



*Σχολή Χημικών Μηχανικών ΕΜΠ  
Ανάλυση Συστημάτων Χημικής Μηχανικής, 2<sup>ο</sup> εξάμηνο*

# Προσαρμογή καμπύλης με τη μέθοδο των ελαχίστων τετραγώνων

*Διδάσκοντες: Χ. Κυρανούδης, Γ. Μαυρωτάς*



# Εισαγωγή

- Με βάση κάποιο δείγμα  $(X, Y)$  ζητούμε να εξάγουμε συμπεράσματα για τη σχέση μεταξύ εξαρτημένης μεταβλητής ( $Y$ ) και της ανεξάρτητης ( $X$ )
- Εξίσωση παλινδρόμησης
- Π.χ. Γραμμική σχέση:  $Y = \alpha X + \beta$ 
  - $Y$ : εξαρτημένη μεταβλητή
  - $X$ : ανεξάρτητη μεταβλητή
  - $\alpha, \beta$ : Παράμετροι της εξίσωσης παλινδρόμησης



- Τα δεδομένα δεν ταιριάζουν ποτέ ακριβώς επάνω στα δεδομένα λόγω παραγόντων λάθους (πειραματικά και άλλου τύπου σφάλματα)
- Μέθοδος ελαχίστων τετραγώνων (least squares)
  - Η μέθοδος ελαχίστων τετραγώνων είναι επιδιώκει την ελαχιστοποίηση του σφάλματος μεταξύ του μοντέλου και των δεδομένων.
  - Ορίζουμε σαν συνάρτηση σφάλματος το άθροισμα των τετραγώνων των σφαλμάτων (SSE).

$$SSE = \sum_{i=1}^n (Y_i - Y_i^{est})^2$$

$Y_i$ : πραγματική τιμή i-παρατήρησης  
 $Y_i^{est}$ : εκτιμούμενη τιμή i-παρατήρησης



# Μέθοδος ελαχίστων τετραγώνων

- Για γραμμικό μοντέλο ισχύει:

$$SSE = \sum_{i=1}^n (Y_i - Y_i^{est})^2 \quad \text{γραμμικό μοντέλο} \quad = \quad \sum_{i=1}^n (Y_i - aX_i - \beta)^2$$

- Για να ελαχιστοποιήσουμε μια παράσταση θέτουμε τις μερικές παραγώγους ίσες με μηδέν
- Οι τιμές των παραμέτρων  $a$ ,  $\beta$  προκύπτουν από τη λύση του συστήματος:

$$\partial SSE / \partial a = 0$$

$$\partial SSE / \partial \beta = 0$$

- Προκύπτουν δύο εξισώσεις (όσες κι οι παράμετροι)



# Γραμμικά μοντέλα

Ειδικά για γραμμικά μοντέλα ισχύει:

$$\partial SSE / \partial \alpha = 0 \Rightarrow \partial \left[ \sum_{i=1}^n (Y_i - \alpha X_i - \beta)^2 \right] / \partial \alpha = 0 \Rightarrow \sum_{i=1}^n (2\alpha X_i^2 - 2X_i Y_i + 2\beta X_i Y_i) = 0$$

$$\partial SSE / \partial \beta = 0 \Rightarrow \partial \left[ \sum_{i=1}^n (Y_i - \alpha X_i - \beta)^2 \right] / \partial \beta = 0 \Rightarrow \sum_{i=1}^n (2\beta - 2Y_i + 2\alpha X_i Y_i) = 0$$

- Υπολογισμός των παραμέτρων  $\alpha$ ,  $\beta$  από το γραμμικό σύστημα που προκύπτει:

$$\sum_{i=1}^n (2X_i^2)\alpha + \sum_{i=1}^n (2X_i Y_i)\beta = \sum_{i=1}^n (2X_i Y_i)$$

$$\sum_{i=1}^n (2X_i Y_i)\alpha + 2n\beta = \sum_{i=1}^n (2Y_i)$$

$$\alpha = \frac{SS_{xy}}{SS_{xx}} = \frac{\sum_{i=1}^n X_i Y_i - n \cdot \bar{X} \cdot \bar{Y}}{\sum_{i=1}^n X_i^2 - n \cdot \bar{X}^2}$$

$$\beta = \bar{Y} - \alpha \cdot \bar{X}$$



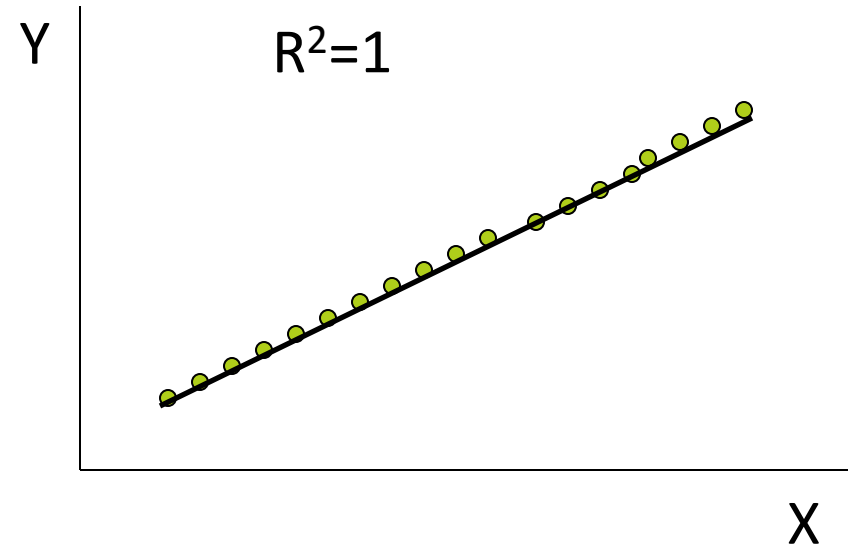
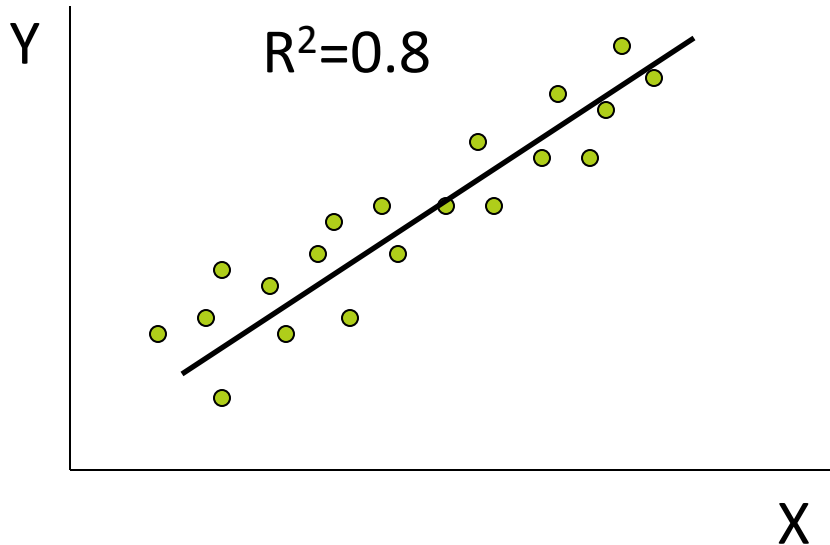
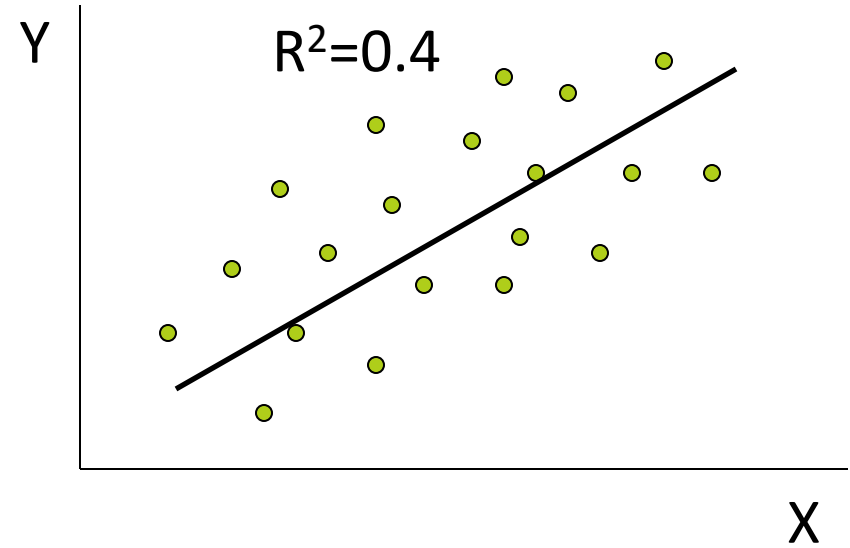
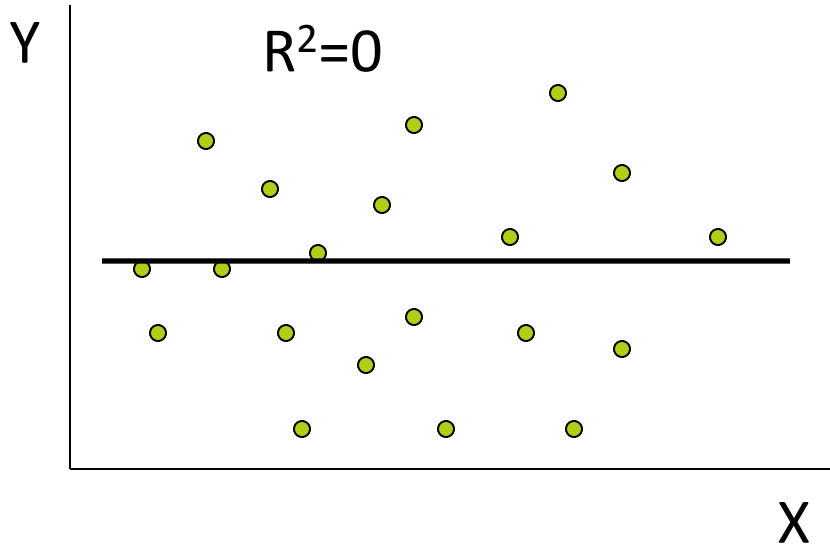
# Συντελεστής προσδιορισμού

- Ενδεικτικό της καλής προσαρμογής του μοντέλου παλινδρόμησης στα δεδομένα είναι ο συντελεστής προσδιορισμού  $R^2$ .

$$R^2 = 1 - \frac{SSE}{SST} = 1 - \frac{\sum_{i=1}^n (Y_i - Y_i^{est})^2}{\sum_{i=1}^n (Y_i - Y^{avg})^2}$$

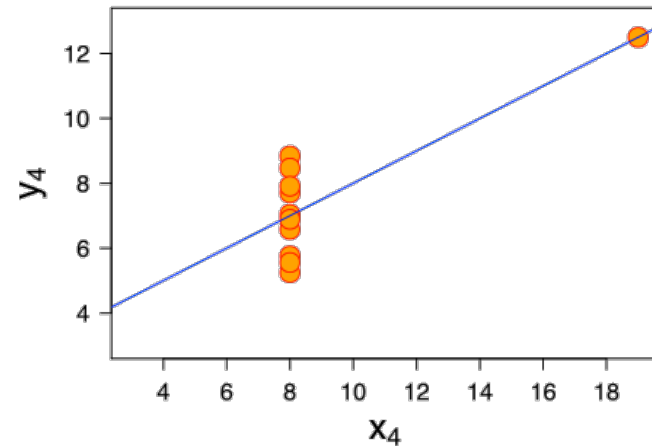
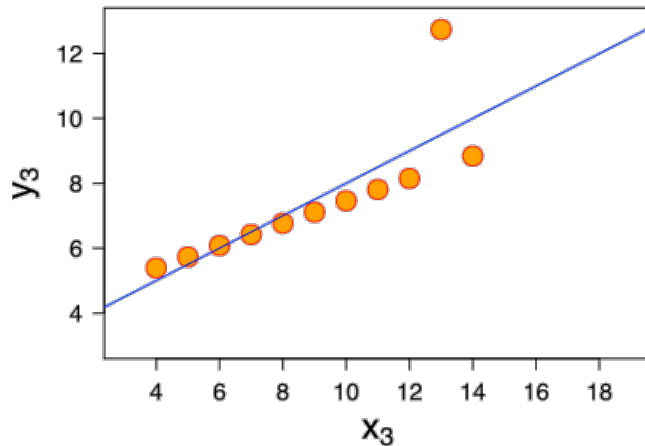
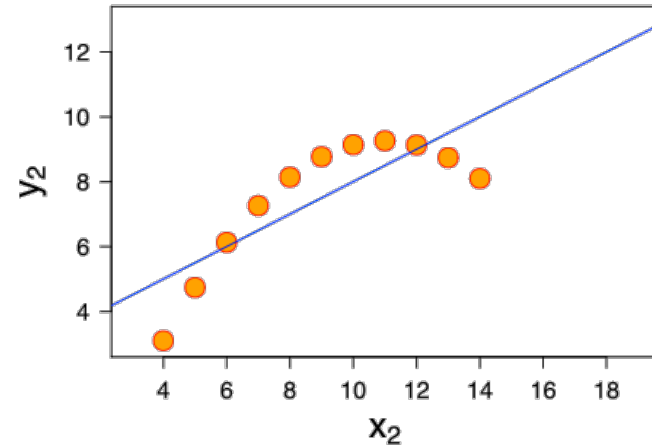
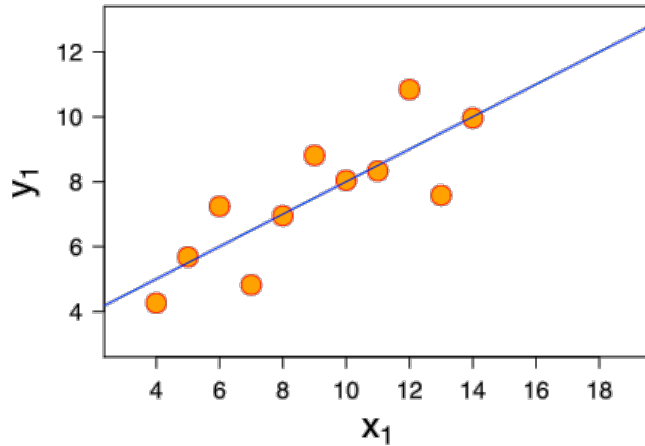


# Συντελεστής προσδιορισμού (2)





# Προσοχή-Κοινή εξίσωση προσέγγισης δεν σημαίνει ίδια δεδομένα

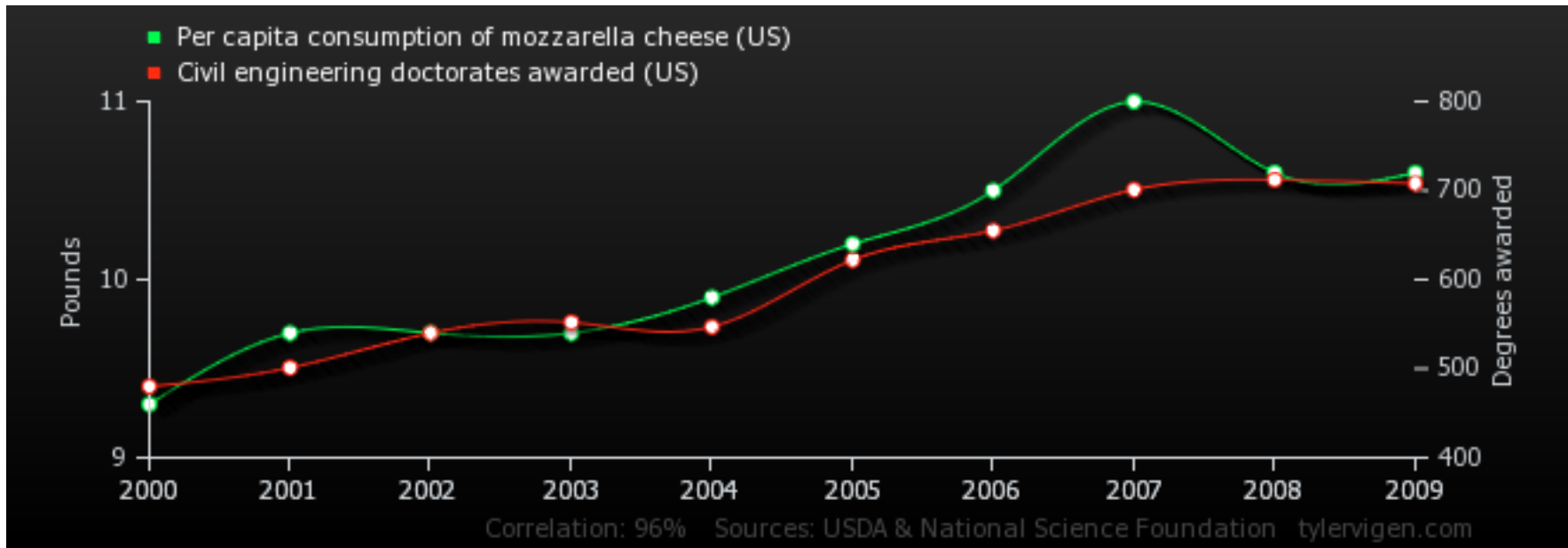


Τα δεδομένα αυτά προτάθηκαν από τον Francis Anscombe για να δείξουν την αξία της βαθύτερης μελέτης των δεδομένων  
[https://en.wikipedia.org/wiki/Anscombe%27s\\_quartet](https://en.wikipedia.org/wiki/Anscombe%27s_quartet)





# Προσοχή - Μαθηματική συσχέτιση δεν σημαίνει πάντα πραγματική σχέση



Κατά κεφαλή κατανάλωση mozzarella (USA) - Pounds (USDA)

Πλήθος νέων διδακτόρων Πολιτικών Μηχανικών που αναγορεύτηκαν (USA) (National Science Foundation)

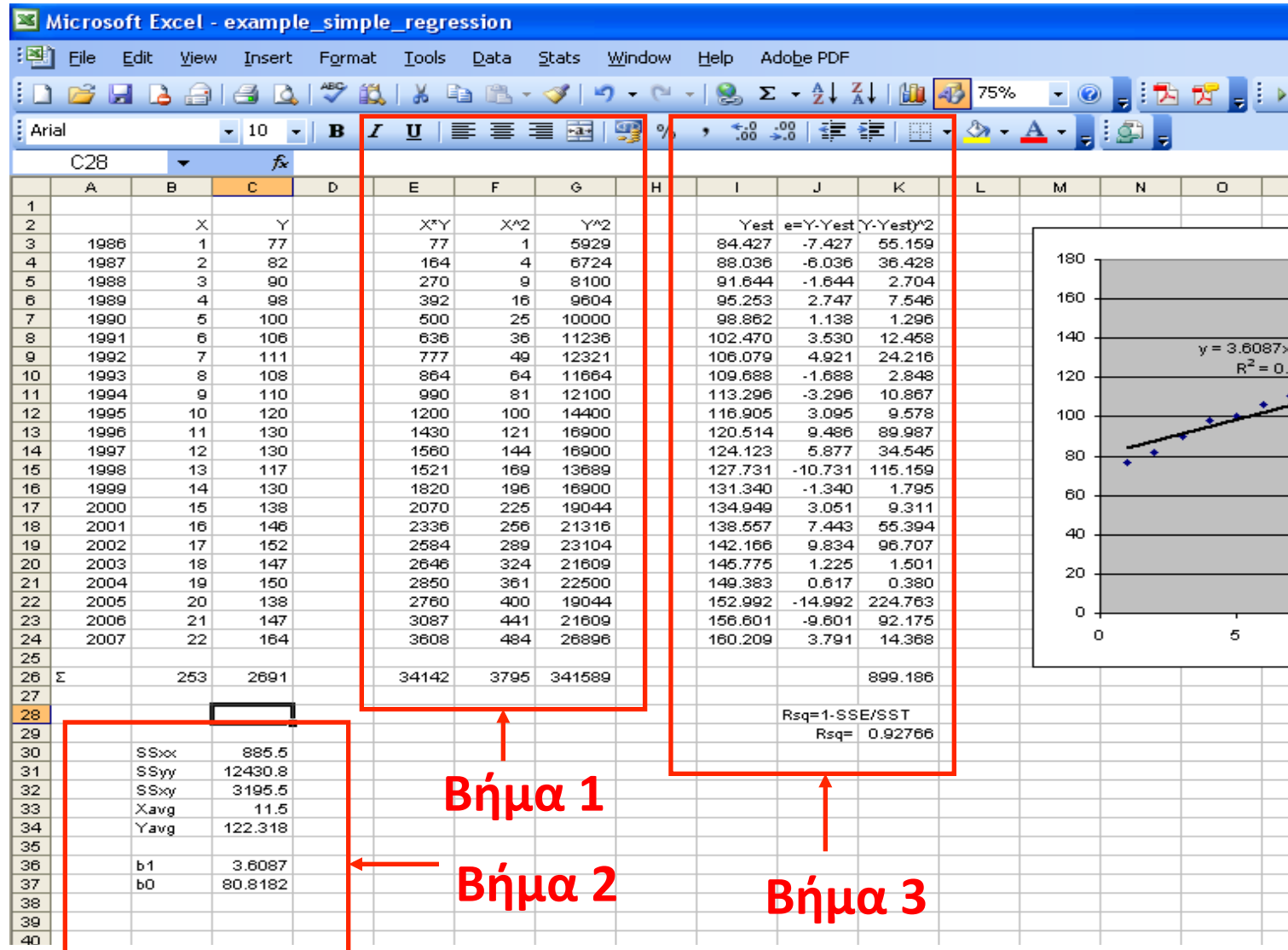
2000	2001	2002	2003	2004	2005	2006	2007	2008	2009
9.3	9.7	9.7	9.7	9.9	10.2	10.5	11	10.6	10.6
480	501	540	552	547	622	655	701	712	708

$R^2: 0.958648$

<http://www.tylervigen.com/>



# Παράδειγμα στο Excel





# Μη γραμμικά μοντέλα

- Π.χ.  $Y = \alpha + \beta X + \gamma X^2$  ή  $Y = \frac{a}{1 + 3e^{-\beta X}}$
- Μετασχηματισμό σε αντίστοιχα γραμμικά
  - Π.χ.  $Y = \alpha \cdot e^{-\beta X}$  (με λογαριθμοποίηση)
- Αν η γραμμικοποίηση δεν είναι δυνατή τότε αναζητούμε τις τιμές των παραμέτρων  $\alpha, \beta$  που ελαχιστοποιούν το SSE με μια μέθοδο βελτιστοποίησης



# Γραμμική παρεμβολή

Για να βρίσκουμε τιμές μεταξύ δύο τιμών πίνακα

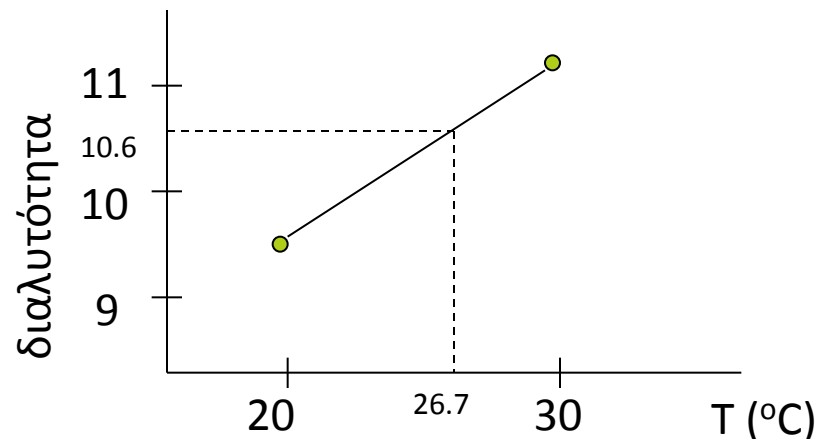
T (°C)	60	50	40	30	20	10
C (mol/l)	16.4	14.45	12.7	11.1	9.6	8.15

Έστω ότι ψάχνω την διαλυτότητα για  $T = 26.7^{\circ}\text{C}$

Η ευθεία που ορίζεται από τα σημεία A( $20^{\circ}\text{C}$ , 9.6) και B( $30^{\circ}\text{C}$ , 11.1) έχει εξίσωση:

$$(x-20)/(30-20) = (y-9.6)/(11.1-9.6)$$

Αντικαθιστώντας όπου x το 26.7  $\Rightarrow (26.7-20)/(30-20) = (y-9.6)/(11.1-9.6) \Rightarrow y = 10.6$





# Επίλυση με Solver Excel

- Δίπλα στη στήλη με τα πραγματικά δεδομένα βάζουμε τα estimated από το μοντέλο (προσοχή: απόλυτη αναφορά με \$ στα κελιά των παραμέτρων του μοντέλου)
- Υπολογίζουμε σε διπλανή στήλη το τετράγωνο των διαφορών για κάθε σημείο
- Κάτω από τη στήλη αυτή βάζουμε την παράσταση με το:  
$$SSE (=sum([actual-estimated]^2)$$
 σε ένα κελί
- Τρέχουμε τον Solver με κελί προς ελαχιστοποίηση αυτό που περιέχει το SSE και μεταβαλλόμενα κελιά αυτά που περιέχουν τις παραμέτρους του μοντέλου
- Υπολογίζονται αυτόματα όλες οι τιμές του μοντέλου