιπή πληροφορία, το μοντέλο στο οποίο θα καταλήξουμε μπορεί να μην ικανοποιείται ακριβώς για ζεύγη τιμών του πληθυσμού που δεν έχουν παρατηρηθεί στο τυχαίο δείγμα. Για τον λόγο αυτό στο μοντέλο προσθέτουμε και ένα **τυχαίο σφάλμα ε** το οποίο θεωρούμε ότι προέρχεται από μια γνωστή κατανομή με άγνωστες παραμέτρους. Το στοχαστικό δηλαδή μοντέλο παίρνει την μορφή $Y = g(X) + \varepsilon$.
- Ακούγεται λογικό να θεωρήσουμε ότι η μέση τιμή του ε είναι 0, δηλαδή κατά μέσο όρο το σφάλμα μας είναι μηδενικό. Συνήθως θεωρούμε $\varepsilon \sim N(0, \sigma^2)$, με σ^2 άγνωστο.

Απλό Γραμμικό Μοντέλο

- Ας θεωρήσουμε τώρα ότι η g είναι μια γραμμική συνάρτηση, $g(X) = a + bX$, με a, b άγνωστες ποσότητες. Δηλαδή θεωρούμε ότι η τυχαία μεταβλητή X επηρεάζει γραμμικά την αναμενόμενη τιμή της τυχαίας μεταβλητής Y .

$$Y = a + bX + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2) \Leftrightarrow E(Y | X = x) = a + bx$$

- Ισοδύναμα σε τυχαίο δείγμα $(Y_1, X_1), \dots, (Y_n, X_n)$ θεωρούμε ότι ισχύει η σχέση

$$Y_i = a + bX_i + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2) \Leftrightarrow E(Y_i | X_i = x_i) = a + bx_i$$

με ε_i ανεξάρτητες (και ισόνομες) τυχαίες μεταβλητές.

- Τα ε_i καλούνται τυχαία σφάλματα, και παριστάνουν την άγνωστη κατακόρυφη απόκλιση της τιμής y_i από την ευθεία $E(Y_i | X_i = x_i) = a + bx_i$ για δοθείσα τιμή x_i .

Απλό Γραμμικό Μοντέλο

- Το παραπάνω μοντέλο καλείται **απλό** (γιατί έχουμε μία μόνο επεξηγηματική μεταβλητή X) **γραμμικό** (λόγω της γραμμικής συνάρτησης που χρησιμοποιούμε) **μοντέλο παλινδρόμησης**. Τα a , b και σ^2 είναι οι άγνωστες παράμετροι του μοντέλου μας (**συντελεστές μοντέλου**), τις οποίες θα εκτιμήσουμε με την βοήθεια των παρατηρήσεων που διαθέτουμε $(y_1, x_1), \dots, (y_n, x_n)$ που είναι οι τιμές ενός τυχαίου δείγματος $(Y_1, X_1), \dots, (Y_n, X_n)$.

Ερμηνεία παραμέτρων

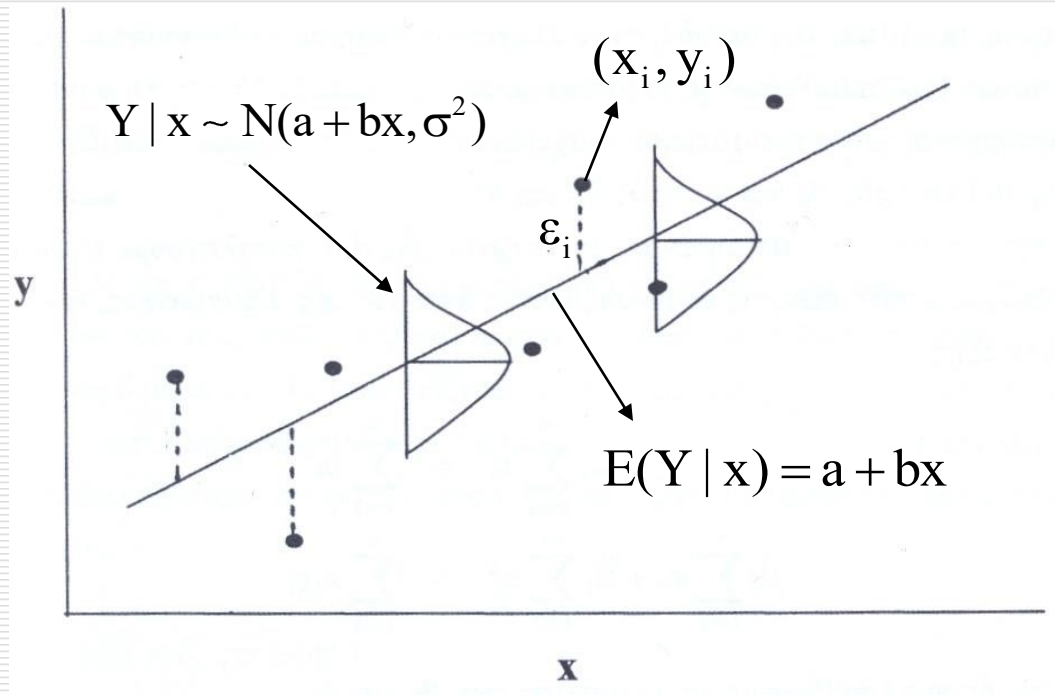
- Η σταθερά a εκφράζει την μέση τιμή της Y όταν το $X=0$.
- Η σταθερά b εκφράζει το πόσο αναμένεται να **μεταβληθεί** η αναμενόμενη τιμή της Y , αν η X αυξηθεί κατά μία μονάδα.

$$\frac{dY}{dX} = b$$

- Παραπάνω (και στις υπόλοιπες διαφάνειες) χρησιμοποιώ την λέξη **μεταβληθεί**, επειδή δεν γνωρίζω αν το b είναι θετικό ή αρνητικό. Όταν το b είναι θετικό χρησιμοποιήστε την λέξη **αυξηθεί**, αν το b είναι αρνητικό την λέξη **μειωθεί**.
- Η ποσότητα σ^2 εκφράζει την διασπορά των σφαλμάτων, την οποία θεωρούμε σταθερή ανεξάρτητα της τιμής της τ.μ. X (**υπόθεση ομοσκεδαστικότητας**). Επειδή η τυχαιότητα της Y δεδομένης μιας τιμής της $X = x$ οφείλεται στα σφάλματα, το σ^2 εκφράζει και την διασπορά της δεσμευμένης κατανομής της τ.μ. $Y|x$.

Απλό Γραμμικό Μοντέλο

$$Y = a + bx + \varepsilon, \varepsilon \sim N(0, \sigma^2) \Leftrightarrow Y | x \sim N(a + bx, \sigma^2)$$



Απλό Γραμμικό Μοντέλο

- Εκτιμώντας λοιπόν τα a και b από τα \hat{a} και \hat{b} αντίστοιχα καταλήγουμε στο

$$\hat{Y} = \hat{a} + \hat{b}x.$$

- Το \hat{Y} καλείται **προβλεπόμενη τιμή** και είναι όπως είδαμε η αναμενόμενη τιμή που θα πάρει η Y **όταν $X=x$** , όπως αυτήν την εκτιμήσαμε με βάση το μοντέλο παλινδρόμησης. **Η προβλεπόμενη τιμή είναι τ.μ., δηλαδή για διαφορετικό δείγμα ενδέχεται να πάρει άλλη τιμή όταν $X=x$.** Η προβλεπόμενη τιμή αποτελεί μία αμερόληπτη εκτιμήτρια της άγνωστης τιμής y που παίρνει η τ.μ. Y όταν $X=x$. Παρακάτω θα δούμε **δύο διαφορετικά Δ.Ε.** για την τιμή αυτή y όταν $X=x$.

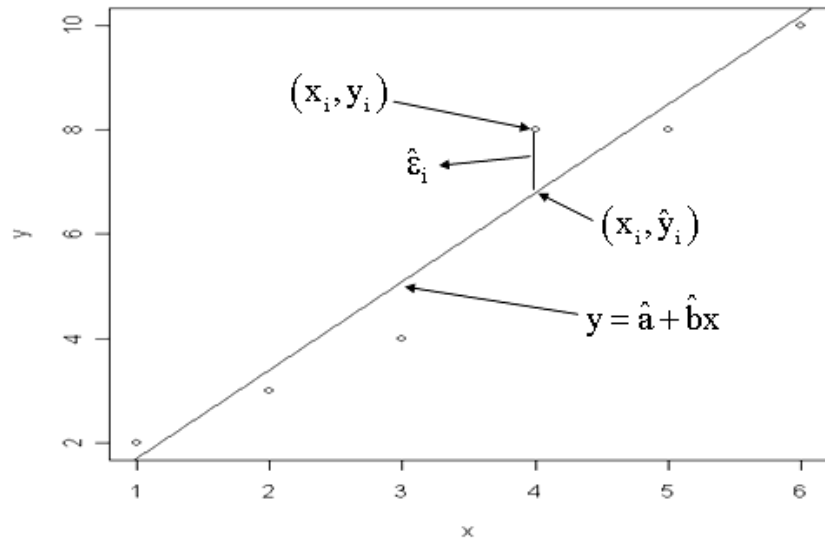
Απλό Γραμμικό Μοντέλο

- Για κάθε παρατήρηση x_i μπορούμε να υπολογίσουμε τις προβλεπόμενες τιμές
$$\hat{y}_i = \hat{a} + \hat{b}x_i.$$
- Αν η ευθεία $\hat{a} + \hat{b}x_i$ δεν περνάει ακριβώς από τα σημεία (y_i, x_i) , περιμένουμε να έχουμε αποκλίσεις μεταξύ των y_i και των \hat{y}_i .
- Οι ποσότητες $\hat{\varepsilon}_i = y_i - \hat{y}_i$ αποτελούν τις εκτιμήσεις των άγνωστων τυχαίων σφαλμάτων ε_i και καλούνται **υπόλοιπα** (*residuals*).

Απλό Γραμμικό Μοντέλο

- Τους συντελεστές του μοντέλου τους εκτιμούμε με βάση τις παρατηρήσεις $(y_1, x_1), \dots, (y_n, x_n)$, εφαρμόζοντας τη μέθοδο ελαχίστων τετραγώνων.
- Με την μέθοδο ελαχίστων τετραγώνων επιλέγουμε την ευθεία (δηλαδή τα \hat{a} και \hat{b}) εκείνη που προσαρμόζεται καλύτερα στα δεδομένα που έχουμε. Πιο συγκεκριμένα επιλέγουμε την ευθεία εκείνη $y = \hat{a} + \hat{b}x$ που ελαχιστοποιεί τα $\hat{\varepsilon}_i$.

Απλό Γραμμικό Μοντέλο



Τα \hat{a} και \hat{b} τέτοια ώστε να

$$\text{ελαχιστοποιείται η συνάρτηση } \sum_{i=1}^n [y_i - (a + bx_i)]^2 = \sum_{i=1}^n \varepsilon_i^2.$$

Απλό Γραμμικό Μοντέλο

- Μετά από πράξεις προκύπτουν τότε οι εκτιμήτριες

$$\hat{b} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \hat{a} = \bar{y} - \hat{b}\bar{x}$$

οι οποίες είναι **τυχαίες μεταβλητές** (από διαφορετικό δείγμα ενδέχεται να προκύψουν διαφορετικές εκτιμήτριες).

Απλό Γραμμικό Μοντέλο

- Το σ^2 το εκτιμούμε από την ποσότητα

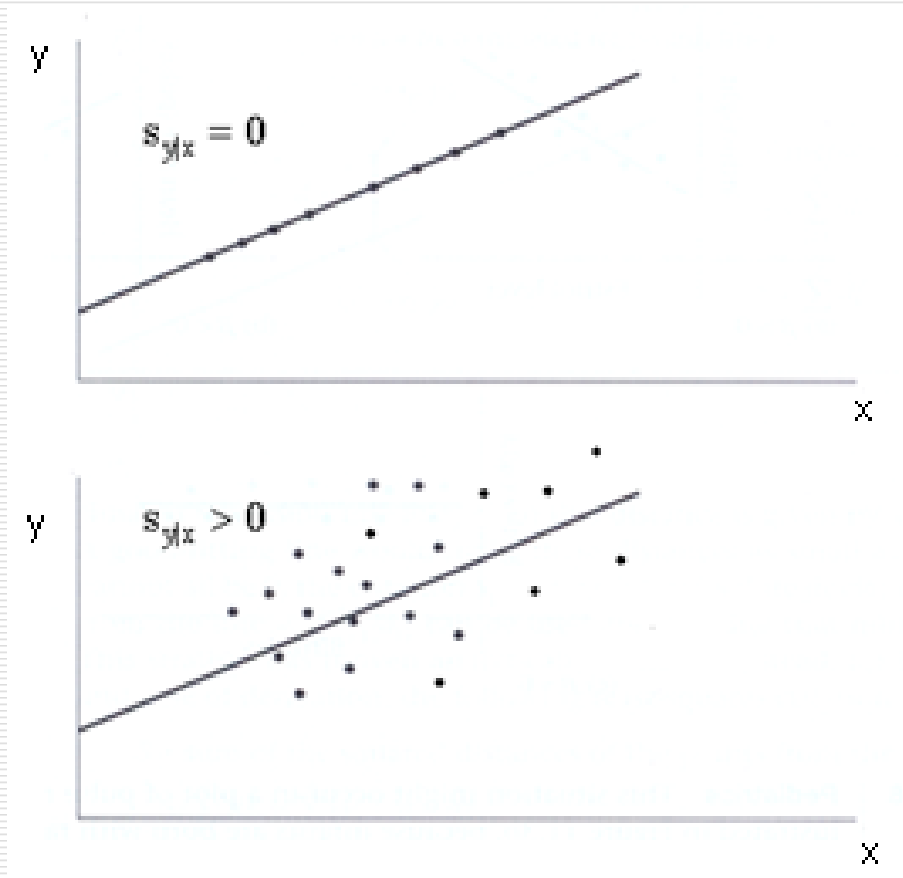
$$s_{y|x}^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \longrightarrow$$

Εκτίμηση διασποράς των σφαλμάτων

Η θετική τετραγωνική της ρίζα της παραπάνω εκτιμήτριας καλείται **τυπικό σφάλμα της παλινδρόμησης** και όσο μικρότερη τιμή έχει τόσο καλύτερη προσαρμογή έχουμε για το μοντέλο παλινδρόμησης.

- Το $s_{y|x}^2$ αποτελεί αμερόληπτη εκτιμήτρια του σ^2 και καλείται **μέσο τετραγωνικό σφάλμα (MSE)**.

Απλό Γραμμικό Μοντέλο



Συντελεστής Συσχέτισης

- Ο **συντελεστής συσχέτισης** (*correlation coefficient*) μεταξύ των τ.μ. X και Y εκφράζει το “βαθμό” στον οποίο μπορούμε να εκτιμήσουμε γραμμικά τη μία τ.μ. όταν γνωρίζουμε την τιμή της άλλης.

$$\rho = \text{Cov}(X, Y) / \{V[X]V[Y]\}^{1/2}$$

- Όταν $\rho=0$ οι τ.μ. X και Y είναι ασυσχέτιστες. Όταν $\rho=1$ υπάρχει τέλεια θετική γραμμική συσχέτιση των δύο τ.μ. ενώ όταν $\rho=-1$ υπάρχει τέλεια αρνητική γραμμική συσχέτιση.

Συντελεστής Συσχέτισης

- Όταν δεν γνωρίζουμε το ρ το εκτιμούμε με την βοήθεια των παρατηρήσεων (y_i, x_i)

$$r = \frac{\sum_1^n (x_i - \bar{x})(y_i - \bar{y})}{\left\{ \sum_1^n (x_i - \bar{x})^2 \sum_1^n (y_i - \bar{y})^2 \right\}^{1/2}}$$

από το **δειγματικό συντελεστή συσχέτισης**.

Συντελεστής Συσχέτισης

- Ο δειγματικός συντελεστής συσχέτισης εκτιμά το βαθμό στον οποίο οι τ.μ. X και Y είναι **γραμμικά συσχετισμένες**, χωρίς να συνεπάγεται κατά ανάγκη κάποιο είδος **αιτιακής σχέσης** μεταξύ των X και Y .
- Αρκετά συχνά μας ενδιαφέρει να ελέγξουμε, σε ε.σ. έστω $\alpha\%$, κατά πόσο οι δύο τ.μ. X και Y είναι ασυσχέτιστες ή όχι, δηλαδή τον έλεγχο $H_0: \rho=0$ με εναλλακτική $H_1: \rho \neq 0$.
- Αποδεικνύεται ότι κάτω από την μηδενική υπόθεση

$$T = \frac{r}{\sqrt{1-r^2}} \sqrt{n-2} \sim \text{St}(n-2).$$

Συντελεστής Συσχέτισης

- Υπολογίζουμε λοιπόν την τιμή του στατιστικού ελέγχου T και η P -τιμή του αμφίπλευρου ελέγχου είναι 2 φορές η πιθανότητα δεξιά του $|T|$ με βάση την $St(n-2)$.
- Στην R μπορούμε να εφαρμόσουμε τον εν λόγω έλεγχο με την βοήθεια της εντολής `cor.test(X,Y)`.
- Αν οι δυο μεταβλητές δεν είναι συνεχείς τότε ο παραπάνω έλεγχος (γνωστός με το όνομα *Pearson correlation coefficient test*) δεν είναι πλέον **έγκυρος**. Πρέπει αντί αυτού να εφαρμοστεί ο μη παραμετρικός *Spearman rank correlation coefficient test*, κατά τον οποίο αντικαθιστούμε τις παρατηρήσεις με την σειρά κατάταξης των τιμών τους (**rank**) και εφαρμόζουμε την προηγούμενη μεθοδολογία.
- Στην R μπορούμε να εφαρμόσουμε τον μη παραμετρικό έλεγχο με την βοήθεια της εντολής `cor.test(X,Y, method="spearman")`.
- Όταν στο δείγμα υπάρχουν **ισοπαλίες** (παρατηρήσεις με την ίδια τιμή) η παραπάνω εντολή μας δίνει προειδοποιητικό μήνυμα και δεν υπολογίζει την ακριβή P -τιμή του ελέγχου αλλά αυτή που προκύπτει από μια προσέγγιση.

Συντελεστής Προσδιορισμού

□ Στο απλό γραμμικό μοντέλο η ποσότητα

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

καλείται **συντελεστής προσδιορισμού**, παίρνει τιμές στο $[0,1]$ και εκφράζει το ποσοστό της διασποράς της τ.μ. Y που εξηγείται με βάση το μοντέλο παλινδρόμησης. Αποδεικνύεται ότι, στο απλό γραμμικό μοντέλο, η παραπάνω ποσότητα ισούται με το τετράγωνο του δειγματικού συντελεστή συσχέτισης r . Γενικά όσο μεγαλύτερες τιμές παίρνει ο συντελεστής προσδιορισμού τόσο ισχυρότερη είναι η γραμμική σχέση εξάρτησης των τ.μ. Y και X , **υπό την προϋπόθεση ότι το γραμμικό μοντέλο είναι το κατάλληλο.**

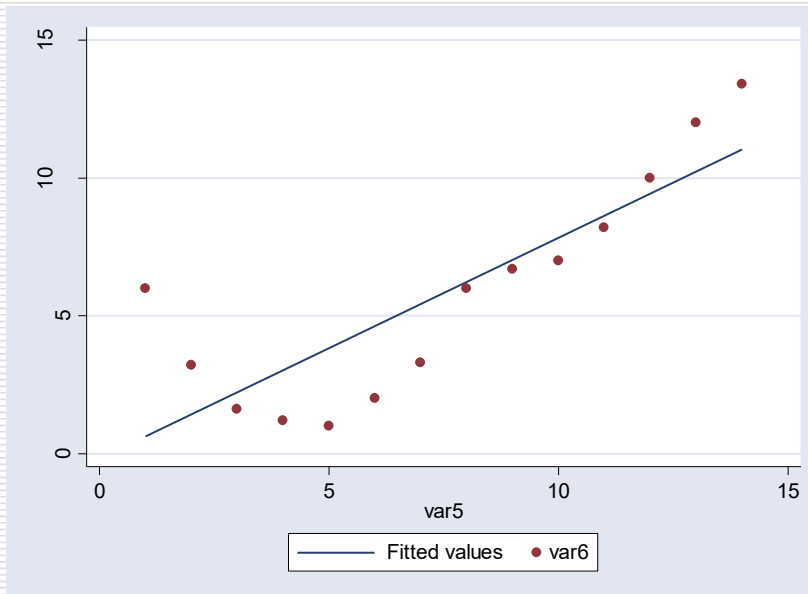
Συντελεστής Προσδιορισμού

- Αρκετές φορές υπολογίζουμε και τον διορθωμένο συντελεστή προσδιορισμού

$$\tilde{R}^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2 / n - 2}{\sum_{i=1}^n (y_i - \bar{y})^2 / n - 1}$$

του οποίου η ερμηνεία δίνεται στο πολλαπλό γραμμικό μοντέλο.

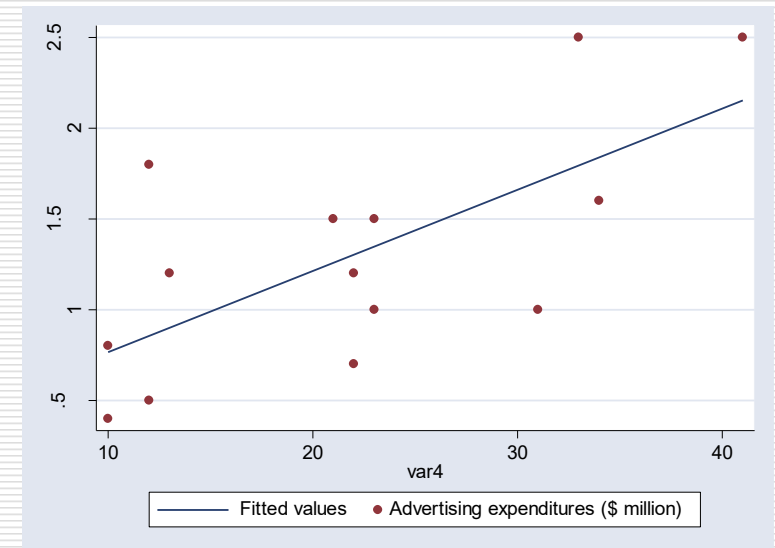
Συντελεστής Προσδιορισμού



Εσφαλμένη υπόθεση γραμμικότητας

$$R^2=0.69$$

$$R^2=0.46$$



Συντελεστής Προσδιορισμού

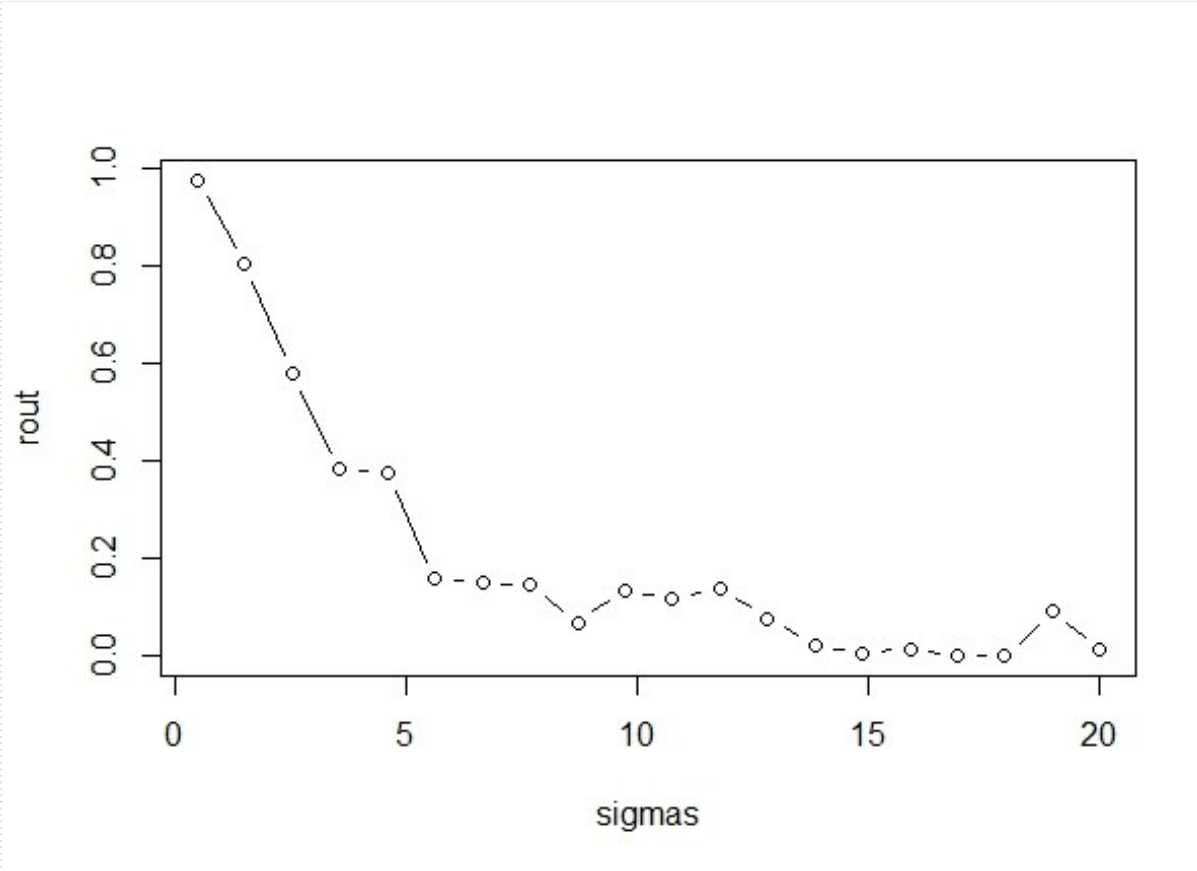
- *Το R^2 δεν είναι κατάλληλος δείκτης καλής προσαρμογής ενός μοντέλου. Μπορεί να είναι πολύ μικρός ακόμα και αν το μοντέλο είναι το σωστό. Στα γραμμικά μοντέλα θεωρώντας με όλες τις προϋποθέσεις να ικανοποιούνται όσο μεγαλώνει το σ^2 τόσο το R^2 μικραίνει.*

Συντελεστής Προσδιορισμού

```
r2.0 <- function(sig){  
  x <- seq(1,10,length.out = 100)      # our predictor  
  y <- 2 + 1.2*x + rnorm(100,0,sd = sig) # our response; a function of x plus  
  #some random noise  
  summary(lm(y ~ x))$r.squared          # print the R-squared value  
}
```

```
sigmas <- seq(0.5,20,length.out = 20)  
rout <- sapply(sigmas, r2.0)            # apply our function to a series of sigma  
values  
plot(rout ~ sigmas, type="b")
```

Συντελεστής Προσδιορισμού



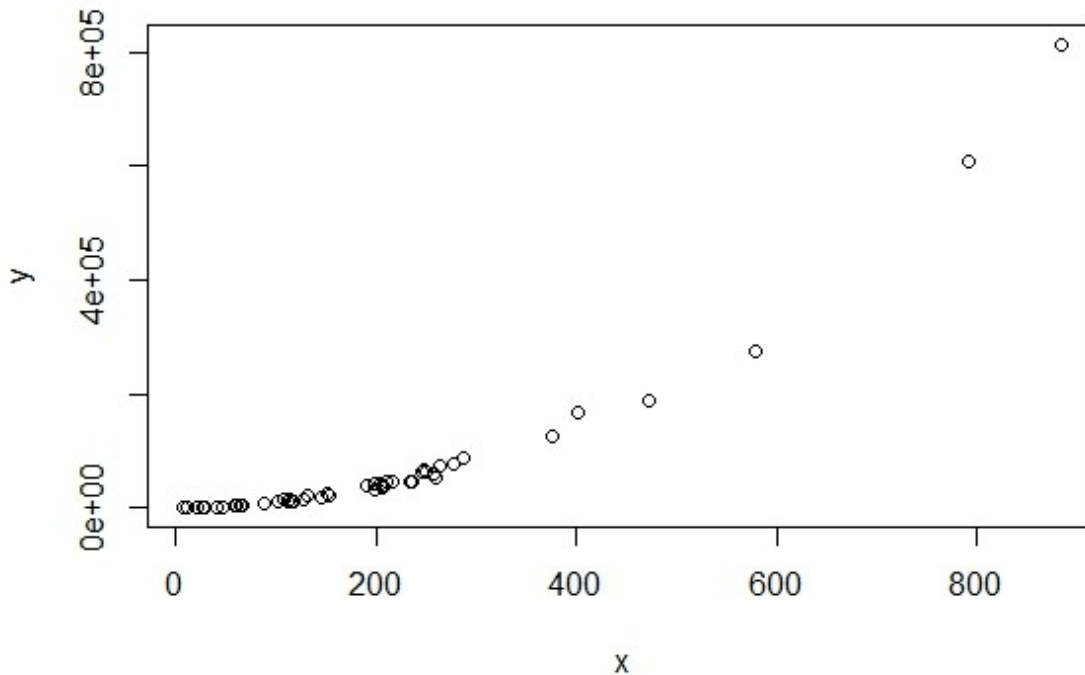
Συντελεστής Προσδιορισμού

- *Το R^2 μπορεί να είναι κοντά στο 1 όταν το μοντέλο είναι εντελώς εσφαλμένο.*

```
set.seed(1)
x <- rexp(50,rate=0.005) # our predictor is data from an exponential
#distribution
y <- (x-1)^2 * runif(50, min=0.8, max=1.2) # non-linear data generation
plot(x,y) # clearly non-linear
```

Συντελεστής Προσδιορισμού

```
summary(lm(y ~ x))$r.squared  
[1] 0.8485146
```



Συντελεστής Προσδιορισμού

- Το R^2 δεν μας λέει κάτι για το σφάλμα πρόβλεψης. Ακόμα και αν το σ^2 είναι ακριβώς το ίδιο και δεν έχουμε διαφορά στις εκτιμήσεις των a και b η τιμή του R^2 διαφοροποιείται αλλάζοντας π.χ. το εύρος τιμών του X
- Για το σφάλμα πρόβλεψης είναι καλύτερα να υπολογίσουμε το **μέσο τετραγωνικό σφάλμα (MSE)**

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Συντελεστής Προσδιορισμού

```
x <- seq(1,10,length.out = 100)
set.seed(1)
y <- 2 + 1.2*x + rnorm(100,0,sd = 0.9)
mod1 <- lm(y ~ x)
summary(mod1)$r.squared
[1] 0.9383379
sum((fitted(mod1) - y)^2)/100 # Mean squared error
[1] 0.6468052
```


Συντελεστής Προσδιορισμού

Επαναλαμβάνουμε τον κώδικα αλλάζοντας μόνο το εύρος τιμών του X.

```
x <- seq(1,2,length.out = 100)      # new range of x
set.seed(1)
y <- 2 + 1.2*x + rnorm(100,0,sd = 0.9)
mod1 <- lm(y ~ x)
summary(mod1)$r.squared
[1] 0.1502448
sum((fitted(mod1) - y)^2)/100      # Mean squared error
[1] 0.6468052
```

Συντελεστής Προσδιορισμού

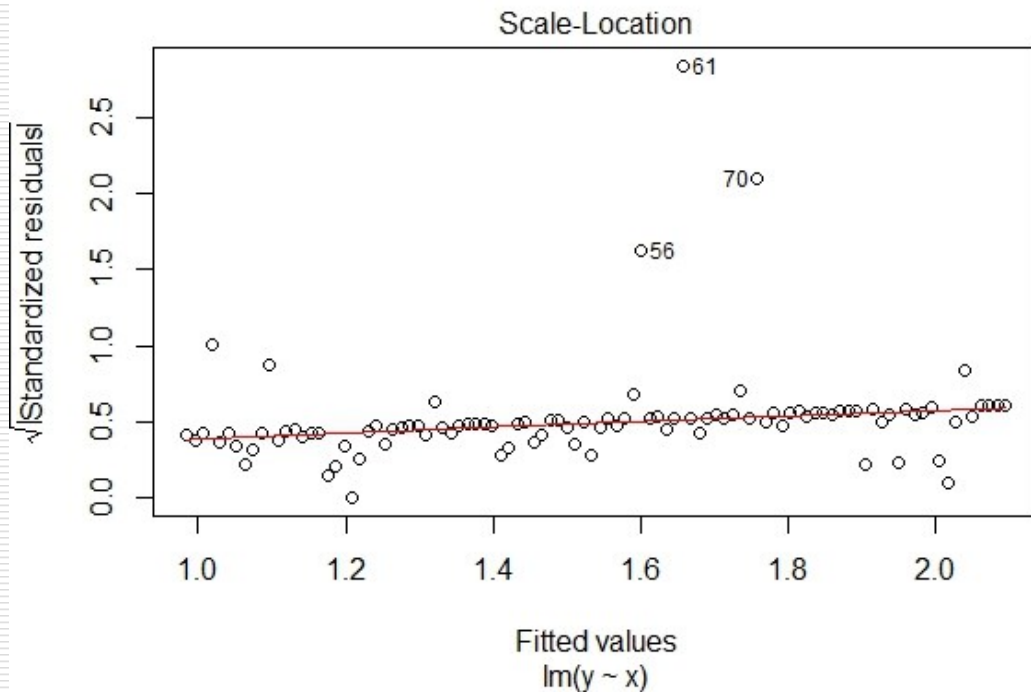
- *Το R^2 δεν πρέπει να χρησιμοποιείται για να συγκρίνουμε δύο μοντέλα στα οποία στο ένα η μεταβλητή απόκρισης είναι Y και στο άλλο ένας μετασχηματισμός της Y .*

Συντελεστής Προσδιορισμού

- Ας δούμε ένα παράδειγμα στο οποίο ένας μετασχηματισμός της Y προφανώς θα βελτιώνει τις προϋποθέσεις.

```
x <- seq(1,2,length.out = 100)
set.seed(1)
y <- exp(-2 - 0.09*x + rnorm(100,0,sd = 2.5))
summary(lm(y ~ x))$r.squared
[1] 0.003281718
```

Συντελεστής Προσδιορισμού

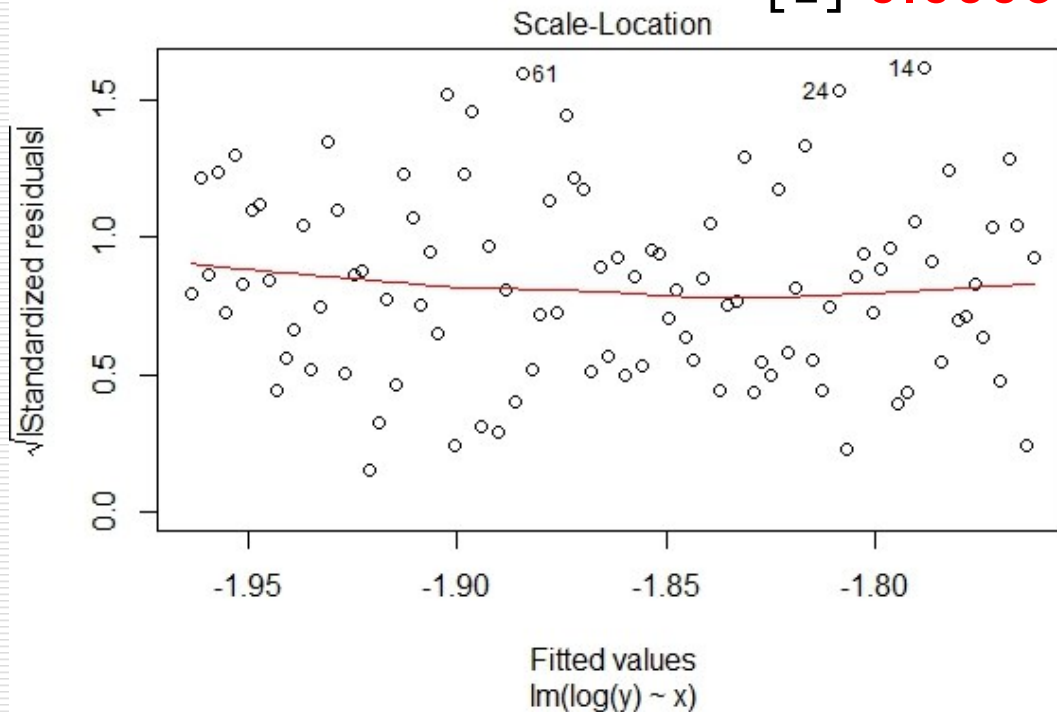


Όπως θα δούμε και παρακάτω από το εν λόγω διάγραμμα καταλήγουμε στο συμπέρασμα ότι έχουμε πρόβλημα ομοσκεδαστικότητας.

Για να επιλύσουμε το εν λόγω πρόβλημα προσαρμόζουμε το γραμμικό μοντέλο με μεταβλητή απόκρισης $\log(Y)$

Συντελεστής Προσδιορισμού

```
summary(lm(log(y) ~ x))$r.squared  
[1] 0.0006921086
```



Μικρότερο από πριν!

Το πρόβλημα
ομοσκεδαστικότητας έχει
λυθεί!

Συντελεστής Προσδιορισμού

- Το R^2 δηλώνει το ποσοστό μεταβλητότητας της Y που εξηγείται από το μοντέλο παλινδρόμησης. Αν αντιστρέψουμε τους ρόλους του Y με το X το R^2 θα παραμείνει το ίδιο! Άρα υψηλή τιμή του R^2 δεν μας λέει τίποτα για το αν μια μεταβλητή εξηγεί μια άλλη (με άλλα λόγια και πάλι συσχέτιση δεν σημαίνει κατά ανάγκη και αιτιακή σχέση).

Συντελεστής Προσδιορισμού

```
x <- seq(1,10,length.out = 100)
y <- 2 + 1.2*x + rnorm(100,0,sd = 2)
summary(lm(y ~ x))$r.squared
[1] 0.7065779
summary(lm(x ~ y))$r.squared
[1] 0.7065779
```

Συμπερασματολογία στο Απλό Γραμμικό Μοντέλο

- Οι εκτιμήσεις των a και b που λαμβάνουμε με την μέθοδο ελαχίστων τετραγώνων βασίζονται στα συγκεκριμένα δεδομένα που διαθέτουμε. Συχνά λοιπόν ενδιαφερόμαστε να ελέγξουμε τις ακόλουθες υποθέσεις, σε ε.σ. έστω a :
 - $H_0: b=0$ έναντι της εναλλακτικής $H_1: b \neq 0$
 - $H_0: a=0$ έναντι της εναλλακτικής $H_1: a \neq 0$

Συμπερασματολογία στο Απλό Γραμμικό Μοντέλο

- Με τον πρώτο έλεγχο θέλουμε να διαπιστώσουμε αν πράγματι αύξηση κατά μια μονάδα της X σημαίνει και μεταβολή της αναμενόμενης τιμής της Y .
- Στο δεύτερο έλεγχο θέλουμε να δούμε κατά πόσο η αναμενόμενη τιμή της Y είναι 0 όταν $X=0$. Πολλές φορές η τιμή αυτή δεν έχει ερμηνεία, διότι η τιμή $X=0$ δεν παρατηρείται ποτέ στην πράξη.

Συμπερασματολογία στο Απλό Γραμμικό Μοντέλο

□ Τα στατιστικά ελέγχου με βάση τις μηδενικές υποθέσεις τότε είναι:

$$T_2 = \frac{\hat{b} - 0}{\text{se}(\hat{b})} = \frac{\hat{b}}{\sqrt{\sigma^2 / \sum_{i=1}^n (x_i - \bar{x})^2}} \approx \frac{\hat{b}}{\sqrt{s_{y|x}^2 / \sum_{i=1}^n (x_i - \bar{x})^2}} \sim \text{St}(n-2)$$

$$T_1 = \frac{\hat{a} - 0}{\text{se}(\hat{a})} = \frac{\hat{a} - 0}{\sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}} \approx \frac{\hat{a}}{\sqrt{s_{y|x}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}} \sim \text{St}(n-2).$$

Συμπερασματολογία στο Απλό Γραμμικό Μοντέλο

- Υπολογίζουμε λοιπόν τα παραπάνω δύο στατιστικά ελέγχου T_1 και T_2 και η P-τιμή των ελέγχων είναι 2 φορές η πιθανότητα της περιοχής της $St(n-2)$ δεξιά από τις τιμές των στατιστικών ελέγχων που παρατηρούμε.
- Ισοδύναμα θα μπορούσαμε να είχαμε κατασκευάσει συμμετρικά $(1-\alpha)\%$ Δ.Ε. για τα a και b και να ελέγξουμε αν η τιμή 0 ανήκει σ' αυτά τα Δ.Ε.

$$\left(\hat{b} \pm t_{n-2, \alpha/2} S_{y|x} / \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \longrightarrow (1-\alpha)\% \text{ Δ.Ε. για το } b$$

$$\left(\hat{a} \pm t_{n-2, \alpha/2} S_{y|x} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \right) \longrightarrow (1-\alpha)\% \text{ Δ.Ε. για το } a$$

Συμπερασματολογία στο Απλό Γραμμικό Μοντέλο

- Ένας άλλος έλεγχος που συνήθως εξετάζουμε στο μοντέλο παλινδρόμησης, γνωστός με την ονομασία **F-test**, είναι και ο παρακάτω, ο οποίος ελέγχει κατά πόσο το προτεινόμενο μοντέλο $y=a+bx$ διαφέρει από το σταθερό $y=a$. Στη απλή γραμμική παλινδρόμηση ο εν λόγω έλεγχος είναι ισοδύναμος με τον έλεγχο για το b που είδαμε πριν.

Συμπερασματολογία στο Απλό Γραμμικό Μοντέλο

- Υπολογίζουμε το στατιστικό ελέγχου

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

το οποίο κάτω από την μηδενική υπόθεση $H_0: b=0$ (με εναλλακτική $H_1: b \neq 0$) ακολουθεί την $F(1, n-2)$. Υπολογίζουμε λοιπόν την τιμή του στατιστικού ελέγχου F και η P -τιμή είναι πιθανότητα της περιοχής της $F(1, n-2)$ δεξιά από το F που παρατηρούμε.

Συμπερασματολογία στο Απλό Γραμμικό Μοντέλο

- Τέλος ένα συμμετρικό $(1-\alpha)\%$ Δ.Ε. για το y όταν $X=x$ είναι το

$$\left(\hat{y} \pm t_{n-2, \alpha/2} S_{y|x} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_1^n (x_i - \bar{x})^2}} \right).$$

- Το παραπάνω διάστημα εμπιστοσύνης καλείται και **διάστημα μέσης πρόβλεψης** (*mean prediction interval*) μιας και στην πραγματικότητα είναι ένα συμμετρικό $(1-\alpha)\%$ Δ.Ε. της τιμής, έστω y , της δεσμευμένης μέσης τιμής της μεταβλητής απόκρισης Y όταν η επεξηγηματική μεταβλητή X ισούται με x , δηλαδή της τιμής της ποσότητας $E[Y|X=x] = a + bx$.

Συμπερασματολογία στο Απλό Γραμμικό Μοντέλο

- Εναλλακτικά, θα μπορούσαμε να κατασκευάζαμε το παρακάτω συμμετρικό $(1-\alpha)\%$ Δ.Ε. για την τιμή, έστω y , της μεταβλητής απόκρισης Y όταν η επεξηγηματική μεταβλητή X ισούται με x είναι το

$$\left(\hat{y} \pm t_{n-2, \alpha/2} S_{y|x} \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_1^n (x_i - \bar{x})^2}} \right).$$

- Το εν λόγω διάστημα, καλείται **διάστημα (ατομικής) πρόβλεψης** (*individual prediction interval*) και αποτελεί ένα συμμετρικό $(1-\alpha)\%$ Δ.Ε. της τιμής, έστω y , της τ.μ. $Y = a + bX + \varepsilon$.

Συμπερασματολογία στο Απλό Γραμμικό Μοντέλο

- Το πρώτο Δ.Ε. παρέχει πληροφορία για τον βαθμό αβεβαιότητας που έχουμε για την εκτίμηση της δεσμευμένης μέσης τιμής $E[Y|X = x]$. Το δεύτερο διάστημα παρέχει πληροφορία για τον βαθμό αβεβαιότητας που έχουμε για την τιμή που θα πάρει η τυχαία μεταβλητή Y όταν $X = x$. Το δεύτερο δηλαδή διάστημα λαμβάνει επιπλέον υπόψιν, πέραν της αβεβαιότητας που έχουμε από την εκτίμηση της δεσμευμένης μέσης τιμής $E[Y|X = x]$ και την μεταβλητότητα της δεσμευμένης κατανομής $Y | (X = x)$. Χρησιμοποιώντας δηλαδή το διάστημα μέσης πρόβλεψης γενικά υποεκτιμούμε την αβεβαιότητά μας για την χρήση της τιμής \hat{y} ως εκτιμήτρια της τιμής που θα πάρει η τυχαία μεταβλητή Y όταν $X = x$.

Συμπερασματολογία στο Απλό Γραμμικό Μοντέλο

- Το πρώτο διάστημα θεωρείται κατάλληλο και χρησιμοποιείται όταν θέλουμε να κατασκευάσουμε διάστημα εμπιστοσύνης για την τιμή, έστω y , της μεταβλητής απόκρισης Y δοσμένης μίας εκ των ήδη παρατηρηθέντων τιμών της επεξηγηματικής μεταβλητής X , για αυτό και λέγεται επίσης και **διάστημα εμπιστοσύνης προσαρμοσμένων (*fitted*) τιμών**. Αντιθέτως αν θέλουμε να χρησιμοποιήσουμε μια μελλοντική παρατήρηση, έστω x , της επεξηγηματικής μεταβλητής X τότε για την κατασκευή του διαστήματος εμπιστοσύνης της τιμής y της μεταβλητής απόκρισης Y χρησιμοποιούμε το **διάστημα (ατομικής) πρόβλεψης**.

Προϋποθέσεις απλού γραμμικού μοντέλου

1. Γραμμικότητα
2. Κανονικότητα Σφαλμάτων
3. Ομοσκεδαστικότητα
4. Ανεξαρτησία Σφαλμάτων

Παράδειγμα απλού γραμμικού μοντέλου στην R

- Μετρήσεις της ποσότητας οξειδίου Y που σχηματίζεται στην επιφάνεια ενός μετάλλου που τίθεται για χρόνο X (min) σε κλίβανο σταθερής θερμοκρασίας δίνεται από τον παρακάτω πίνακα.

x	10	20	30	40	50	60	70	80	90
y	2.0	5.0	6.5	9.5	11.0	13.5	15.0	17.5	19.0

Ζητείται να προσαρμόσουμε ένα μοντέλο απλής γραμμικής παλινδρόμησης και να ελέγξουμε κατά πόσο ο χρόνος X όπου το μέταλλο τίθεται σε κλίβανο σταθερής θερμοκρασίας επηρεάζει την ποσότητα οξειδίου Y που σχηματίζεται στην επιφάνειά του μετάλλου.

Παράδειγμα απλού γραμμικού μοντέλου στην R

```
> x<-seq(10,90,by=10)
```

```
> x
```

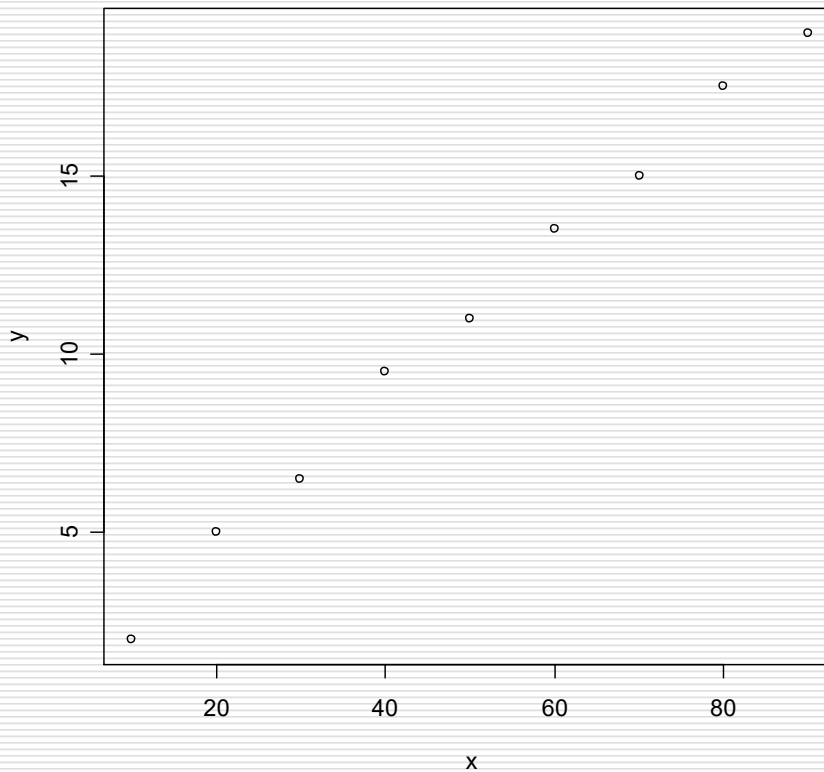
```
[1] 10 20 30 40 50 60 70 80 90
```

```
> y<-c(2,5,6.5,9.5,11,13.5,15,17.5,19)
```

Δημιουργούμε ένα διάγραμμα διασποράς, μια απεικόνιση δηλαδή των σημείων (x_i, y_i) , για να ελέγξουμε αν η γραμμική συνάρτηση φαίνεται να είναι η κατάλληλη.

```
> plot(x,y)
```

Παράδειγμα απλού γραμμικού μοντέλου στην R



Υπόθεση γραμμικότητας λογική

Παράδειγμα απλού γραμμικού μοντέλου στην R

```
> results<-lm(y~x)  
> results
```

Συνάρτηση στην R που προσαρμόζει το γραμμικό μοντέλο με y τα δεδομένα για την μεταβλητή απόκρισης και x τα δεδομένα για την εξηγηματική μεταβλητή.

Call:

```
lm(formula = y ~ x)
```

Coefficients:

(Intercept)

0.4583

x

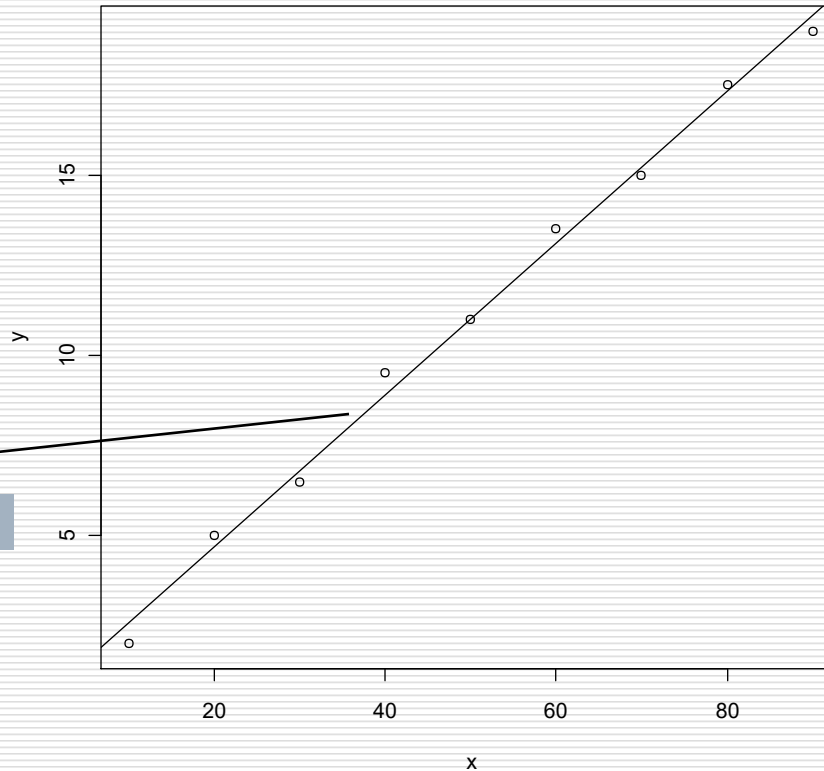
0.2108

\hat{a}

\hat{b}

Παράδειγμα απλού γραμμικού μοντέλου στην R

```
> plot(x,y)  
> abline(results)
```



Ευθεία ελαχίστων τετραγώνων

Παράδειγμα απλού γραμμικού μοντέλου στην R

```
> summary(results)
```

```
Call:
lm(formula = y ~ x)
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-5.667e-01 -2.833e-01  1.559e-16  3.250e-01  6.083e-01
```

Περιγραφικοί δείκτες υπολοίπων

```
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.458333   0.312573   1.466   0.186
x            0.210833   0.005555  37.957 2.29e-09 ***
---
             $\hat{a}$            $se(\hat{a})$         $T_1$ 
```

P-τιμή για τον έλεγχο του a

P-τιμή για τον έλεγχο του b

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
             $\hat{b}$            $se(\hat{b})$         $T_2$ 
```

$S_{y|x}$

```
Residual standard error: 0.4303 on 7 degrees of freedom
Multiple R-squared: 0.9952, Adjusted R-squared: 0.9945
F-statistic: 1441 on 1 and 7 DF, p-value: 2.292e-09
```

διορθωμένος συντελεστής προσδιορισμού

F

P-τιμή για τον F έλεγχο του b (ίδιο με το παραπάνω)

R^2

Παράδειγμα απλού γραμμικού μοντέλου στην R

```
> confint(results)
                2.5 %    97.5 %
(Intercept) -0.2807838 1.1974504
x             0.1976989 0.2239678
```

Συμμετρικά 95% Δ.Ε. για τις παραμέτρους

Συμμετρικό 95% Δ.Ε. για το a

Συμμετρικό 95% Δ.Ε. για το b (δεν περιέχει το 0)

```
> confint(results, level=0.99)
                0.5 %    99.5 %
(Intercept) -0.6355098 1.5521764
x             0.1913952 0.2302714
```

Συμμετρικά 99% Δ.Ε. για τις παραμέτρους

```
> predict(results)
```

Προβλεπόμενες τιμές για κάθε x_i . Το ίδιο αποτέλεσμα προκύπτει και με την εντολή *fitted(results)*, όπως και με την εντολή *results\$fitted*.

```
 1      2      3      4      5      6      7      8
2.566667 4.675000 6.783333 8.891667 11.000000 13.108333 15.216667 17.325000
 9
19.433333
```

Υπόλοιπα. Το ίδιο αποτέλεσμα προκύπτει και με την εντολή *results\$res*.

```
> residuals(results)
 1      2      3      4      5
-5.666667e-01 3.250000e-01 -2.833333e-01 6.083333e-01 1.558947e-16
 6      7      8      9
3.916667e-01 -2.166667e-01 1.750000e-01 -4.333333e-01
```

Παράδειγμα απλού γραμμικού μοντέλου στην R

```
> predict(results, int="c")
  fit   lwr   upr
1 2.566667 1.941342 3.191992
2 4.675000 4.155123 5.194877
3 6.783333 6.354364 7.212303
4 8.891667 8.527990 9.255343
5 11.000000 10.660870 11.339130
6 13.108333 12.744657 13.472010
7 15.216667 14.787697 15.645636
8 17.325000 16.805123 17.844877
9 19.433333 18.808008 20.058658
```

Αναμενόμενες τιμές των y και συμμετρικά 95% Δ.Ε. για κάθε x_i (διαστήματα μέσης πρόβλεψης)

```
> predict(results, list(x=c(15,47)), int="c")
  fit   lwr   upr
1 3.620833 3.049573 4.192094
2 10.367500 10.026088 10.708912
```

Προβλεπόμενες τιμές των y και συμμετρικά 95% Δ.Ε. όταν $x=15$ και $x=47$ (διαστήματα μέσης πρόβλεψης)

```
> predict(results, list(x=c(15,47)), int="c", level=0.99)
  fit   lwr   upr
1 3.620833 2.775406 4.46626
2 10.367500 9.862234 10.87277
```

Προβλεπόμενες τιμές των y και συμμετρικά 99% Δ.Ε. όταν $x=15$ και $x=47$ (διαστήματα μέσης πρόβλεψης)

Παράδειγμα απλού γραμμικού μοντέλου στην R

```
> predict(results, int="p")
```

	fit	lwr	upr
1	2.566667	1.372466	3.760867
2	4.675000	3.532479	5.817521
3	6.783333	5.679205	7.887461
4	8.891667	7.811230	9.972104
5	11.000000	9.927576	12.072424
6	13.108333	12.027896	14.188770
7	15.216667	14.112539	16.320795
8	17.325000	16.182479	18.467521
9	19.433333	18.239133	20.627534

Αναμενόμενες τιμές των y και συμμετρικά 95% Δ.Ε. για κάθε x_i (διαστήματα ατομικής πρόβλεψης)

προειδοποιητικό μήνυμα λάθους διότι κατασκευάζουμε διαστήματα πρόβλεψης για τα παρατηρηθέντα δεδομένα αντί για διαστήματα εμπιστοσύνης των προσαρμοσμένων τιμών.

Warning message:

In predict.lm(results, int = "p") :
Predictions on current data refer to future responses

```
> predict(results, list(x=c(15,47)), int="p", level=0.99)
```

	fit	lwr	upr
1	3.620833	1.894049	5.347618
2	10.367500	8.779315	11.955685

Προβλεπόμενες τιμές των y και συμμετρικά 99% Δ.Ε. όταν $x=15$ και $x=47$ (διαστήματα ατομικής πρόβλεψης)

Παράδειγμα απλού γραμμικού μοντέλου στην R

- Παρατηρούμε λοιπόν ότι

$$\hat{y} = 0.46 + 0.21x.$$

- Ακόμα παρατηρούμε ότι ο συντελεστής προσδιορισμού $R^2 = 0.99$, έχουμε δηλαδή σχεδόν τέλεια προσαρμογή του μοντέλου, ενώ $s_{y|x} = 0.43$ (αρκετά μικρό).

- Τέλος θα μπορούσαμε να είχαμε υπολογίσει τον δειγματικό συντελεστή συσχέτισης με την βοήθεια της εντολής `cor(x,y)`

```
> cor(x,y)
[1] 0.9975795
> cor(x,y)^2
[1] 0.9951648
```

→ =R²

Παράδειγμα απλού γραμμικού μοντέλου στην R

- Στον έλεγχο για το b βλέπουμε ότι η P -τιμή είναι πολύ μικρή οπότε πράγματι αύξηση κατά μια μονάδα της X σημαίνει και μεταβολή της αναμενόμενης τιμής της Y . Πιο συγκεκριμένα αύξηση του χρόνου κατά 1 λεπτό σημαίνει αύξηση της αναμενόμενης ποσότητας οξειδίου κατά 0.21.
- Αντιθέτως η σταθερά δεν φαίνεται από την p -τιμή του αντίστοιχου ελέγχου να είναι στατιστικά διάφορη του μηδενός.
- Για να ελέγξουμε τις προϋποθέσεις του μοντέλου κάνουμε τα εξής:

Παράδειγμα απλού γραμμικού μοντέλου στην R

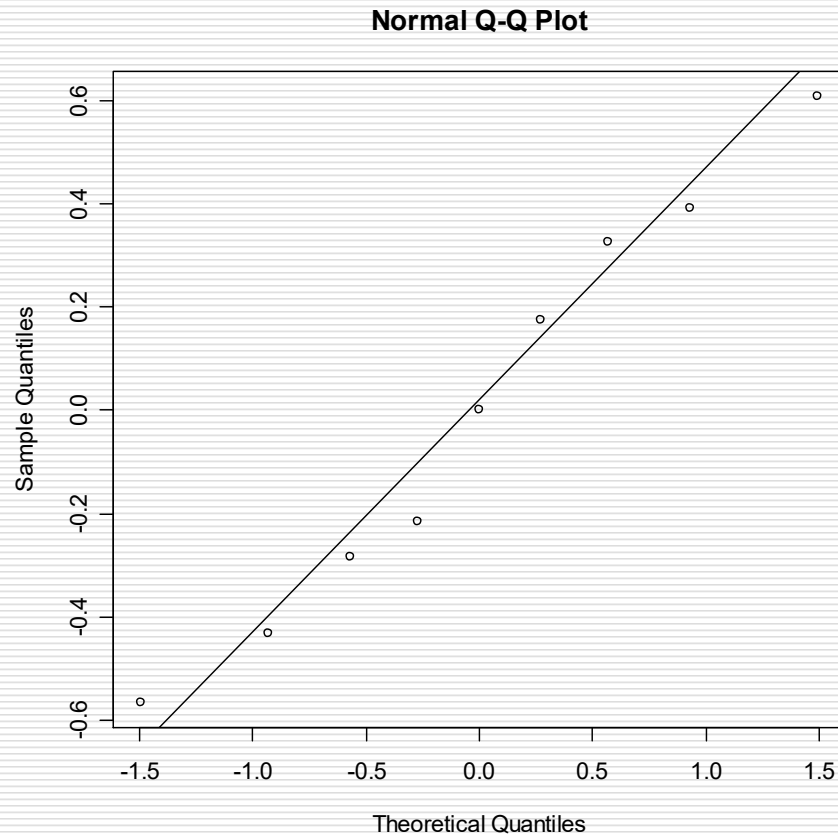
- **Γραμμικότητα.** Ήδη ελέγχθηκε με το διάγραμμα διασποράς.
- **Κανονικότητα σφαλμάτων.** Ιστογράμματα και QQplots για τα υπόλοιπα. Π.χ.

```
> qqnorm(results$res)  
> qqline(results$res)
```



υπόλοιπα

Παράδειγμα απλού γραμμικού μοντέλου στην R



Παράδειγμα απλού γραμμικού μοντέλου στην R

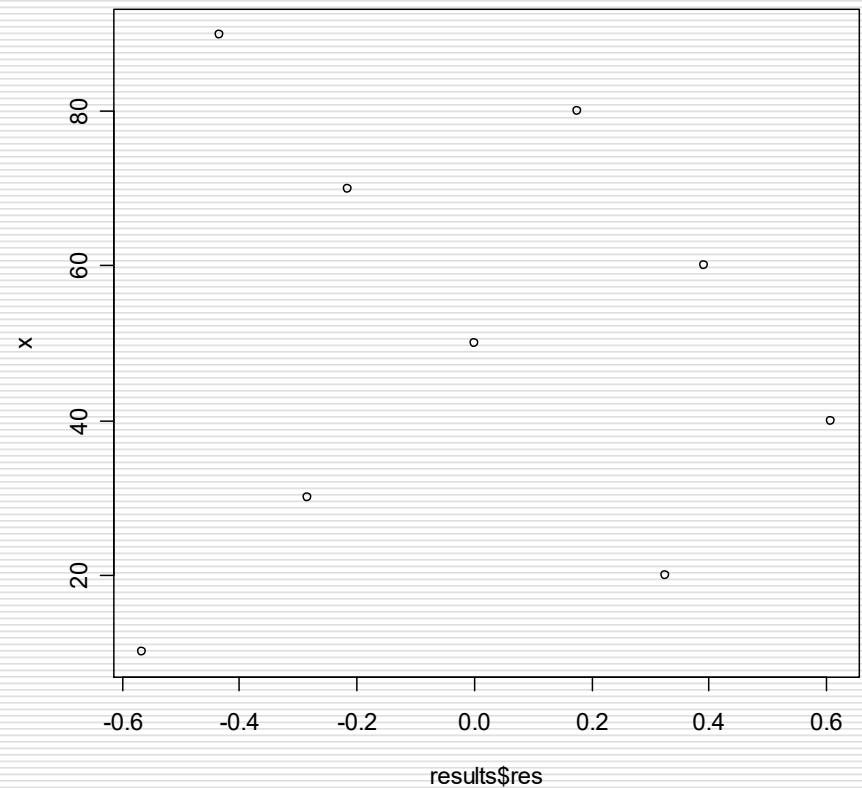
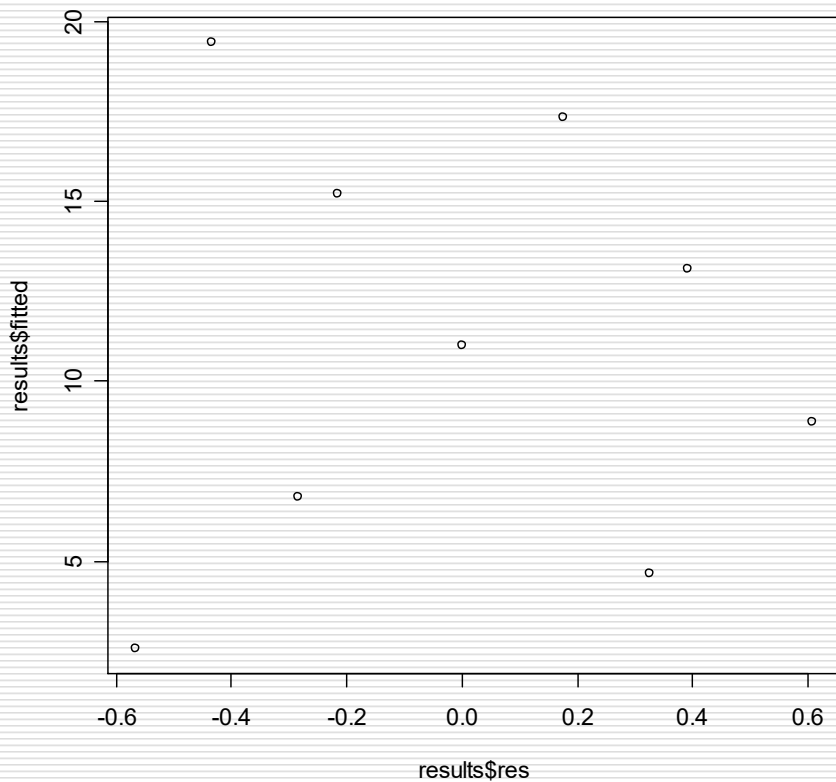
- **Ομοσκεδαστικότητα.** Γραφική παράσταση των υπολοίπων συναρτήσει των προβλεπόμενων τιμών ή συναρτήσει των τιμών της X . Τα ζεύγη αυτών των τιμών δεν πρέπει να εμφανίζουν κάποιο συστηματικό τρόπο συμπεριφοράς.

```
> plot(results$res, results$fitted)  
> plot(results$res, x)
```



προβλεπόμενες τιμές

Παράδειγμα απλού γραμμικού μοντέλου στην R

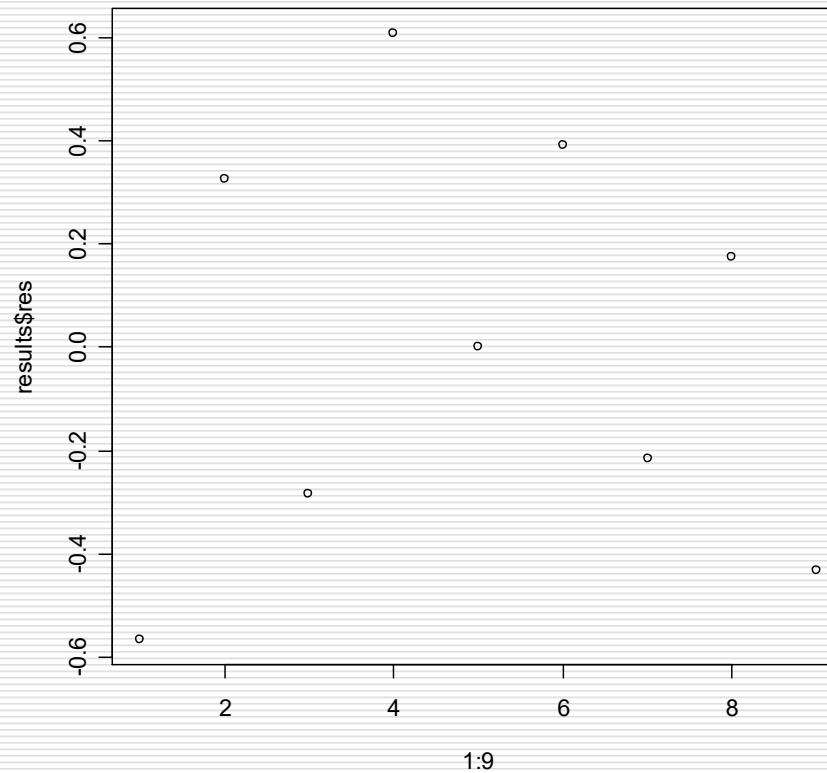


Παράδειγμα απλού γραμμικού μοντέλου στην R

- **Ανεξαρτησία σφαλμάτων.** Κατασκευάζουμε ένα διάγραμμα υπολοίπων σε σχέση με την σειρά των δεδομένων, στο οποίο δεν πρέπει να παρουσιάζεται κάποια σχέση και τα υπόλοιπα να συμπεριφέρονται τυχαία.

```
> plot(1:9,results$res)
```

Παράδειγμα απλού γραμμικού μοντέλου στην R



Παράδειγμα απλού γραμμικού μοντέλου στην R

- Διάφορα διαγνωστικά διαγράμματα μπορούν να γίνουν στην R με την βοήθεια της εντολής

```
> plot(results)
```

Πολλαπλασιαστικό Μοντέλο

- Μη γραμμικότητα.

$$Y = \beta_0 X_1^{\beta_1} \varepsilon$$

Πολλαπλασιαστικό Μοντέλο

Όταν το X_1 είναι ίσο με 1 **μονάδα** η αναμενόμενη τιμή του Y είναι β_0

- Παίρνουμε log και στα δύο μέλη

$$\log(Y) = \log(\beta_0) + \beta_1 \log(X_1) + \log(\varepsilon)$$

Γραμμικό Μοντέλο!!!!!!

$$\varepsilon \sim \text{LogNormal}(0, \sigma^2) \Rightarrow \log(\varepsilon) \sim \text{Normal}(0, \sigma^2)$$

Πολλαπλασιαστικό Μοντέλο

$$(\log(Y) | \mathbf{x}_1) \sim N(\log(\beta_0) + \beta_1 \log(\mathbf{x}_1), \sigma^2)$$

$$(Y | X_1) \sim \text{LN}(\log(\beta_0) + \beta_1 \log(\mathbf{x}_1), \sigma^2)$$

$$E(\log(Y) | \mathbf{x}_1) = \log(\beta_0) + \beta_1 \log(\mathbf{x}_1)$$

$$E(Y | \mathbf{x}_1) = \exp(\log(\beta_0) + \beta_1 \log(\mathbf{x}_1) + \sigma^2 / 2)$$

$$\text{Median}(\log(Y) | \mathbf{x}_1) = \log(\beta_0) + \beta_1 \log(\mathbf{x}_1)$$

$$\text{Median}(Y | \mathbf{x}_1) = \exp(\log(\beta_0) + \beta_1 \log(\mathbf{x}_1))$$

Πολλαπλασιαστικό Μοντέλο

$$\log Y_i = \log \beta_0 + \beta_1 \log X_{1i} + \log \varepsilon_i$$

Λογαριθμίζω τα δεδομένα μου και προσαρμόζω ένα απλό γραμμικό μοντέλο

% μεταβολή στο X

% μεταβολή στο Y

$$\frac{dY}{Y} = \beta_1 \frac{dX_1}{X_1} \Rightarrow \beta_1 = \frac{X_1}{Y} \frac{dY}{dX_1} = \frac{dY / Y}{dX_1 / X_1}$$

1% αύξηση στο X_1 επιφέρει β_1 % μεταβολή στη διάμεσο του Y.

Ο συντελεστής β_1 ονομάζεται η **ελαστικότητα**.

Άλλα Λογαριθμικά Μοντέλα

$$Y = \beta_0 + \beta_1 \log(X_1) + \varepsilon \Leftrightarrow \exp(Y) = \exp(\beta_0) X_1^{\beta_1} \exp(\varepsilon)$$

$$dY = \beta_1 \frac{dX_1}{X_1}$$

Όταν το $\log(X_1)$ αυξηθεί κατά 1 μονάδα (δηλαδή το X_1 πολλαπλασιαστεί με $e \approx 2.72$ ή **ισοδύναμα** το X_1 αυξηθεί κατά 172% (= $100 \times (2.72-1) = 172$)) η αναμενόμενη τιμή του Y θα μεταβληθεί κατά β_1 μονάδες.

Αν το X_1 αυξηθεί κατά 1% η αναμενόμενη τιμή του Y θα μεταβληθεί κατά $\beta_1 \times \log((100+1)/100) = \beta_1 \times \log(1.01)$ μονάδες.

Προσεγγιστικά αν το X_1 αυξηθεί κατά 1% η αναμενόμενη τιμή του Y θα μεταβληθεί κατά $\beta_1/100$ μονάδες.

Όταν το X_1 είναι 1 μονάδα η διάμεσος του $\exp(Y)$ είναι $\exp(\beta_0)$.

Άλλα Λογαριθμικά Μοντέλα

$$\log(Y) = \beta_0 + \beta_1 X_1 + \varepsilon$$

$$\frac{dY}{Y} = \beta_1 dX_1$$

Όταν το X_1 αυξηθεί κατά 1 μονάδα η διάμεσος του $\log(Y)$ θα μεταβληθεί κατά β_1 μονάδες ή **ισοδύναμα** η διάμεσος του Y θα πολλαπλασιαστεί με $\exp(\beta_1)$.

Αν ο συντελεστής β_1 είναι θετικός, **ισοδύναμα** μπορούμε να πούμε ότι όταν το X_1 αυξηθεί κατά 1 μονάδα η διάμεσος του Y θα αυξηθεί κατά $(\exp(\beta_1) - 1) \times 100$ %. Για μικρά β_1 η ποσότητα $(\exp(\beta_1) - 1) \times 100$ % προσεγγίζεται από το $\beta_1 \times 100$ %.

Αν ο συντελεστής β_1 είναι αρνητικός, **ισοδύναμα** μπορούμε να πούμε ότι όταν το X_1 αυξηθεί κατά 1 μονάδα η διάμεσος του Y θα μειωθεί κατά $(1 - \exp(\beta_1)) \times 100$ %.

Όταν το X_1 είναι μηδέν η διάμεσος του Y είναι $\exp(\beta_0)$.

Πολυωνυμική Παλινδρόμηση

- Μη γραμμικότητα.
- Παράδειγμα: Τετραγωνικό Μοντέλο

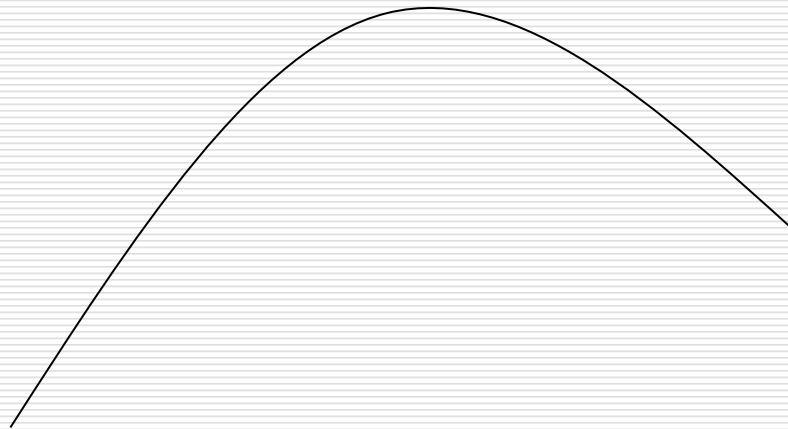
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2).$$

$$z_1 = x_1 \quad \text{and} \quad z_2 = x_1^2$$

$$y_i = \beta_0 + \beta_1 z_{1i} + \beta_2 z_{2i} + \varepsilon_i$$

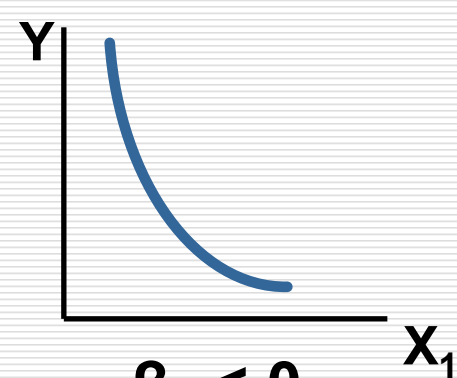
Πολυωνυμική Παλινδρόμηση

- Για παράδειγμα $Y =$ μισθός και $X =$ ηλικία.



Πολυωνυμική Παλινδρόμηση

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{1i}^2 + \varepsilon_i$$



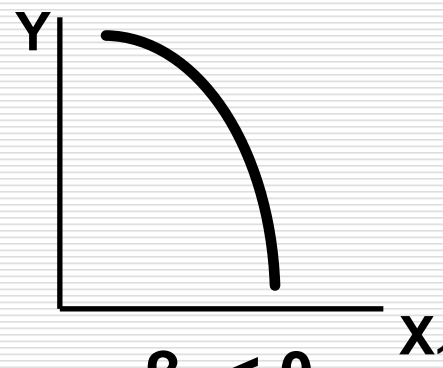
$$\beta_1 < 0$$

$$\beta_2 > 0$$



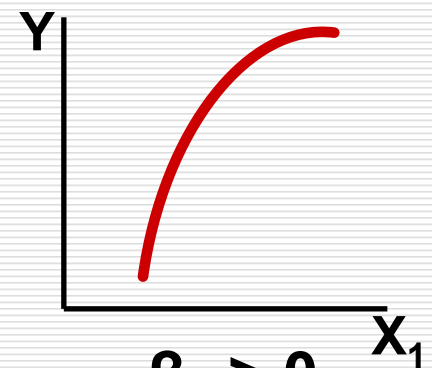
$$\beta_1 > 0$$

$$\beta_2 > 0$$



$$\beta_1 < 0$$

$$\beta_2 < 0$$



$$\beta_1 > 0$$

$$\beta_2 < 0$$

Πολυωνυμική Παλινδρόμηση

$$E(Y|X) = 1.539 + 1.567 X + 0.245 X^2$$

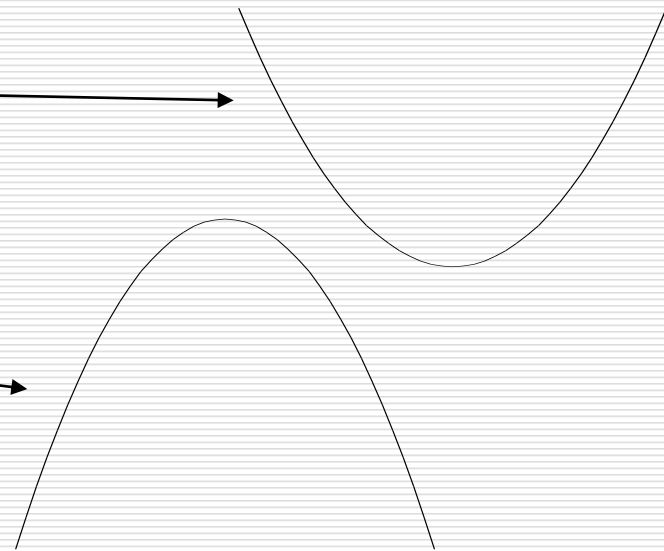
Μέση τιμή του Y
όταν $X = x+1$

Μέση τιμή του Y
όταν $X = x$

➤ Αν αυξηθεί το X κατά μία μονάδα $\Leftrightarrow \mu_{x+1} - \mu_x = 1.567 + 0.245 (2 X + 1)$
[Εξαρτάται από την τιμή του X]

➤ If $b_2 > 0 \Leftrightarrow$ ελάχιστο για $x = -b_1 / (2b_2)$

➤ If $b_2 < 0 \Leftrightarrow$ μέγιστο για $x = -b_1 / (2b_2)$



Τυποποίηση Μεταβλητών

```
Y<-c(64,71,53,67,55,58,77,57,56,51,76,68)
X<-c(57,59,49,62,51,50,55,48,52,42,61,57)
```

```
> mean(Y)
[1] 62.75
> mean(X)
[1] 53.58333
> sd(Y)
[1] 8.9861
> sd(X)
[1] 5.946096
```

```
# Y hourly compensation in euros
# X age
```

```
# 1) Original Variables
```

```
result1<-lm(Y~X)
summary(result1)
```

```
Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-3.6235	14.7524	-0.246	0.8109
X	1.2387	0.2738	4.524	0.0011

**

```
# interpretation
```

If age is increased by 1 year the hourly compensation will increase by 1.23 euros.
If age = 0 (!!!!) then the expected hourly compensation will be -3.62 euros (!!!).

Τυποποίηση Μεταβλητών

2) Centering X

#centering the covariate does not have any effect in b1 but only in b0

```
Xcentered<-X-mean(X)
```

```
result2<-lm(Y~Xcentered)  
summary(result2)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	62.7500	1.5586	40.261	2.14e-12	***
Xcentered	1.2387	0.2738	4.524	0.0011	**

interpretation

If age is increased by 1 year the hourly compensation will increase by 1.23 euros (as before).

If age = 53.6 (mean of age) then the expected hourly compensation will be 62.75 euros.

Τυποποίηση Μεταβλητών

```
# 3) Centering both X & Y
#centering the response and the covariate does not have any effect in b1 but only in b0
(when centering both Y and X, b_0 always will be 0)
```

```
Ycentered<-Y-mean(Y)
result3<-lm(Ycentered~Xcentered)
summary(result3)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2.952e-15	1.559e+00	0.000	1.0000
Xcentered	1.239e+00	2.738e-01	4.524	0.0011 **

```
# interpretation
```

If age is increased by 1 year the hourly compensation will increase by 1.23 euros (as before).

If age = 53.6 (mean of age) then the expected centered hourly compensation will be 0.

Τυποποίηση Μεταβλητών

```
# 4) Centering only Y
```

```
#centering the response does not have any effect in b1 but only in b0
```

```
result4<-lm(Ycentered~X)
summary(result4)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-66.3735	14.7524	-4.499	0.00114	**
X	1.2387	0.2738	4.524	0.00110	**

```
# interpretation
```

If age is increased by 1 year the hourly compensation will increase by 1.23 euros (as before).

If age = 0 (!!!) then the expected centered hourly compensation will be -66.3735 euros (!!!!).

Τυποποίηση Μεταβλητών

```
#5) standardize X  
# interpretation now for b1 is for sd change in X
```

```
Xstd<-(X-mean(X))/sd(X)  
result5<-lm(Y~Xstd)  
summary(result5)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	62.750	1.559	40.261	2.14e-12	***
Xstd	7.365	1.628	4.524	0.0011	**

```
# interpretation
```

If age is increased by 1 sd (i.e. 5.946 years) the hourly compensation will increase by 7.365 euros (this is equal to 5.946×1.23).

If age = 53.6 (mean of age) then the expected hourly compensation will be 62.75 euros.

Τυποποίηση Μεταβλητών

#6) standardize Y

interpretation now for b1 is for sd change in X. New interpretation for b0.

```
Ystd<-(Y-mean(Y))/sd(Y)
result6<-lm(Ystd~X)
summary(result6)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-7.38624	1.64170	-4.499	0.00114	**
X	0.13785	0.03047	4.524	0.00110	**

interpretation

If age is increased by 1 year the hourly compensation will increase by 0.13785 sds (this is equal to $8.98 \times 0.13785 = 1.23$ euros).

If age = 0 years (!!!) then the expected hourly compensation will be (-7.28624 sds + mean(Y)) (this is equal to $-7.28624 \times 8.98 + 62.75$).

Τυποποίηση Μεταβλητών

```
#7) standardize Y & X  
# interpretation now for b1 is for sd change in X.  
# b0 will be zero now.
```

```
result7<-lm(Ystd~Xstd)  
summary(result7)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	3.150e-16	1.734e-01	0.000	1.0000
Xstd	8.196e-01	1.812e-01	4.524	0.0011 **

```
# interpretation
```

If age is increased by 1 sd (i.e. 5.946 years) the hourly compensation will increase by 0.8196 sds (this is equal to $8.98 \times 0.8196 = 7.36$ euros).

If age = 53.6 (mean of age) then the expected standardized hourly compensation will be 0.

Κανονικοποίηση (Normalization)

- ❑ Οι νέες τιμές είναι στο $[0,1]$.
- ❑ Χρήσιμο αν οι μεταβλητές Y και X έχουν πολύ διαφορετικό εύρος τιμών, πχ. το X παίρνει τιμές στο $[0,100]$ και το Y στο $[0, 100000]$.
- ❑ Μειώνει την επίδραση των έκτροπων παρατηρήσεων.
- ❑ Αν χρησιμοποιείς προσεγγιστική μέθοδο ελαχιστοποίησης, όπως τον `gradient descent algorithm`, θα συγκλίνει πιο γρήγορα.

```
> summary(lm(y~x)) # original data (διαφορετικά από πριν)
```

```
Residuals:
```

```
   Min       1Q   Median       3Q      Max
-2.82418 -0.72666 -0.07032  0.65312  2.60991
```

```
Coefficients:
```

```
      Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.9719     0.1058  18.64 <2e-16 ***
x            3.2492     0.1192  27.26 <2e-16 ***
---

```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 1.055 on 98 degrees of freedom
```

```
Multiple R-squared:  0.8835,    Adjusted R-squared:  0.8823
```

```
F-statistic: 743.1 on 1 and 98 DF,  p-value: < 2.2e-16
```

$$dY = \beta dX \Rightarrow \beta = \frac{dY}{dX}$$

Όταν το X αυξηθεί κατά μία μονάδα η αναμενόμενη τιμή του Y θα αυξηθεί κατά 3.2492 μονάδες.

Κανονικοποίηση (Normalization)

$$Y_{\text{new}} = \frac{Y - \min(Y)}{\max(Y) - \min(Y)}$$

$$X_{\text{new}} = \frac{X - \min(X)}{\max(X) - \min(X)}$$

$$\frac{dY}{\max(Y) - \min(Y)} = \beta \frac{dX}{\max(X) - \min(X)}$$

$$\beta = \frac{\max(X) - \min(X)}{\max(Y) - \min(Y)} \frac{dY}{dX}$$

```

> xnew<-(x-min(x))/(max(x)-min(x))
> ynew<-(y-min(y))/(max(y)-min(y))
> min(y)
[1] -4.286472
> min(x)
[1] -1.775145
> max(x)-min(x)
[1] 4.458834
> max(y)-min(y)
[1] 14.1019
> (max(y)-min(y))/(max(x)-min(x))
[1] 3.162688
> summary(lm(ynew~xnew)) # normalized data
Residuals:
    Min       1Q   Median       3Q      Max
-0.200270 -0.051529 -0.004987  0.046314  0.185075

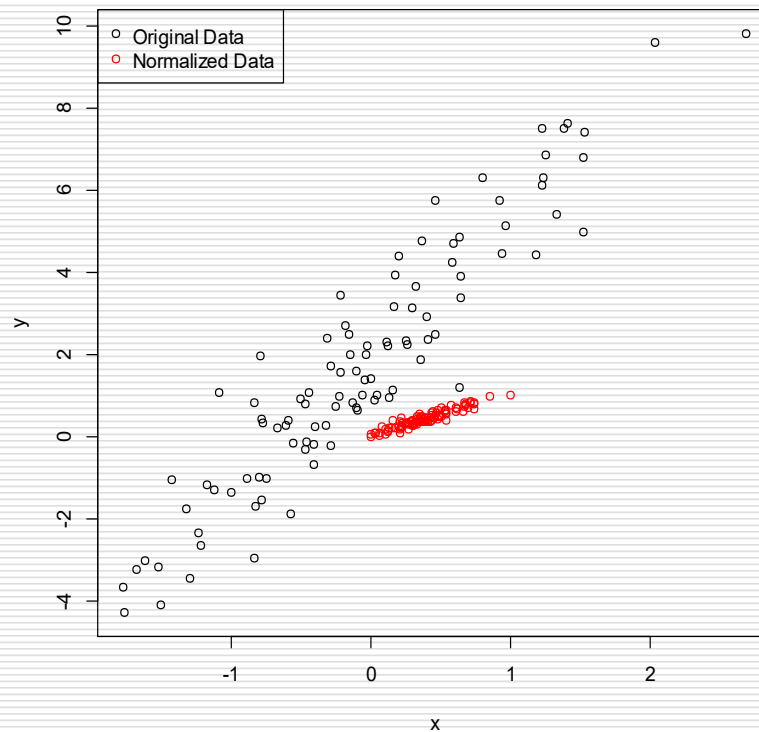
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.03479   0.01625   2.141  0.0347 *
xnew         1.02735   0.03769  27.260 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.07479 on 98 degrees of freedom
Multiple R-squared:  0.8835,    Adjusted R-squared:  0.8823
F-statistic: 743.1 on 1 and 98 DF, p-value: < 2.2e-16

> 1.02735*(max(y)-min(y))/(max(x)-min(x))
[1] 3.2492
    
```

- Όταν το X αυξηθεί κατά μία μονάδα η αναμενόμενη τιμή του Y θα αυξηθεί κατά $1.02735 * (\max(y) - \min(y)) / (\max(x) - \min(x)) = 3.2492$ μονάδες (όπως πριν).
- Όταν το X = min(X) = -1.775145 η αναμενόμενη τιμή του Y_{new} είναι 0.03479 ή η αναμενόμενη τιμή του Y είναι $0.03479 * (\max(y) - \min(y)) + \min(y) = -3.795867$ (= 1.9719 + 3.2492 * (-1.775145)), χρησιμοποιώντας του συντελεστές του μοντέλου με τις αρχικές τιμές).

Κανονικοποίηση (Normalization)



Πολλαπλό Γραμμικό Μοντέλο

- Ας θεωρήσουμε τώρα ότι διαθέτουμε p επεξηγηματικές μεταβλητές $\mathbf{X}=(X_1, \dots, X_p)$, έστω όλες **ποσοτικές**, οι οποίες συνδέονται πάλι γραμμικά με την **ποσοτική** μεταβλητή απόκρισης Y .

$$Y = a + b_1 X_1 + \dots + b_p X_p + \varepsilon, \varepsilon \sim N(0, \sigma^2) \Leftrightarrow$$

$$E(Y | X_1, \dots, X_p) = a + b_1 X_1 + \dots + b_p X_p$$



Πολλαπλό γραμμικό
μοντέλο παλινδρόμησης

Πολλαπλό Γραμμικό Μοντέλο

- Έστω $(Y_1, X_{11}, \dots, X_{1p}), \dots, (Y_n, X_{n1}, \dots, X_{np})$ τυχαίο δείγμα. Τότε ισοδύναμα:

$$Y_i = a + b_1 X_{i1} + \dots + b_p X_{ip} + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2) \Leftrightarrow$$

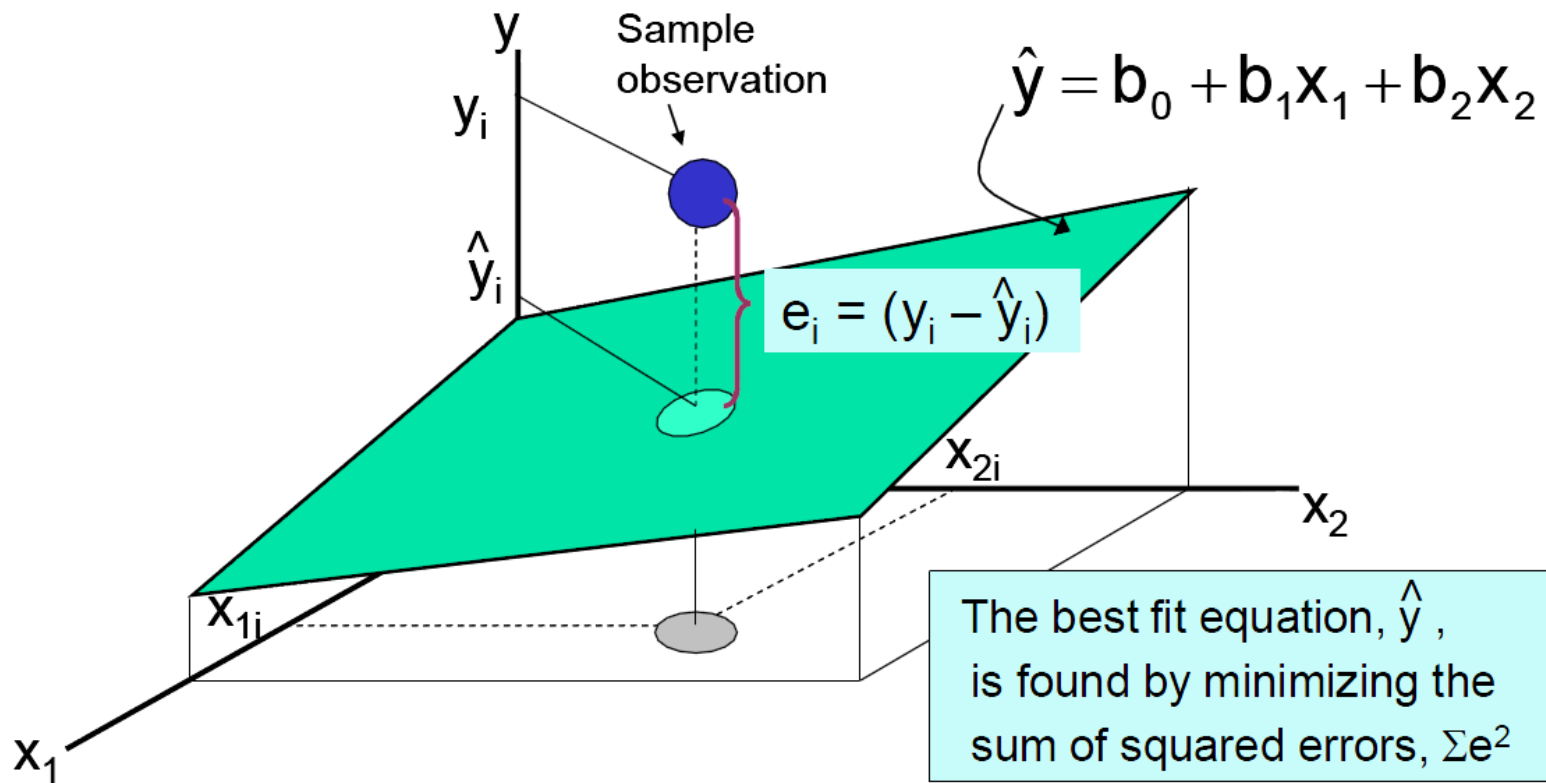
$$E(Y_i | X_{i1}, \dots, X_{ip}) = a + b_1 X_{i1} + \dots + b_p X_{ip}$$

όπου ε_i ανεξάρτητα.

- Με την βοήθεια των παρατηρήσεων $(Y_1, X_{11}, \dots, X_{1p}), \dots, (Y_n, X_{n1}, \dots, X_{np})$ και με βάση την μέθοδο ελαχίστων τετραγώνων παίρνουμε πάλι εκτιμήτριες για τις παραμέτρους του μοντέλου $a, b_1, \dots, b_p, \sigma^2$.

Πολλαπλό Γραμμικό Μοντέλο

Two variable model



Ερμηνεία Παραμέτρων

- Η παράμετρος a εκφράζει την μέση τιμή της τ.μ. Y όταν όλα τα X_j , $j=1, \dots, p$, είναι μηδέν.
- Η παράμετρος b_j , $j=1, \dots, p$, εκφράζει την αναμενόμενη μεταβολή της τιμή της τ.μ. Y όταν η X_j αυξηθεί κατά μία μονάδα και οι υπόλοιπες X_k , $k \neq j$, παραμείνουν σταθερές.
- Η ποσότητα σ^2 εκφράζει την διασπορά των σφαλμάτων, την οποία θεωρούμε σταθερή ανεξάρτητα των τιμών των τ.μ. X_1, \dots, X_p (**υπόθεση ομοσκεδαστικότητας**). Επειδή η τυχαιότητα της Y δεδομένου των τιμών των X οφείλεται στα σφάλματα, το σ^2 εκφράζει και την διασπορά της δεσμευμένης κατανομής της τ.μ. $Y|\mathbf{X}$.

Ερμηνείες συντελεστών

- Οι ερμηνείες είναι ένα από τα πιο σημαντικά κομμάτια της στατιστικής ανάλυσης. Είναι η επικοινωνία σας με αυτόν που σας έδωσε το πρόβλημα, ο οποίος συνήθως δεν γνωρίζει στατιστική.
- Θα πρέπει να είναι **απλές και κατανοητές** για κάποιον που δεν γνωρίζει καν ίσως το πρόβλημα που έχετε και ούτε καν τα δεδομένα που διαθέτετε.

Ερμηνείες συντελεστών

- Στην πολλαπλή γραμμική παλινδρόμηση ερμηνεύετε κάθε εκτιμώμενη τιμή των συντελεστών a, b_1, \dots, b_p ξεχωριστά σε μία απλή πρόταση.
- Μην χρησιμοποιείται σύμβολα ή κωδικοποιήσεις ή συντομεύσεις για τις μεταβλητές αλλά τα πραγματικά τους ονόματα στα Ελληνικά.
- Χρησιμοποιήστε πάντα μονάδες μέτρησης (όποτε διαθέτετε).
- Αναφέρετε, με τα πραγματικά τους ονόματα, τις υπόλοιπες μεταβλητές που μένουν σταθερές λεπτομερώς (μην λέτε δηλαδή ότι οι υπόλοιπες μεταβλητές μένουν σταθερές, αυτός που το διαβάζει ίσως να μην ξέρει ποιες είναι).

Πολλαπλό Γραμμικό Μοντέλο με πίνακες

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \quad \tilde{\mathbf{X}} = \begin{pmatrix} 1 & X_{11} & \dots & X_{1p} \\ \vdots & X_{21} & \ddots & \vdots \\ \vdots & \vdots & & \vdots \\ 1 & X_{n1} & \dots & X_{np} \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} a \\ b_1 \\ \vdots \\ b_p \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\mathbf{Y} = \tilde{\mathbf{X}}\mathbf{b} + \boldsymbol{\varepsilon}$$

- Η εκτιμήτριες με την μέθοδο ελαχίστων τετραγώνων είναι:

$$\hat{\mathbf{b}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{Y} \quad \text{και} \quad s_{y|x_1, \dots, x_p}^2 = \frac{1}{n-p-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

οι οποίες είναι **τυχαίες μεταβλητές** (από διαφορετικό δείγμα ενδέχεται να προκύψουν διαφορετικές εκτιμήτριες).

Πολλαπλό Γραμμικό Μοντέλο

- Εκτιμώντας λοιπόν το \mathbf{b} από το $\hat{\mathbf{b}}$ καταλήγουμε στο

$$\hat{\mathbf{Y}} = \tilde{\mathbf{X}}\hat{\mathbf{b}}.$$

- Το $\hat{\mathbf{Y}}$ καλείται **προβλεπόμενη τιμή** και είναι η αναμενόμενη που θα πάρει η τ.μ. Y όταν $\mathbf{X}=\mathbf{x}$, όπως αυτήν την εκτιμήσαμε με βάση το μοντέλο παλινδρόμησης. Η **προβλεπόμενη τιμή είναι τ.μ., δηλαδή για διαφορετικό δείγμα ενδέχεται να πάρει άλλη τιμή όταν $\mathbf{X}=\mathbf{x}$.** Η προβλεπόμενη τιμή αποτελεί μία αμερόληπτη εκτιμήτρια της άγνωστης τιμής y που παίρνει η τ.μ. Y όταν $\mathbf{X}=\mathbf{x}$. Παρακάτω θα δούμε ένα Δ.Ε. για την τιμή αυτή y όταν $\mathbf{X}=\mathbf{x}$.

Πολλαπλό Γραμμικό Μοντέλο

- Για κάθε x_{i1}, \dots, x_{ip} μπορούμε να υπολογίσουμε τις **προβλεπόμενες τιμές**

$$\hat{y}_i = \hat{a} + \hat{b}_1 x_{i1} + \dots + \hat{b}_p x_{ip}.$$

- Οι ποσότητες $\hat{\varepsilon}_i = y_i - \hat{y}_i$ αποτελούν τις εκτιμήσεις των σφαλμάτων των μετρήσεων και καλούνται όπως και πριν **υπόλοιπα**.

Συντελεστής πολλαπλού προσδιορισμού

- Στο πολλαπλό γραμμικό μοντέλο η ποσότητα
$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

καλείται **πολλαπλός συντελεστής προσδιορισμού**, παίρνει τιμές στο $[0,1]$ και εκφράζει το ποσοστό της διασποράς της τ.μ. Y που εξηγείται με βάση το μοντέλο παλινδρόμησης.

- Κάθε φορά που προσθέτουμε μια επεξηγηματική μεταβλητή στο μοντέλο ο πολλαπλός συντελεστής προσδιορισμού αυξάνεται. Όταν δεν θέλουμε να υπάρχει αυτή η εξάρτηση από το p υπολογίζουμε τον **διορθωμένο συντελεστής προσδιορισμού**

$$\tilde{R}^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2 / n - p - 1}{\sum_{i=1}^n (y_i - \bar{y})^2 / n - 1}$$

Διορθωμένος Συντελεστής Προσδιορισμού

- Ο διορθωμένος συντελεστής προσδιορισμού συνδέεται με το R^2 με βάση τον τύπο:

$$\tilde{R}^2 = 1 - \frac{n-1}{n-p-1} R^2$$

- Μπορεί να αποδειχθεί ότι πάντα $\tilde{R}^2 \leq R^2 \leq 1$.
- Ο διορθωμένος συντελεστής είναι **αρνητικός** αν και μόνο αν $R^2 \leq \frac{p}{n-1}$.
- Η παραπάνω ανισότητα ισχύει πάντα όταν $p+1 > n$, ενώ μπορεί και να ισχύει όταν $p+1 < n$ και το R^2 είναι πολύ μικρό.

Συμπερασματολογία στο πολλαπλό γραμμικό μοντέλο

- Ενδιαφερόμαστε όπως και πριν για τους εξής ελέγχους υποθέσεων:
 - $H_0: a=0$ έναντι της εναλλακτικής $H_1: a \neq 0$
 - $H_0: b_1=0$ έναντι της εναλλακτικής $H_1: b_1 \neq 0$
 - · · ·
 - · · ·
 - · · ·
 - $H_0: b_p=0$ έναντι της εναλλακτικής $H_1: b_p \neq 0$
- Επίσης ενδιαφέρον έχει και το F-test το οποίο ελέγχει την μηδενική υπόθεση $H_0: b_1 = b_2 = \dots = b_p = 0$, με εναλλακτική ότι τουλάχιστον ένα από τα b_j δεν είναι 0.

Συμπερασματολογία στο πολλαπλό γραμμικό μοντέλο

- Τέλος ένα συμμετρικό $(1-\alpha)\%$ **διάστημα μέσης πρόβλεψης** για την αναμενόμενη τιμή, έστω y , της μεταβλητής απόκρισης Y όταν η επεξηγηματική μεταβλητή \mathbf{X} ισούται με $\mathbf{x} = (x_1, \dots, x_p)^T$ είναι το

$$\left(\hat{y} \pm t_{n-p-1, \alpha/2} S_{y|x} \sqrt{\tilde{\mathbf{x}}^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{x}}} \right), \text{ όπου } \tilde{\mathbf{x}} = (1, x_1, \dots, x_p)^T.$$

Συμπερασματολογία στο πολλαπλό γραμμικό μοντέλο

□ Ένα συμμετρικό $(1-\alpha)\%$ **διάστημα (ατομικής) πρόβλεψης** για την προβλεπόμενη τιμή, έστω y , της μεταβλητής απόκρισης Y όταν η επεξηγηματική μεταβλητή \mathbf{X} ισούται με $\mathbf{x} = (x_1, \dots, x_p)^T$ είναι το

$$\left(\hat{y} \pm t_{n-p-1, \alpha/2} S_{y|x} \sqrt{1 + \tilde{\mathbf{x}}^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{x}}} \right), \text{ όπου } \tilde{\mathbf{x}} = (1, x_1, \dots, x_p)^T.$$

Προϋποθέσεις πολλαπλού γραμμικού μοντέλου

1. Γραμμικότητα
2. Κανονικότητα Σφαλμάτων
3. Ομοσκεδαστικότητα
4. Ανεξαρτησία Σφαλμάτων

Προϋποθέσεις πολλαπλού γραμμικού μοντέλου

- **Γραμμικότητα.** Στην απλή γραμμική παλινδρόμηση ο έλεγχος της γραμμικότητας γινόταν με το διάγραμμα διασποράς των (y_i, x_i) . Στην πολλαπλή παλινδρόμηση, με p επεξηγηματικές μεταβλητές, θα μπορούσαμε να δημιουργήσουμε p διαφορετικά διαγράμματα διασποράς, ένα για κάθε επεξηγηματική μεταβλητή, και να εξετάσουμε την υπόθεση της γραμμικότητας. Βέβαια με τον τρόπο αυτόν δεν ελέγχουμε την εγκυρότητα του γενικού γραμμικού μοντέλου αλλά ελέγχουμε την εγκυρότητα των παρακάτω απλών γραμμικών μοντέλων:

$$E[Y|x_j] = a + b_j x_j, \quad j = 1, \dots, p.$$

Προϋποθέσεις πολλαπλού γραμμικού μοντέλου

- Αν οι επεξηγηματικές μεταβλητές είναι **ασυσχέτιστες**, τότε δεν υπάρχει διαφορά από το να ελέγξουμε τα παραπάνω μοντέλα ή το πολλαπλό γραμμικό μοντέλο. Αλλά στις περισσότερες περιπτώσεις οι επεξηγηματικές μεταβλητές συσχετίζονται. Τότε αυτό που πρέπει να ελέγξουμε είναι αν πράγματι η επεξηγηματική μεταβλητή X_j , $j = 1, \dots, p$, συνδέεται γραμμικά με την δεσμευμένη μέση τιμή της Y , αν όλες οι άλλες επεξηγηματικές τ.μ. $X_1, \dots, X_{j-1}, X_{j+1}, \dots, X_p$ συνδέονται γραμμικά με την δεσμευμένη μέση τιμή της Y . Με άλλα λόγια αν θεωρήσουμε την σχέση

$$y_i \approx \hat{a} + \hat{b}_1 x_{i1} + \dots + \hat{b}_{j-1} x_{i(j-1)} + p_j(x_{ij}) + \hat{b}_{j+1} x_{i(j+1)} + \dots + \hat{b}_p x_{ip} \quad (i = 1, \dots, n),$$

χρησιμοποιώντας τις εκτιμήτριες ελαχίστων τετραγώνων που έχουμε από το πολλαπλό γραμμικό μοντέλο, αρκεί να δείξουμε ότι η συνάρτηση p_j είναι γραμμική.

Προϋποθέσεις πολλαπλού γραμμικού μοντέλου

- Αν τώρα θεωρήσουμε την σχέση που μας δίνει τα υπόλοιπα

$$y_i = \hat{a} + \hat{b}_1 x_1 + \hat{b}_2 x_2 + \dots + \hat{b}_p x_p + \hat{\varepsilon}_i \quad (i = 1, \dots, n)$$

και την αφαιρέσουμε από την σχέση στην προηγούμενη διαφάνεια έχουμε

$$p_j(x_{ij}) \approx \hat{b}_j x_{ij} + \hat{\varepsilon}_i \equiv P_{ij} \quad (i = 1, \dots, n).$$

Προϋποθέσεις πολλαπλού γραμμικού μοντέλου

- Οι όροι P_{ij} καλούνται **j -μερικά υπόλοιπα** (*partial residuals*). Από την τελευταία σχέση είναι εμφανές ότι μπορούμε να ελέγξουμε αν η συνάρτηση p_j είναι γραμμική, με το διάγραμμα διασποράς των σημείων (x_{ij}, P_{ij}) , $i=1, \dots, n$. Επαναλαμβάνοντας την διαδικασία για κάθε $j=1, \dots, p$ ελέγχουμε την γραμμικότητα στο πολλαπλό γραμμικό μοντέλο.
- Οι υπόλοιπες προϋποθέσεις ελέγχονται όπως και στο απλό γραμμικό μοντέλο.

Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

- Η αντοχή Y ξύλινων δοκών, σε σχέση με την περιεκτικότητα τους σε νερό X_1 και το ειδικό τους βάρος X_2 , αποτέλεσε το αντικείμενο μιας μελέτης. Με βάση τα παρακάτω δεδομένα, να προσαρμοστεί το μοντέλο $Y = a + b_1X_1 + b_2X_2$ και να ερμηνευτούν τα αποτελέσματά του.

Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

Δοκός	Αντοχή Y	Περιεκτικότητα σε Νερό X_1	Ειδικό Βάρος X_2
1	11.14	11.1	0.499
2	12.74	8.9	0.558
3	13.13	8.8	0.604
4	11.51	8.9	0.441
5	12.38	8.8	0.550
6	12.60	9.9	0.528
7	11.13	10.7	0.418
8	11.70	10.5	0.480
9	11.02	10.5	0.406
10	11.41	10.7	0.467

Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

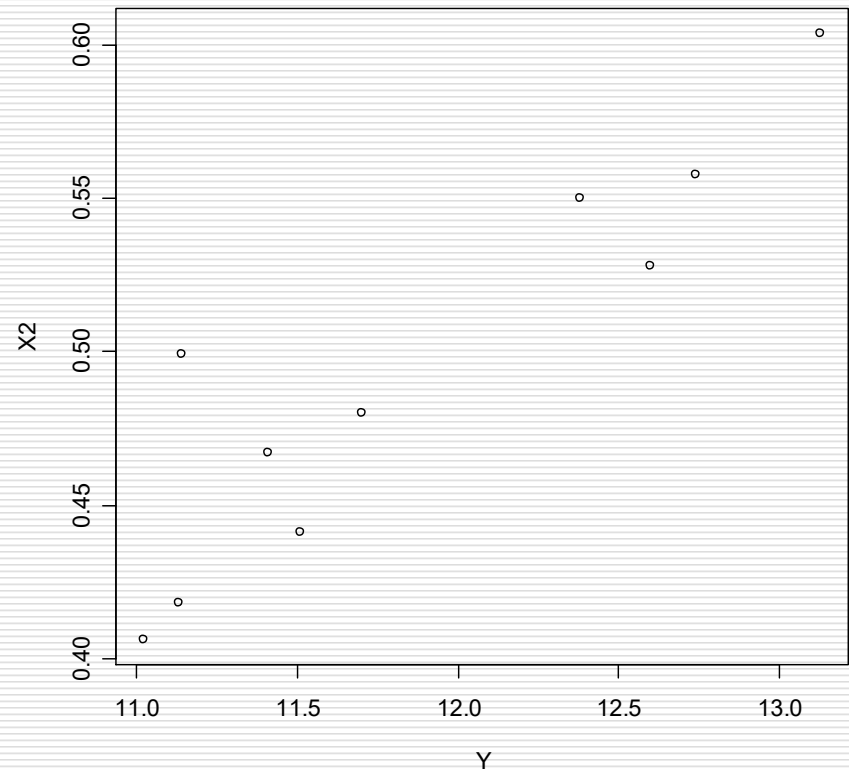
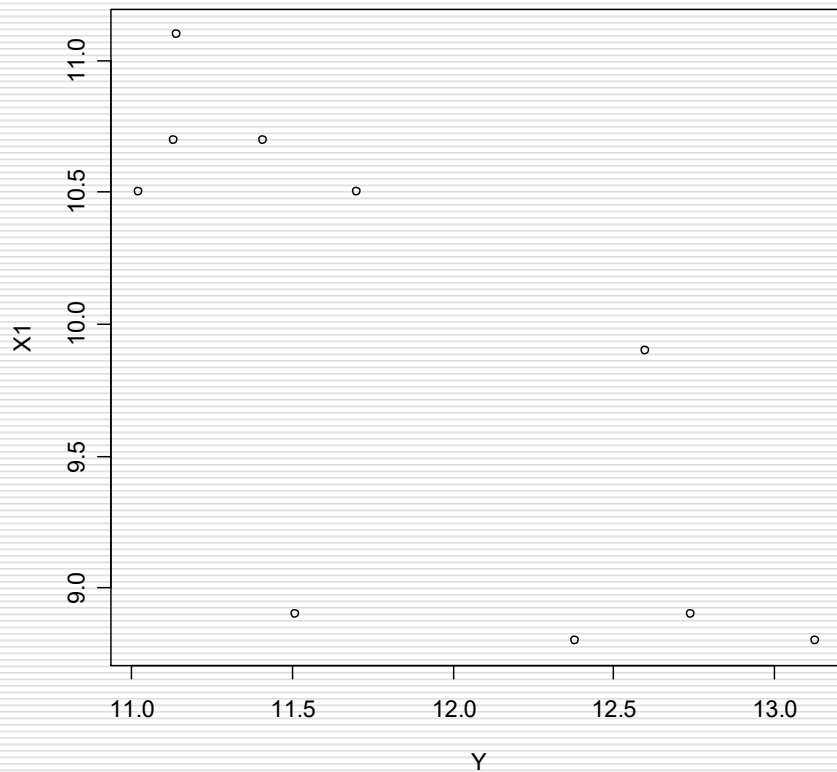
```
> Y<-c(11.14, 12.74, 13.13, 11.51, 12.38,  
12.60, 11.13, 11.70, 11.02, 11.41)  
> X1<-c(11.1, 8.9, 8.8, 8.9, 8.8, 9.9, 10.7,  
10.5, 10.5, 10.7)  
> X2<-c(0.499, 0.558, 0.604, 0.441, 0.55,  
0.528, 0.418, 0.48, 0.406, 0.467)  
> cor(X1,X2)  
[1] -0.6077351
```

παρατηρούμε ότι X_1 και X_2
συσχετίζονται

Με τα επόμενα γραφήματα δεν μπορούμε να ελέγξουμε υπόθεση γραμμικότητας

Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

```
> plot(Y,X1)  
> plot(Y,X2)
```



Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

```
> results<-lm(Y~X1+X2)
```

→ Προσαρμόζουμε το μοντέλο $Y=a+b_1X_1+b_2X_2$

```
> results
```

Call:

```
lm(formula = Y ~ X1 + X2)
```

Coefficients:

(Intercept)	\hat{a}	X1	\hat{b}_1	X2	\hat{b}_2
10.3015		-0.2663		8.4947	

```
> residuals(results)
```

→ Υπόλοιπα

→ Ισοδύναμο με την εντολή `results$res`

```
-0.44421731  0.06868776  0.04129893 -0.16743109  
-0.24998669  0.44985044  0.12732572  0.11738938  
0.06599797 -0.00891511
```

Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

```
> residuals(results, "partial")
      X1      X2
1 -0.76912937 -0.41108794
2  0.32968269  0.60300506
3  0.32892600  0.96637293
4  0.09356385 -0.62699494
5  0.03764038  0.21637293
6  0.44452402  0.72932643
7 -0.09105780 -0.52761648
8 -0.04772987 -0.01088075
9 -0.09912127 -0.69088075
10 -0.22729863 -0.24761648
attr(,"constant")
[1] 11.876
```

Μερικά υπόλοιπα. Για να τα υπολογίσει η R μετασχηματίζει πρώτα τις X_1 και X_2 αφαιρώντας τον μέσο τους και διαιρώντας με την τυπική τους απόκλιση

```
> newx1<-(X1-mean(X1))/sd(X1)
> newx2<-(X2-mean(X2))/sd(X2)
> lm(Y~newx1+newx2)
Call:
lm(formula = Y ~ newx1 + newx2)
```

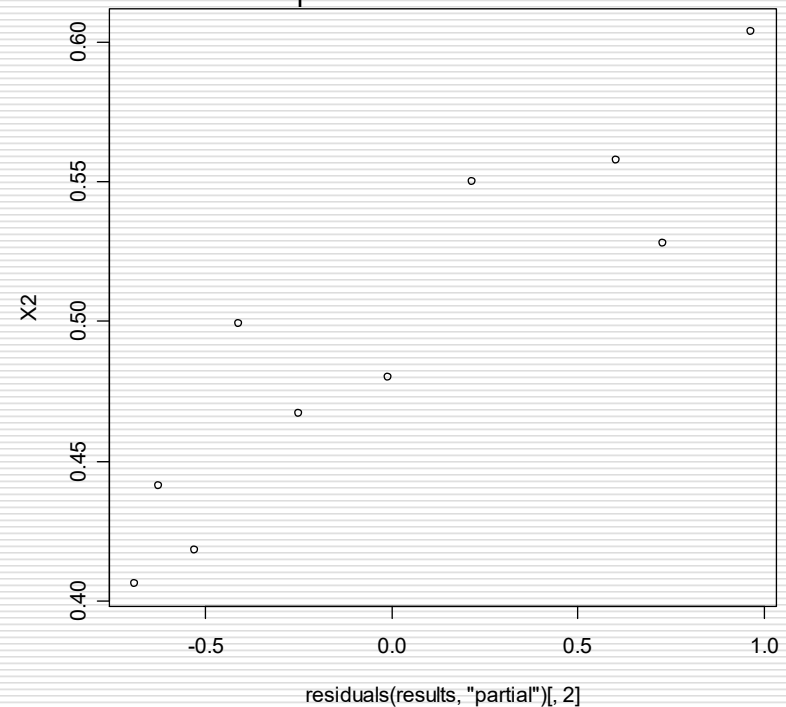
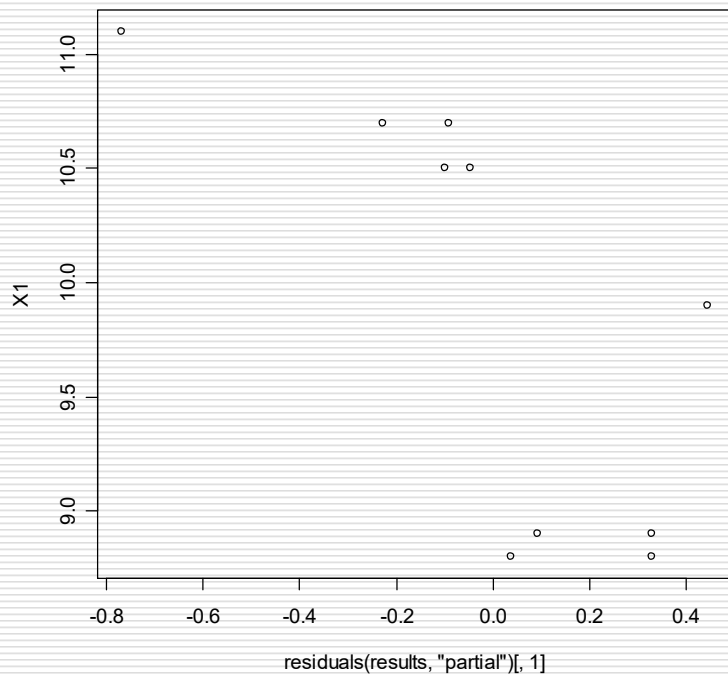
```
Coefficients:
(Intercept)      newx1      newx2
  11.8760      -0.2488      0.5502
```

```
> -0.44421731-0.2488*newx1[1] →  $\hat{\varepsilon}_i + \hat{b}_1 x_{1i}$ 
[1] -0.7691031
```


Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

```
> plot(residuals(results, "partial")[,1],X1)  
> plot(residuals(results, "partial")[,2],X2)
```

υπόθεση γραμμικότητας λογική



Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

```
> summary(results)
```

Call:

```
lm(formula = Y ~ X1 + X2)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.44422	-0.12780	0.05365	0.10521	0.44985

Περιγραφικοί δείκτες υπολοίπων

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	10.3015	1.8965	5.432	0.000975 ***
X1	-0.2663	0.1237	-2.152	0.068394 .
X2	8.4947	1.7850	4.759	0.002062 **

P-τιμή για τον έλεγχο του a

P-τιμή για τον έλεγχο του b_1

P-τιμή για τον έλεγχο του b_2

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.2754 on 7 degrees of freedom

Multiple R-squared: 0.9, Adjusted R-squared: 0.8714

F-statistic: 31.5 on 2 and 7 DF, p-value: 0.0003163

$S_{y|x}$

διορθωμένος συντελεστής προσδιορισμού

P-τιμή για τον F έλεγχο

Παράδειγμα γενικού γραμμικού μοντέλου στην R

```
> confint(results)
```

	2.5 %	97.5 %
(Intercept)	5.8170162	14.7860314
X1	-0.5589333	0.0262906
X2	4.2738001	12.7156214

Συμμετρικά 95% Δ.Ε. για τις παραμέτρους

Συμμετρικό 95% Δ.Ε. για το a

Συμμετρικό 95% Δ.Ε. για το b_1 (περιέχει το 0)

Συμμετρικό 95% Δ.Ε. για το b_2 (δεν περιέχει το 0)

```
> predict(results)
```

1	2	3	4	5	6	7	8
11.58422	12.67131	13.08870	11.67743	12.62999	12.15015	11.00267	11.58261
9	10						
10.95400	11.41892						

Προβλεπόμενες τιμές για κάθε x_i

```
> residuals(results)
```

1	2	3	4	5	6
-0.44421731	0.06868776	0.04129893	-0.16743109	-0.24998669	0.44985044
7	8	9	10		
0.12732572	0.11738938	0.06599797	-0.00891511		

Υπόλοιπα

Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

```
> predict(results, int="c")
      fit   lwr   upr
1 11.58422 11.16318 12.00526
2 12.67131 12.35100 12.99163
3 13.08870 12.66797 13.50943
4 11.67743 11.17108 12.18378
5 12.62999 12.30291 12.95706
6 12.15015 11.89970 12.40060
7 11.00267 10.66952 11.33582
8 11.58261 11.32699 11.83823
9 10.95400 10.58816 11.31985
10 11.41892 11.13701 11.70082
```

Προβλεπόμενες τιμές των y και συμμετρικά 95% Δ.Ε. για κάθε x_i (διαστήματα μέσης πρόβλεψης)

Προβλεπόμενες τιμές των y και συμμετρικά 95% Δ.Ε. όταν $x_1=10$ και $x_2=0.42$ και όταν $x_1=11$ και $x_2=0.48$ (διαστήματα μέσης πρόβλεψης)

```
> predict(results,list(X1=c(10,11), X2=c(0.42, 0.48)), int="c")
```

```
      fit   lwr   upr
1 11.20609 10.84469 11.56749
2 11.44945 11.09098 11.80792
```

```
> predict(results,list(X1=c(10,11), X2=c(0.42, 0.48)), int="p")
```

```
      fit   lwr   upr
1 11.20609 10.46123 11.95095
2 11.44945 10.70601 12.19289
```

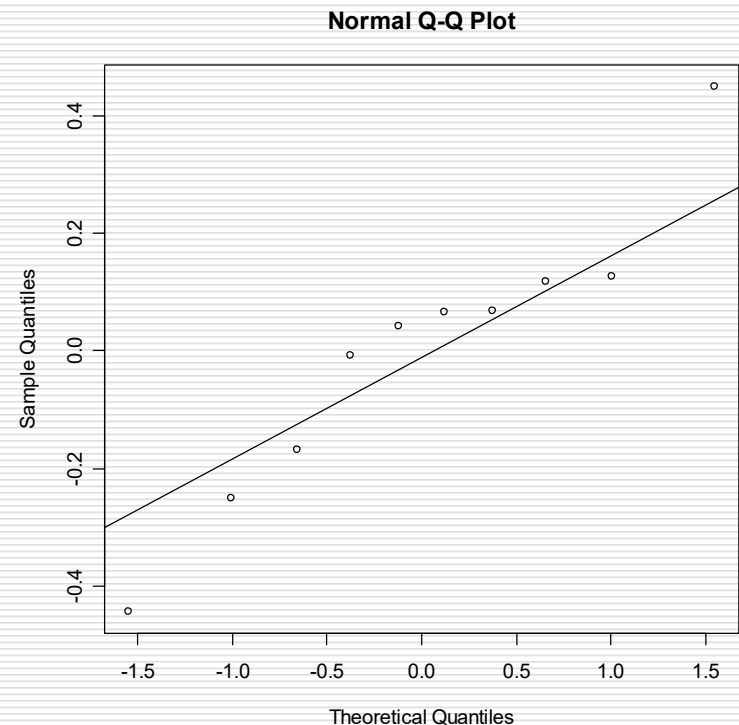
Προβλεπόμενες τιμές των y και συμμετρικά 95% Δ.Ε. όταν $x_1=10$ και $x_2=0.42$ και όταν $x_1=11$ και $x_2=0.48$ (διαστήματα ατομικής πρόβλεψης)

Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

□ Προϋποθέσεις

```
> qqnorm(residuals(results))  
> qqline(residuals(results))
```

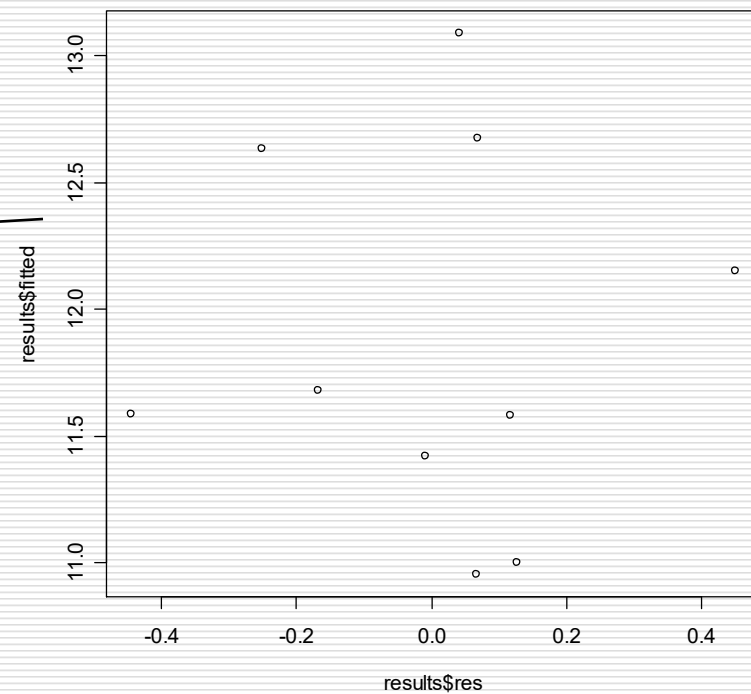
κανονικότητα υπολοίπων



Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

```
> plot(results$res, results$fitted)
```

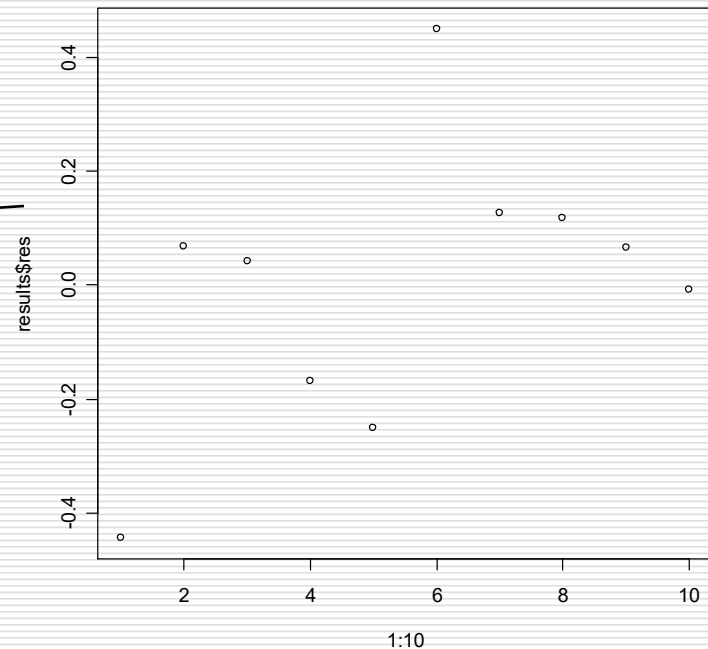
ομοσκεδαστικότητα



Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

```
> plot(1:10,results$res)
```

ανεξαρτησία υπολοίπων



Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

- Παρατηρούμε λοιπόν ότι

$$\hat{y} = 10.30 - 0.27x_1 + 8.49x_2.$$

- Ακόμα παρατηρούμε ότι ο συντελεστής προσδιορισμού $R^2 = 0.9$, έχουμε δηλαδή σχεδόν τέλεια προσαρμογή του μοντέλου, ενώ ο διορθωμένος είναι 0.87 επίσης πολύ υψηλός, ενώ

$$s_{y|x} = 0.27 \text{ (αρκετά μικρό).}$$

Παράδειγμα πολλαπλού γραμμικού μοντέλου στην R

- Στον έλεγχο για το b_1 βλέπουμε ότι η P-τιμή δεν είναι πολύ μικρή οπότε σε ε.σ. 5% δεν μπορούμε να απορρίψουμε την μηδενική υπόθεση ότι αύξηση κατά μια μονάδα της X_1 δεν σημαίνει μεταβολή της αναμενόμενης τιμής της Y (υπό την προϋπόθεση ότι η X_2 παραμένει σταθερή).
- Στον έλεγχο για το b_2 βλέπουμε ότι η P-τιμή είναι πολύ μικρή οπότε πράγματι αύξηση κατά μια μονάδα της X_2 σημαίνει και μεταβολή της αναμενόμενης τιμής της Y (υπό την προϋπόθεση ότι η X_1 παραμένει σταθερή). Πιο συγκεκριμένα αύξηση του ειδικού βάρους των δοκών κατά 1 μονάδα σημαίνει αύξηση της αναμενόμενης αντοχής τους κατά 8.49, υπό την προϋπόθεση βέβαια ότι η περιεκτικότητα σε νερό παραμένει σταθερή.
- Τέλος η σταθερά φαίνεται από την p-τιμή του αντίστοιχου ελέγχου να είναι στατιστικά διάφορη του μηδενός. Πιο συγκεκριμένα η αναμενόμενη αντοχή των δοκών με μηδενική περιεκτικότητα σε νερό και μηδενικό ειδικό βάρος είναι 10.30. Φυσικά η συγκεκριμένη δήλωση δεν έχει ερμηνεία!
- Τέλος από το F-test καταλήγουμε να απορρίψουμε την υπόθεση ότι $b_1=b_2=0$, δηλαδή μπορούμε να εξηγήσουμε την μεταβλητότητα της αντοχής των δοκών από το πολλαπλό γραμμικό μοντέλο.

Διάφορα θέματα στο πολλαπλό γραμμικό μοντέλο

- **Επιλογή Μεταβλητών:** Ένα πρόβλημα που αρκετές φορές αντιμετωπίζουμε είναι η επιλογή των κατάλληλων ανεξάρτητων μεταβλητών που θα χρησιμοποιήσουμε στο μοντέλο μας. Όταν στην διάθεσή μας έχουμε στοιχεία από πολλές επεξηγηματικές μεταβλητές δεν σημαίνει κατά ανάγκη ότι στο τελικό μοντέλο θα πρέπει να χρησιμοποιηθούν όλες αυτές. Υπάρχουν πολλά κριτήρια επιλογής μεταβλητών, η παρουσίαση των οποίων ξεφεύγει από τα όρια αυτών των σημειώσεων.

Διάφορα θέματα στο πολλαπλό γραμμικό μοντέλο

- **Συγχυτικοί Παράγοντες:** Συγχυτικός παράγοντας (*confounder*) ονομάζεται μια μεταβλητή, έστω Z , η οποία διαστρεβλώνει τη σχέση μεταξύ της μεταβλητής απόκρισης Y και μιας επεξηγηματικής μεταβλητής X . Ο λόγος αυτής της διαστρέβλωσης είναι ότι ο συγχυτικός παράγοντας Z σχετίζεται και με την επεξηγηματική μεταβλητή X αλλά και με την μεταβλητή απόκρισης Y . Άρα μπορεί η επεξηγηματική μεταβλητή X να είναι **στατιστικά σημαντική** (δηλαδή αύξησή της κατά μια μονάδα πράγματι να προκαλεί μεταβολή στην αναμενόμενη τιμή της μεταβλητής απόκρισης Y) στο απλό γραμμικό μοντέλο $E(Y|X) = a + bX$, αλλά όταν στο μοντέλο προσθέσουμε και την μεταβλητή Z η επίδρασή της να εξαλείφεται και η Z να είναι αυτή που περιγράφει την μεταβλητότητα της Y .

Διάφορα θέματα στο πολλαπλό γραμμικό μοντέλο

Ως παράδειγμα θεωρήστε τα ακόλουθα δεδομένα τα οποία μας δίνουν το ποσοστό των ενηλίκων που χρησιμοποιούν το διαδίκτυο (Internet), το Ακαθάριστο Εγχώριο Προϊόν (GDP) και το μέσο αριθμό παιδιών ανά ενήλικη γυναίκα (Fertility), για 39 χώρες.

Διάφορα θέματα στο πολλαπλό γραμμικό μοντέλο

INTERNET	GDP	FERTILITY	INTERNET	GDP	FERTILITY	INTERNET	GDP	FERTILITY
0.65	6.09	2.8	37.36	25.35	1.4	0.34	1.89	5.1
10.08	11.32	2.4	13.21	17.44	1.3	2.56	3.84	3.2
37.14	25.37	1.7	0.68	2.84	3.0	2.93	7.10	1.1
38.7	26.73	1.3	1.56	6.00	2.3	1.34	13.33	4.5
31.04	25.52	1.7	23.31	32.41	1.9	6.49	11.29	2.6
4.66	7.36	2.2	27.66	19.79	2.7	18.27	20.15	1.2
46.66	27.13	1.5	38.42	25.13	1.3	51.63	24.18	1.6
20.14	9.19	2.4	27.31	8.75	2.9	30.70	28.10	1.4
2.57	4.02	1.8	3.62	8.43	2.5	6.04	5.89	2.4
42.95	29.00	1.8	49.05	27.19	1.7	32.96	24.16	1.6
0.93	3.52	3.3	46.12	19.16	2.0	50.15	34.32	2.1
43.03	24.43	1.7	0.10	0.85	5.4	1.24	2.07	2.3
26.38	23.99	1.9	46.38	29.62	1.8	0.09	0.79	7.0

Διάφορα θέματα στο πολλαπλό γραμμικό μοντέλο

```
> summary(lm(data$INTERNET~data$FERTILITY))
```

```
Call:  
lm(formula = data$INTERNET ~ data$FERTILITY)
```

```
Residuals:
```

```
   Min     1Q  Median     3Q      Max  
-28.662 -13.890  1.424 10.436 26.727
```

```
Coefficients:
```

```
              Estimate Std. Error t value Pr(>|t|)  
(Intercept)  40.577     5.451  7.444 7.32e-09 ***  
data$FERTILITY -8.169     2.035 -4.013 0.00028 ***
```

```
---  
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 15.62 on 37 degrees of freedom  
Multiple R-squared:  0.3033,    Adjusted R-squared:  0.2845  
F-statistic: 16.11 on 1 and 37 DF,  p-value: 0.0002803
```

```
>summary(lm(data$INTERNET~data$FERTILITY+data$GDP))
```

```
Call:  
lm(formula = data$INTERNET ~ data$FERTILITY + data$GDP)
```

```
Residuals:
```

```
   Min     1Q  Median     3Q      Max  
-23.155 -4.758  0.454  2.761 20.061
```

```
Coefficients:
```

```
              Estimate Std. Error t value Pr(>|t|)  
(Intercept)  -3.1970     5.6922 -0.562  0.578  
data$FERTILITY -0.1180     1.4399 -0.082  0.935  
data$GDP       1.5392     0.1692  9.099 7.28e-11 ***
```

```
---  
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 8.718 on 36 degrees of freedom  
Multiple R-squared:  0.7889,    Adjusted R-squared:  0.7771  
F-statistic: 67.25 on 2 and 36 DF,  p-value: 6.957e-13
```

Διάφορα θέματα στο πολλαπλό γραμμικό μοντέλο

- Παρατηρήστε πως η επίδραση της μεταβλητής Fertility στο γραμμικό μοντέλο εξαλείφεται όταν προσθέσουμε τον συγχυτικό παράγοντα GDP.
- Παρατηρήστε πως σχετίζονται μεταξύ τους οι 3 αυτές μεταβλητές.

```
> cor(data)
      INTERNET      GDP FERTILITY
INTERNET 1.0000000 0.8881536 -0.5507279
GDP      0.8881536 1.0000000 -0.6145094
FERTILITY -0.5507279 -0.6145094 1.0000000
```

Διάφορα θέματα στο πολλαπλό γραμμικό μοντέλο

- **Πολυσυγγραμικότητα.** Δεν θα πρέπει στο μοντέλο παλινδρόμησης να συμπεριλαμβάνουμε επεξηγηματικές μεταβλητές με **υψηλή συσχέτιση** μεταξύ τους, διότι τότε τα αποτελέσματα δεν είναι έγκυρα. Ως παράδειγμα θεωρήστε τα ακόλουθα δεδομένα τα οποία μας δίνουν το ύψος 25 ανδρών σε inches (Height) καθώς και το μήκος σε cm του αριστερού (LeftFoot) καθώς και δεξιού τους ποδιού (RtFoot). Σκοπός είναι να προσαρμόσουμε ένα γραμμικό μοντέλο με μεταβλητή απόκρισης Height και επεξηγηματικές μεταβλητές LeftFoot και RtFoot.

Διάφορα θέματα στο πολλαπλό γραμμικό μοντέλο

Height	LeftFoot	RtFoot	Height	LeftFoot	RtFoot
69.0	27.0	26.5	74.0	29.0	30.0
79.0	29.0	27.5	75.0	28.0	29.0
75.0	31.0	32.0	71.0	27.0	27.5
69.0	25.5	25.5	72.0	26.5	27.5
65.0	23.5	23.0	66.0	25.5	26.0
79.0	28.0	28.0	71.0	29.0	28.0
72.0	28.5	28.5	67.0	27.2	27.0
69.5	27.0	27.0	71.0	29.0	28.5
73.0	30.6	31.4	72.0	28.0	28.0
71.5	27.4	28.5	72.0	28.5	29.0
69.5	27.0	27.0	73.5	29.0	30.0
73.0	28.5	27.5	73.0	23.5	24.0
71.0	29.0	27.0			

Διάφορα θέματα στο πολλαπλό γραμμικό μοντέλο

```
>summary(lm(data$Height~data$LeftFoot+data$RtFoot))
```

```
Call:
lm(formula = data$Height ~ data$LeftFoot + data$RtFoot)
```

```
Residuals:
    Min     1Q  Median     3Q      Max
-4.1330 -1.6460 -0.5014  0.6246  6.9920
```

```
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  43.9334     8.9983   4.882 7e-05 ***
data$LeftFoot  0.6379     0.7730   0.825  0.418
data$RtFoot   0.3647     0.7096   0.514  0.612
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 2.905 on 22 degrees of freedom
Multiple R-squared:  0.3072, Adjusted R-squared:  0.2442
F-statistic: 4.877 on 2 and 22 DF, p-value: 0.01765
```

```
> summary(lm(data$Height~data$LeftFoot))
```

```
Call:
lm(formula = data$Height ~ data$LeftFoot)
```

```
Residuals:
    Min     1Q  Median     3Q      Max
-4.2327 -2.0302 -0.5309  0.4698  6.9684
```

```
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  44.0701     8.8493   4.980 4.9e-05 ***
data$LeftFoot  0.9986     0.3189   3.131 0.00469 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 2.859 on 23 degrees of freedom
Multiple R-squared:  0.2989, Adjusted R-squared:  0.2684
F-statistic: 9.804 on 1 and 23 DF, p-value: 0.004689
```

Διάφορα θέματα στο πολλαπλό γραμμικό μοντέλο

```
> summary(lm(data$Height~data$RtFoot))
```

Call:

```
lm(formula = data$Height ~ data$RtFoot)
```

Residuals:

```
   Min     1Q  Median     3Q    Max
-4.1460 -1.5424 -0.5241  0.2686  7.5095
```

Coefficients:

```
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  46.8408    8.2224   5.697 8.43e-06 ***
data$RtFoot   0.8964    0.2955   3.033 0.00591 **
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Residual standard error: 2.885 on 23 degrees of freedom

Multiple R-squared: 0.2857, Adjusted R-squared: 0.2547

F-statistic: 9.201 on 1 and 23 DF, p-value: 0.00591

Διάφορα θέματα στο πολλαπλό γραμμικό μοντέλο

- Παρατηρήστε πως όταν συμπεριλαμβάνουμε και τις δύο επεξηγηματικές μεταβλητές στο πολλαπλό γραμμικό μοντέλο καμία δεν είναι στατιστικά σημαντική, ενώ αντίθετα κάθε μία εξ αυτών έχει πολύ μεγάλη επίδραση στην μεταβλητή απόκρισης (όπως προφανώς αναμενόταν) όπως παρατηρούμε από τα απλά γραμμικά μοντέλα.
- Το παραπάνω φαινόμενο οφείλεται στην υψηλή συσχέτιση των δύο επεξηγηματικών μεταβλητών.

```
> cor(data$RtFoot,data$LeftFoot)  
[1] 0.9078141
```

Διάφορα θέματα στο πολλαπλό γραμμικό μοντέλο

- **Πρόβλεψη:** Οι προβλέψεις είναι σωστές και αξιόπιστες μόνο για τις τιμές του X που είναι σχετικά κοντά με αυτές που έχουμε παρατηρήσει. Στο παράδειγμα με τις ξύλινες δοκούς δεν θα ήταν λογικό να εκτιμούσαμε την αντοχή των δοκών με περιεκτικότητα σε νερό 30.

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

- Μέχρι τώρα οι επεξηγηματικές μεταβλητές του μοντέλου παλινδρόμησης θεωρήσαμε ότι είναι ποσοτικές τυχαίες μεταβλητές. Ωστόσο, μερικές φορές έχουμε επεξηγηματικές μεταβλητές που αποδίδουν έναν παράγοντα με δύο ή περισσότερες κατηγορίες, όπως π.χ. το φύλο ή την φυσική κατάσταση ενός ατόμου (κακή, μέτρια, καλή και άριστη). Ας υποθέσουμε παραδείγματος χάρη, ότι θέλουμε να δημιουργήσουμε ένα μοντέλο γραμμικής παλινδρόμησης με μεταβλητή απόκρισης Y το καθαρό εισόδημα και επεξηγηματική μεταβλητή X το επίπεδο εκπαίδευσης (δημοτικό, γυμνάσιο, λύκειο, πανεπιστήμιο). Τότε δημιουργούμε τις λεγόμενες **εικονικές μεταβλητές** ή **ψευδομεταβλητές**, οι οποίες είναι $(m-1)$ στον αριθμό, όπου m είναι ο αριθμός των κατηγοριών της κατηγορικής μεταβλητής X , και λαμβάνουν τις τιμές 0 ή 1. Στο παράδειγμα με το εισόδημα και το επίπεδο εκπαίδευσης (4 κατηγορίες), ορίζουμε τις εξής 3 εικονικές μεταβλητές:

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

$$X_1 = \begin{cases} 1, & \text{δημοτικό} \\ 0, & \text{διαφορετικά} \end{cases} \quad X_2 = \begin{cases} 1, & \text{γυμνάσιο} \\ 0, & \text{διαφορετικά} \end{cases} \quad X_3 = \begin{cases} 1, & \text{λύκειο} \\ 0, & \text{διαφορετικά} \end{cases}$$

και εκτιμούμε με την μέθοδο ελαχίστων τετραγώνων τους συντελεστές του γραμμικού μοντέλου

$$E[Y|X_1, X_2, X_3] = a + b_1x_1 + b_2x_2 + b_3x_3.$$

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

- Ισοδύναμα με το παραπάνω γραμμικό μοντέλο είναι σαν να έχουμε τα ακόλουθα 4 γραμμικά μοντέλα:

$$E(Y | X_1, X_2, X_3) = \begin{cases} a, & \text{πανεπιστήμιο} \\ a + b_1, & \text{δημοτικό} \\ a + b_2, & \text{γυμνάσιο} \\ a + b_3, & \text{λύκειο} \end{cases} .$$

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

- Η κατηγορία για την οποία δεν ορίσαμε ψευδομεταβλητή ονομάζεται **κατηγορία αναφοράς** (*reference category*), στο παράδειγμα μας η κατηγορία αυτή είναι η πανεπιστημιακή μόρφωση, και η επιλογή της εξαρτάται συχνά από το υπό μελέτη πρόβλημα. Αν στο παράδειγμά μας θέλουμε να εξετάσουμε αν η πανεπιστημιακή μόρφωση αυξάνει κατά μέσο όρο το εισόδημα, τότε είναι φυσικό να θέσουμε ως κατηγορία αναφοράς την εν λόγω κατηγορία.
- Οι παράμετροι του παραπάνω γραμμικού μοντέλου μπορούν να ερμηνευτούν αν θεωρήσουμε τις πιθανές τιμές των εικονικών μεταβλητών. Έτσι η παράμετρος a στο παραπάνω παράδειγμα εκφράζει το μέσο καθαρό εισόδημα όταν $X_1 = X_2 = X_3 = 0$, άρα όταν η μόρφωση του ατόμου είναι πανεπιστημιακή. Η παράμετρος b_1 αντιπροσωπεύει την διαφορά μεταξύ του μέσου καθαρού εισοδήματος ατόμων με μόρφωση δημοτικού από το μέσο καθαρό εισόδημα ατόμων με πανεπιστημιακή μόρφωση. Ανάλογα ερμηνεύονται και οι υπόλοιποι συντελεστές.
- Τα γραμμικά μοντέλα με μία κατηγορική επεξηγηματική μεταβλητή είναι ισοδύναμα με τα μοντέλα **Ανάλυσης Διασποράς με έναν Παράγοντα** (*One Way ANOVA*). Αν η κατηγορική μεταβλητή είναι δίτιμη τα συμπεράσματα του γραμμικού μοντέλο είναι ισοδύναμα με αυτά του **independent two sample t-test** που είδαμε στο προηγούμενο κεφάλαιο.

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

- **Παράδειγμα:** Μετρήσεις του ετήσιου μισθού Y (σε χιλιάδες €) των υπαλλήλων μιας εταιρείας, συναρτήσει της εμπειρίας τους X_1 (σε έτη υπηρεσίας) και του μορφωτικού τους επιπέδου X_2 (1 = απολυτήριο λυκείου, 2 = απόφοιτος ΑΕΙ και 3 = μεταπτυχιακός τίτλος) δίνονται από τον πίνακα της επόμενης διαφάνειας. Να εκτιμηθούν και να ερμηνευτούν με την μέθοδο ελαχίστων τετραγώνων οι συντελεστές του μοντέλου

$$E(Y | X_1, X_2, X_3) = a + b_1x_1 + b_2x_2 + b_3x_3$$

$$\text{όπου } X_2 = \begin{cases} 1, & \text{απόφοιτος ΑΕΙ} \\ 0, & \text{διαφορετικά} \end{cases} \quad X_3 = \begin{cases} 1, & \text{μεταπτυχιακός τίτλος} \\ 0, & \text{διαφορετικά} \end{cases}$$

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

Y	25	22	20	32	35	28	19	17	20	26	28	34	40	42	28	19
X_1	2	1	2	5	7	3	2	1	2	2	5	5	7	6	3	1
X_2	2	2	1	3	2	3	1	1	2	2	2	3	3	3	2	1

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

- Έστω τα δεδομένα είναι αποθηκευμένα στο αρχείο data.txt (στον φάκελο από όπου τρέχουμε την R) υπό μορφή πίνακα διάστασης 16×3 .

```
> data<-matrix(scan("data.txt"), ncol=16, byrow=T)
Read 48 items
> data
  [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10] [,11] [,12] [,13] [,14]
[1,] 25 22 20 32 35 28 19 17 20 26 28 34 40 42
[2,]  2  1  2  5  7  3  2  1  2  2  5  5  7  6
[3,]  2  2  1  3  2  3  1  1  2  2  2  3  3  3
  [,15] [,16]
[1,] 28 19
[2,]  3  1
[3,]  2  1
> y<-data[1,]
> x1<-data[2,]
> x2<-data[3,]
```

Διαβάζουμε τα
δεδομένα από
το αρχείο

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

```
> x2<-as.factor(x2)
> x2
[1] 2 2 1 3 2 3 1 1 2 2 2 3 3 3 2 1
Levels: 1 2 3
> results<-lm(y~x1+x2)
> summary(results)
Call:
lm(formula = y ~ x1 + x2)
Residuals:
    Min     1Q   Median     3Q     Max
-3.5729 -1.1971 -0.1892  1.4296  4.9010

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 15.1894    1.3709  11.080 1.17e-07 ***
x1           2.3737    0.4058   5.850 7.84e-05 ***
x22          3.6360    1.6780   2.167  0.0511 .
x23          7.6672    2.2294   3.439  0.0049 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1


Residual standard error: 2.457 on 12 degrees of freedom
Multiple R-squared:  0.9187,    Adjusted R-squared:  0.8983
F-statistic: 45.18 on 3 and 12 DF,  p-value: 8.193e-07
```

μετατρέπουμε την μεταβλητή σε κατηγορική

- Η R παίρνει αυτόματα την 1^η κατηγορία ως κατηγορία αναφοράς

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

```
> model.matrix(results)
(Intercept) x1 x22 x23
1           1  2  1  0
2           1  1  1  0
3           1  2  0  0
4           1  5  0  1
5           1  7  1  0
6           1  3  0  1
7           1  2  0  0
8           1  1  0  0
9           1  2  1  0
10          1  2  1  0
11          1  5  1  0
12          1  5  0  1
13          1  7  0  1
14          1  6  0  1
15          1  3  1  0
16          1  1  0  0
```



βλέπουμε τις εικονικές μεταβλητές που δημιουργήθηκαν

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

- Παρατηρούμε λοιπόν ότι: (α) ο αναμενόμενος μισθός ενός υπαλλήλου της εν λόγω εταιρείας με απολυτήριο λυκείου χωρίς προϋπηρεσία είναι €15190, (β) η αύξηση ενός έτους στα χρόνια εμπειρίας μεταφράζεται σε αύξηση του αναμενόμενου μισθού κατά €2370 για υπαλλήλους **ανεξαρτήτως πανεπιστημιακής μόρφωσης**, (γ) οι απόφοιτοι ΑΕΙ έχουν αναμενόμενο μισθό κατά €3640 υψηλότερο σε σχέση με τους συναδέλφους τους με τα ίδια χρόνια εμπειρίας που έχουν απολυτήριο λυκείου και (δ) οι κάτοχοι μεταπτυχιακών τίτλων έχουν αναμενόμενο μισθό κατά €7670 υψηλότερο σε σχέση με τους συναδέλφους τους με τα ίδια χρόνια εμπειρίας που έχουν απολυτήριο λυκείου.

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

□ Παρατηρούμε λοιπόν ότι

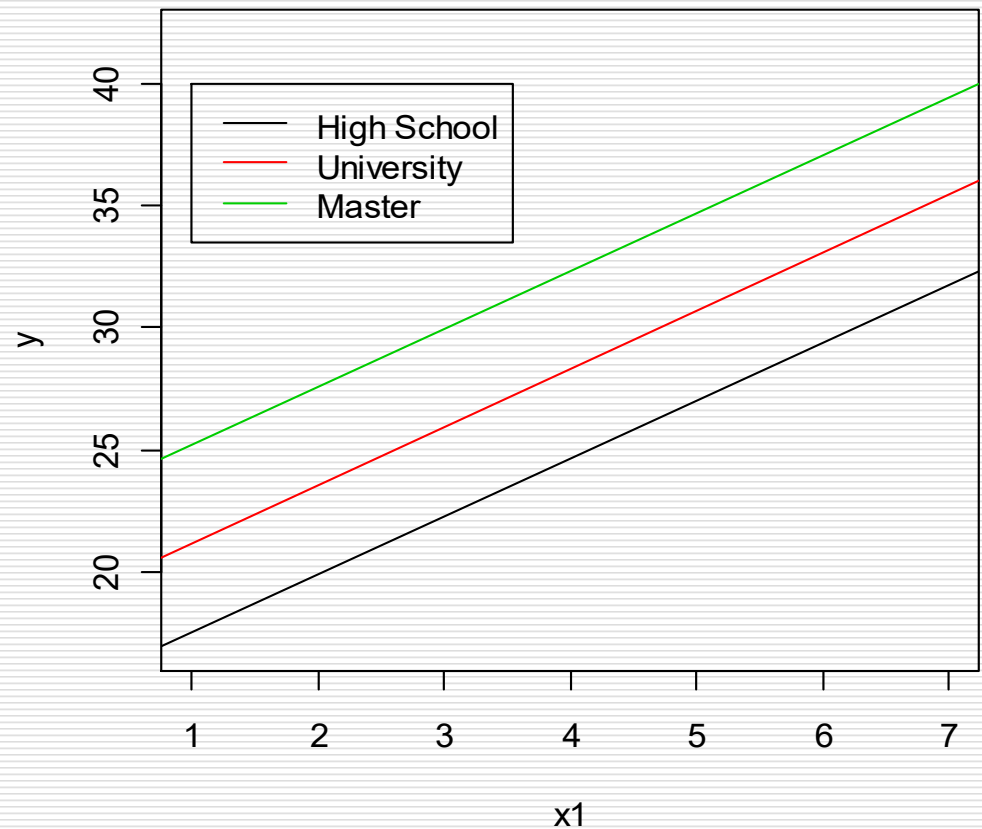
$$E(Y | X_1, X_{22}, X_{23}) = 15.19 + 2.37x_1 + 3.64x_{22} + 7.67x_{23}$$

□ Ισοδύναμα

$$E(Y | X_1, X_{22}, X_{23}) = \begin{cases} 15.19 + 2.37X_1, & \text{απολυτήριο λυκείου} \\ (15.19 + 3.64) + 2.37X_1, & \text{απόφοιτος ΑΕΙ} \\ (15.19 + 7.67) + 2.37X_1, & \text{μεταπτυχιακός τίτλος} \end{cases}$$

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

```
> plot(x1,y,type='n')  
> abline(15.19,2.37, col=1)  
> abline(18.83,2.37, col=2)  
> abline(22.86,2.37, col=3)  
> legend(1,40, lty=1,  
        col=1:3,legend=c("High School",  
                          "University", "Master"))
```



Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

- Αν θέλαμε να χρησιμοποιήσετε άλλη κατηγορία αναφοράς, αλλάζουμε την σειρά των κατηγοριών της μεταβλητής, έτσι ώστε η πρώτη να είναι αυτή που θέλουμε να είναι κατηγορία αναφοράς. Π.χ. αν θέλουμε η δεύτερη κατηγορία (απόφοιτος ΑΕΙ) να είναι η κατηγορία αναφοράς κάνουμε τα εξής:

```
> x2 <- factor(x2, levels = c(2, 1, 3))
> x2
[1] 2 2 1 3 2 3 1 1 2 2 2 3 3 3 2 1
Levels: 2 1 3
> lm(y~x1+x2)
Call:
lm(formula = y ~ x1 + x2)

Coefficients:
(Intercept)      x1      x21      x23
    18.825     2.374    -3.636     4.031
```

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

```
> model.matrix(lm(y~x1+x2))
(Intercept) x1 x21 x23
1           1 2  0  0
2           1 1  0  0
3           1 2  1  0
4           1 5  0  1
5           1 7  0  0
6           1 3  0  1
7           1 2  1  0
8           1 1  1  0
9           1 2  0  0
10          1 2  0  0
11          1 5  0  0
12          1 5  0  1
13          1 7  0  1
14          1 6  0  1
15          1 3  0  0
16          1 1  1  0
```

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

- **Προσαρμογή Μοντέλου Χωρίς τη Σταθερά:** Εναλλακτικά θα μπορούσαμε να είχαμε προσαρμόσει το μοντέλο χωρίς τη σταθερά με όλες τις εικονικές μεταβλητές. Καταλήγουμε τότε στο ίδιο γραμμικό μοντέλο με πριν.

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

```
> summary(lm(y~x1+x2-1))
```

```
Call:
```

```
lm(formula = y ~ x1 + x2 - 1)
```

```
Residuals:
```

```
   Min       1Q   Median       3Q      Max
-3.5729 -1.1971 -0.1892  1.4296  4.9010
```

```
Coefficients:
```

```
   Estimate Std. Error t value Pr(> |t|)
x1    2.3737    0.4058   5.850 7.84e-05 ***
x21   15.1894    1.3709  11.080 1.17e-07 ***
x22   18.8254    1.5775  11.933 5.15e-08 ***
x23   22.8566    2.3790   9.608 5.51e-07 ***
```

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 2.457 on 12 degrees of freedom
```

```
Multiple R-squared:  0.9943,    Adjusted R-squared:  0.9924
```

```
F-statistic: 523.7 on 4 and 12 DF,  p-value: 2.378e-13
```

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

□ Παρατηρούμε λοιπόν ότι

$$E(Y | X_1, X_{22}, X_{23}) = 2.37x_1 + 15.19x_{12} + 18.83x_{22} + 22.86x_{23}$$

□ Ισοδύναμα

$$E(Y | X_1, X_{22}, X_{23}) = \begin{cases} 15.19 + 2.37X_1, & \text{απολυτήριο λυκείου} \\ 18.83 + 2.37X_1, & \text{απόφοιτος ΑΕΙ} \\ 22.86 + 2.37X_1, & \text{μεταπτυχιακός τίτλος} \end{cases}$$

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

- **Μοντέλο με Αλληλεπίδραση των Δύο Επεξηγηματικών Μεταβλητών:** Για το προηγούμενο μοντέλο που προσαρμόσαμε έχουμε υποθέσει ότι τα έτη υπηρεσίας έχουν την ίδια επίδραση στο μισθό για άτομα με διαφορετικό μορφωτικό επίπεδο. Αρκετές φορές αυτή η προϋπόθεση δεν είναι ρεαλιστική και πρέπει να προσθέσουμε στο μοντέλο την **αλληλεπίδραση** (*interaction*) μεταξύ των δύο επεξηγηματικών μεταβλητών.

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

$$E(Y | X_1, X_{22}, X_{23}) = a + b_1x_1 + b_2x_{22} + b_3x_{23} + b_4x_1x_{22} + b_5x_1x_{23}$$

$$E(Y | X_1, X_{22}, X_{23}) = \begin{cases} a + b_1X_1, & \text{απολυτήριο λυκείου} \\ (a + b_2) + (b_1 + b_4)X_1, & \text{απόφοιτος ΑΕΙ} \\ (a + b_3) + (b_1 + b_5)X_1, & \text{μεταπτυχιακός τίτλος} \end{cases}$$

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

```
> x2<- factor(x2, levels=c(1, 2, 3))  
> results2<-lm(y~x1+x2+x1*x2)  
> summary(results2)
```

Call:

```
lm(formula = y ~ x1 + x2 + x1 * x2)
```

Residuals:

```
   Min     1Q  Median     3Q      Max  
-3.9574 -1.1250  0.2899  1.0106  4.0000
```

Coefficients:

```
            Estimate Std. Error t value Pr(>|t|)  
(Intercept) 16.5000     3.7902   4.353 0.00144 **  
x1           1.5000     2.3971   0.626 0.54550  
x2           3.3830     4.1593   0.813 0.43496  
x23          0.5000     5.7595   0.087 0.93253  
x1:x22       0.5372     2.4414   0.220 0.83026  
x1:x23       2.0000     2.5297   0.791 0.44751
```

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.397 on 10 degrees of freedom

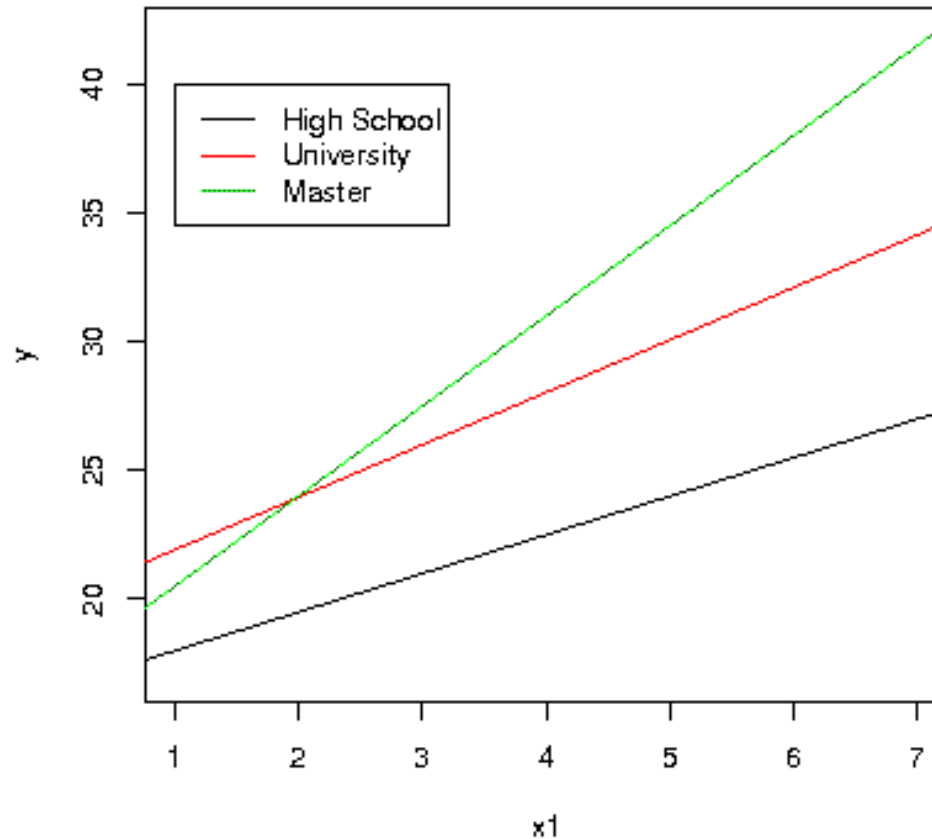
Multiple R-squared: 0.9355, Adjusted R-squared: 0.9032

F-statistic: 28.99 on 5 and 10 DF, p-value: 1.208e-05

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

- Από τα παραπάνω έχουμε ότι τώρα πλέον: (α) ο αναμενόμενος μισθός ενός υπαλλήλου της εν λόγω εταιρείας με απολυτήριο Λυκείου χωρίς προϋπηρεσία είναι 16500 €, (β) η τιμή 3383 € εκφράζει τη διαφορά μεταξύ του μέσου μισθού υπαλλήλων με απολυτήριο Λυκείου χωρίς προϋπηρεσία από τον μέσο μισθό υπαλλήλων που είναι απόφοιτοι ΑΕΙ χωρίς προϋπηρεσία, (γ) η τιμή 500 € εκφράζει τη διαφορά μεταξύ του μέσου μισθού υπαλλήλων με απολυτήριο Λυκείου χωρίς προϋπηρεσία από τον μέσο μισθό υπαλλήλων με μεταπτυχιακό τίτλο χωρίς προϋπηρεσία, (δ) η τιμή 1500 € εκφράζει την αναμενόμενη μεταβολή του μισθού ενός υπαλλήλου της εν λόγω εταιρείας με απολυτήριο Λυκείου όταν αυξηθούν κατά ένα έτος τα χρόνια εμπειρίας, (ε) η τιμή 537 € εκφράζει τη διαφορά μεταξύ των αναμενόμενων μεταβολών του μισθού υπαλλήλων με απολυτήριο Λυκείου και του μισθού υπαλλήλων που είναι απόφοιτοι ΑΕΙ, όταν αυξηθούν κατά ένα έτος τα χρόνια εμπειρίας και (στ) η τιμή 2000 € εκφράζει τη διαφορά μεταξύ των αναμενόμενων μεταβολών του μισθού υπαλλήλων με απολυτήριο Λυκείου και του μισθού υπαλλήλων που είναι κάτοχοι μεταπτυχιακού διπλώματος, όταν αυξηθούν κατά ένα έτος τα χρόνια εμπειρίας.

Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο



Εικονικές Μεταβλητές στο πολλαπλό γραμμικό μοντέλο

- **Σύγκριση των Μοντέλων με και Χωρίς Αλληλεπίδραση:** Τα δύο μοντέλα (με και χωρίς αλληλεπιδράσεις) φαίνεται να προσαρμόζονται εξίσου καλά στα δεδομένα και με τη χρήση τους να εξηγείται το ίδιο περίπου ποσοστό μεταβλητότητας της μεταβλητής απόκρισης (συγκρίνετε τους δύο προσαρμοσμένους συντελεστές προσδιορισμού). Επιπλέον παρατηρούμε ότι το **AIC** έχει χαμηλότερη τιμή για το απλούστερο μοντέλο (χωρίς αλληλεπιδράσεις), το οποίο και προτιμάται τελικά.

```
> AIC(results)
[1] 79.56651
> AIC(results2)
[1] 79.86281
```

Σύνδεση με “two independent sample t-test”

- Πίσω στο παράδειγμα ελέγχου μέσης διάρκειας ζωής των λαμπτήρων πυρακτώσεως και λαμπτήρων φθορίου.
- Δύο ανεξάρτητα δείγματα.
- Υποθέτουμε ισότητα διασπορών.
- Αμφίπλευρος έλεγχος.
- Τα δεδομένα μας είναι σε δύο στήλες, στην πρώτη έχουμε την διάρκεια ζωής 20 λαμπτήρων πυρακτώσεως, στη δεύτερη την διάρκεια ζωής 18 λαμπτήρων φθορίου.

Σύνδεση με t-test

```
> t.test(X,Y, var.equal=TRUE)
```

Two Sample t-test

data: X and Y

t = -30.9836, df = 36, **p-value < 2.2e-16**

alternative hypothesis: true difference in means is not equal to 0

95 percent confidence interval:

-4168.201 -3656.048

sample estimates:

mean of x mean of y

1059.036 4971.161

Διαφορά:

4971.161-1059.036 = **3912.124**

Σύνδεση με t-test

```
> duration<-c(X,Y)
> type<-c(rep(0,20), rep(1,18))
> data.frame(duration, type)
```

	duration	type
1	1101.0377	0
2	1786.3295	0
3	1277.6971	0
4	1109.6934	0
5	1197.7701	0
6	1686.8745	0
7	861.6086	0
8	133.7824	0
9	1261.7946	0
10	457.8806	0
11	400.9930	0
12	796.5819	0
13	1101.5362	0
14	664.1593	0
15	1235.7157	0
16	1439.7528	0
17	2140.8110	0
18	688.0261	0

19	700.0734	0
20	1138.6054	0
21	5275.4390	1
22	4894.0540	1
23	4992.2390	1
24	4739.0010	1
25	4946.2540	1
26	4493.2450	1
27	4842.2370	1
28	5256.0850	1
29	4780.9120	1
30	4727.2350	1
31	5357.2850	1
32	5143.2490	1
33	5020.9290	1
34	4790.3500	1
35	5013.4820	1
36	5230.8480	1
37	4929.0220	1
38	5049.0230	1

Σύνδεση με t-test

```
> summary(lm(duration~type))
```

Call:

```
lm(formula = duration ~ type)
```

Residuals:

Min	1Q	Median	3Q	Max
-925.25	-223.48	42.16	196.24	1081.77

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1059.0	86.9	12.19	2.45e-14 ***
type	3912.1	126.3	30.98	< 2e-16 ***

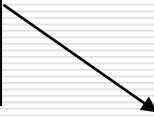
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 388.6 on 36 degrees of freedom

Multiple R-squared: 0.9639, Adjusted R-squared: 0.9629

F-statistic: 960 on 1 and 36 DF, p-value: < 2.2e-16

Πολλαπλασιαστικό Μοντέλο με Εικονικές Μεταβλητές

$$Y \approx \beta_0 X_1^{\beta_1} X_2^{\beta_2} e^{\beta_3 X_3}$$


$$\ln(Y) \approx \ln(\beta_0) + \beta_1 \ln(X_1) + \beta_2 \ln(X_2) + \beta_3 X_3$$

Y = Εβδομαδιαίος Αριθμός Πωλήσεων

X_1 = Τιμή προϊόντος την τρέχουσα εβδομάδα (\$)

X_2 = Κόστος Διαφήμισης για το προϊόν την τρέχουσα εβδομάδα (σε \$ 100s)

X_3 = Αργία ($X_3 = 1$ αν η εν λόγω εβδομάδα είχε αργία)

($X_3 = 0$ αν όχι)

Πολλαπλασιαστικό Μοντέλο με Εικονικές Μεταβλητές

$$\ln(Y) \approx 6.10 - 0.37 \ln(X_1) + 0.42 \ln(X_2) + 0.12 X_3$$

- Για τους συντελεστές -0.37 και 0.42 έχουμε τις γνωστές ερμηνείες ως ελαστικότητες, δηλαδή ποσοστιαία μεταβολή στη διάμεσο του Y όταν το αντίστοιχο X αυξηθεί κατά 1% και τα υπόλοιπα μένουν σταθερά.
- **Προφανώς στην εικονική μεταβλητή δεν βάζουμε λογάριθμο!**
- Όταν η εβδομάδα έχει μέρα αργίας, η διάμεσος των πωλήσεων του προϊόντος θα αυξηθεί κατά 13% ($= (1.13-1)$, $1.13 = \exp(0.12)$), σε σχέση με μια βδομάδα που δεν έχει μέρα αργίας στην οποία η τιμή πώλησης και το κόστος διαφήμισης του προϊόντος παραμένουν σταθερά.
- Αν ο συντελεστής ήταν μικρότερος του μηδενός (και άρα η εκθετική τιμή αυτού ήταν αριθμός μικρότερος του 1, π.χ. 0.90) θα λέγαμε το εξής. Όταν η εβδομάδα έχει αργία, και δεδομένου ότι η τιμή πώλησης και το κόστος διαφήμισης παραμένουν σταθερά η διάμεσος των πωλήσεων θα μειωθεί κατά 10% ($= 1 - 0.90$).

Σύγκριση Μοντέλων

- Έστω δύο μοντέλα M_1 και M_2 .
- Με τα κριτήρια AIC και BIC λαμβάνεται υπόψη η **καλή προσαρμογή** των μοντέλων αλλά και η **πολυπλοκότητά** τους

$$AIC(m) = -2 \log L(m) + 2 \dim(m)$$

$$BIC(m) = -2 \log L(m) + \dim(m) \log(n)$$

Μοντέλο

Αριθμός παραμέτρων

Μεγιστοποιημένη τιμή της λογαριθμικής πιθανοφάνειας

Σύγκριση Μοντέλων

- ❑ Τα υπό σύγκριση μοντέλα μπορεί και να μην είναι εμφωλευμένα καθώς και να έχουν και διαφορετικό αριθμό επεξηγηματικών μεταβλητών.
- ❑ Πρέπει να έχουν **την ίδια μεταβλητή απόκρισης**.
- ❑ Στην R υπάρχουν οι έτοιμες συναρτήσεις AIC() και BIC().
- ❑ Μπορούν να χρησιμοποιηθούν και για μη γραμμικά μοντέλα.
- ❑ Διαλέγουμε το μοντέλο με το μικρότερο AIC/BIC.
- ❑ Το BIC ποινικοποιεί περισσότερο τα πιο πολύπλοκα μοντέλα.

Σύγκριση Μοντέλων

- Έστω ότι θέλουμε να συγκρίνουμε τα μοντέλα

$$M_1 : Y = \alpha + \beta X + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

$$M_2 : \log Y = \alpha + \beta \log X + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

- **Πρόβλημα:** Δεν μπορούμε να χρησιμοποιήσουμε το AIC ή το BIC. Οι πρώτοι όροι (log likelihood) είναι μη συγκρίσιμοι (χρησιμοποιούμε τις τιμές y_i στο πρώτο και τις τιμές $\log y_i$ στο δεύτερο μοντέλο).

- **Λύση:** Στην τιμή του κριτηρίου για το δεύτερο μοντέλο πρόσθεσε τον όρο $\rightarrow 2 \sum_{i=1}^n \log y_i$

(= απόλυτη τιμή Ιακωβιανής $\rightarrow \prod_1^n \left| \frac{\partial \log y_i}{\partial y_i} \right| = \prod_1^n \left| \frac{1}{y_i} \right|$
πολλαπλασιασμένη επί 2 σε λογαριθμική κλίμακα).

Σύγκριση Μοντέλων

□ Άλλα κριτήρια

1. **Adjusted R²**. Θα πρέπει να έχουμε την ίδια Y. Μόνο γραμμικά μοντέλα.
2. **Μέσο τετραγωνικό σφάλμα.**

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{RMSE} = \sqrt{\text{MSE}}$$

Αν συγκρίνουμε δύο μοντέλα θα πρέπει τα μοντέλα να έχουν προσαρμοστεί στην ίδια Y. Σε διαφορετική περίπτωση οι εκτιμώμενες τιμές να είναι στην αρχική κλίμακα. Για τα μοντέλα M_1 και M_2 π.χ., βρίσκουμε αρχικά τις προσαρμοσμένες τιμές (σε λογαριθμική κλίμακα), και εν συνεχεία με το \exp της ανάγουμε στην αρχική κλίμακα.

Σύγκριση Μοντέλων

3. Μέσο απόλυτο ποσοστιαίο σφάλμα

$$\text{MAPE} = \frac{1}{n} \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{|y_i|} \times 100\%$$

- Εκφράζει ποσοστό. Οπότε μπορεί να χρησιμοποιηθεί και όταν η μεταβλητή απόκρισης είναι σε άλλη κλίμακα.
- Π.χ. για τα μοντέλα M_1 και M_2 μπορούμε να συγκρίνουμε (καλή προσαρμογή)

$$\text{MAPE}_1 = \frac{1}{n} \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{|y_i|} \times 100\%, \quad \hat{y}_i \rightarrow M_1$$

$$\text{MAPE}_2 = \frac{1}{n} \sum_{i=1}^n \frac{|\log y_i - \log \hat{y}_i|}{|\log y_i|} \times 100\%, \quad \log \hat{y}_i \rightarrow M_2$$

Σύγκριση Μοντέλων

- ή να συγκρίνουμε (προβλεπτική ικανότητα)

$$\text{MAPE}_1 = \frac{1}{n} \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{|y_i|} \times 100\%, \quad \hat{y}_i \rightarrow M_1$$

$$\text{MAPE}_2 = \frac{1}{n} \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{|y_i|} \times 100\%, \quad \hat{y}_i = \exp(\widehat{\log y_i}), \quad \widehat{\log y_i} \rightarrow M_2$$

Σύγκριση Μοντέλων

- ❑ Τα MSE και MAPE δεν ποινικοποιούν πιο πολύπλοκα μοντέλα.
- ❑ Προτού συγκρίνουμε τα μοντέλα ελέγχουμε τις **προϋποθέσεις** τους. Αν κάποιο δεν τις ικανοποιεί δεν τα συγκρίνουμε.
- ❑ Προσοχή στην υπερπροσαρμογή (**overfitting**)! Καλύτερες τεχνικές είναι αυτές που χρησιμοποιούν **διασταυρωμένη επικύρωση** (**cross validation**) και δεν χρησιμοποιούν τα δεδομένα δύο φορές.

Επίλογος

- Αν η μεταβλητή απόκρισης Y είναι **κατηγορική** μεταβλητή, τότε το γραμμικό μοντέλο δεν είναι πλέον το κατάλληλο, και χρησιμοποιούμε **γενικευμένα γραμμικά μοντέλα**, η παρουσίαση των οποίων ξεφεύγει από τους σκοπούς αυτών των σημειώσεων.
- Π.χ. έστω Y δίτιμη \rightarrow **Λογιστική Παλινδρόμηση**.