

1

## Πάλιν ερώτηση και συσχέτιση

Τεχνικές που επιτρέπουν την μέτρηση της αλληλεπίδρασης ανάμεσα σε δύο μεταβλητές

Ποσοτικές τεχνικές (συγείου ή διαστημικός)

1) Διερεύνηση σχέσης μεταξύ δύο μεταβλητών

Η σχέση ανάμεσα σε μια ανεξάρτητη μεταβλητή  $X$  και σε μια εξαρτημένη μεταβλητή  $Y$

υπορεί να είναι:

ΜΑΘΗΜΑ  
5

- Γραμμική:  $y = b_0 + b_1 x$
- Παραβολική:  $y = b_0 + b_1 x + b_2 x^2$
- Εκθετική:  $y = b_0 b_1^x$  ή  
 $\ln y = \ln b_0 + x \ln b_1$
- Λογαριθμική:  $y = b_0 + b_1 \ln x$

Γραμμική σχέση με πολλές ανεξάρτητες μεταβλητές

$$y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

2) Μέτρηση του βαθμού συσχέτισης

Συντελεστής συσχέτισης Pearson.

## 2) Ποιοτικές γέθοσες

Ιεραρχικός κλίμακας (γέτρηση έντασης

συσχέτισης: συντελεστής συσχέτισης

Spearman, συντελεστής Gamma,

συντελεστές tau-b, tau-c)

Ονομαστικός κλίμακας

$\chi^2$ , λόγος πιθανοφάνειας, συντελεστής phi

### ΠΑΛΙΝΔΡΟΜΗΣΗ

Απλή (για ανεξάρτητη μεταβλητή)

Πολλαπλή (πολλές ανεξάρτητες μεταβλητές

προσδιορίζει τη σχέση ανάμεσα σε δύο μεταβλητές με σκοπό τη δυνατότητα πρόβλεψης και ελέγχου.

(ηχ πρόβλεψη εξόδων γα βάση το εισόδημα)

### Απλή παλινδρόμηση

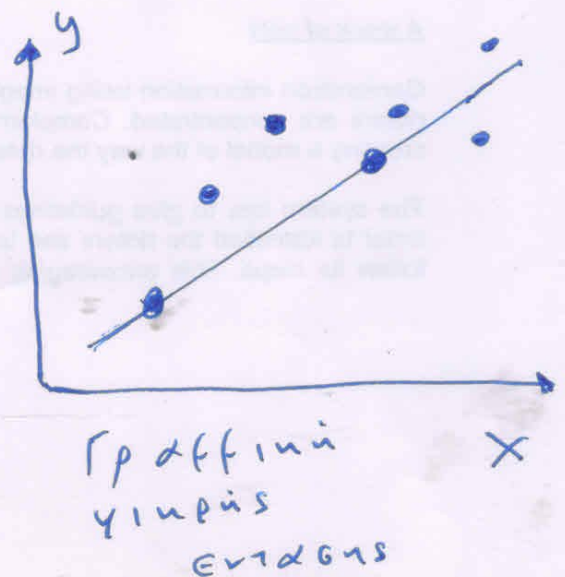
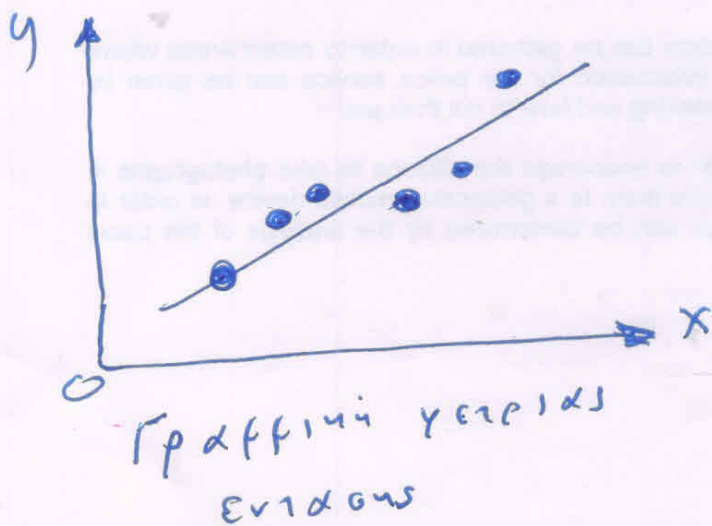
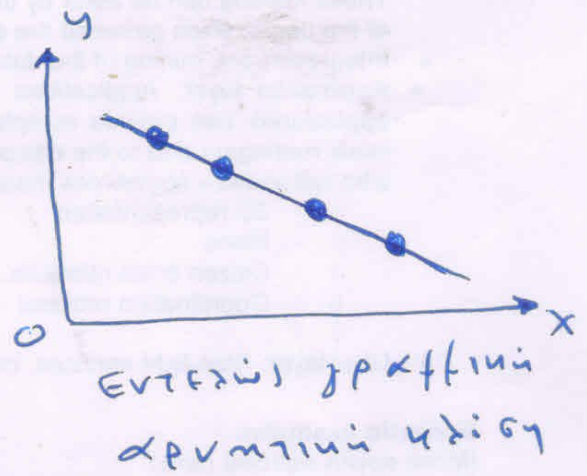
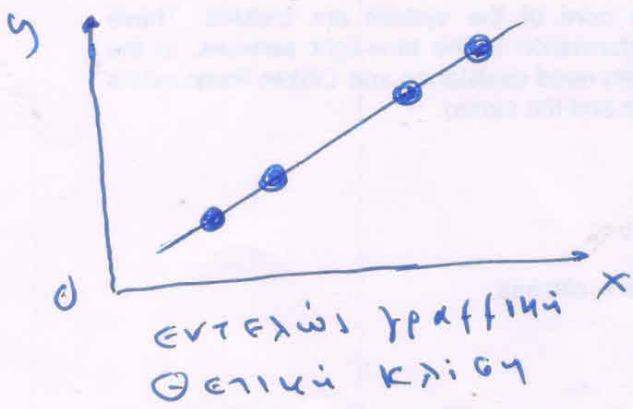
Έστω δύο μεταβλητές  $X$  και  $Y$  με την εξαρτημένη μεταβλητή  $Y$  να εξαρτάται από την ανεξάρτητη μεταβλητή  $X$  μέσω μιας σχέσης 1-προς-1. (Εάν σε ένα  $X$  αντιστοιχού πολλὰ  $Y$  η εξάρτηση λέγεται στοχαστική ή στατιστική)

# Γραμμική παλινδρόμηση

(3)

Έστω  $n$  διατεταγμένα ζεύγη της μορφής  $(x_i, y_i)$  ( $i=1, 2, 3, \dots, n$ ). Εάν τα ζεύγη αυτά θεωρηθούν ως σημεία κλειστού στο επίπεδο  $x-y$  λαμβάνουμε ένα νέφος σημείων ή διαγράμμα διασποράς (scatter plot).

Εάν η εξέλιξη ανάμεσα στα σημεία είναι γραμμικής φύσης η καμπύλη που προκύπτει έχει θετική ή αρνητική κλίση γα ισχύει η σχέση  $B \propto \theta$  συσχετισμός. Παράδειγμα



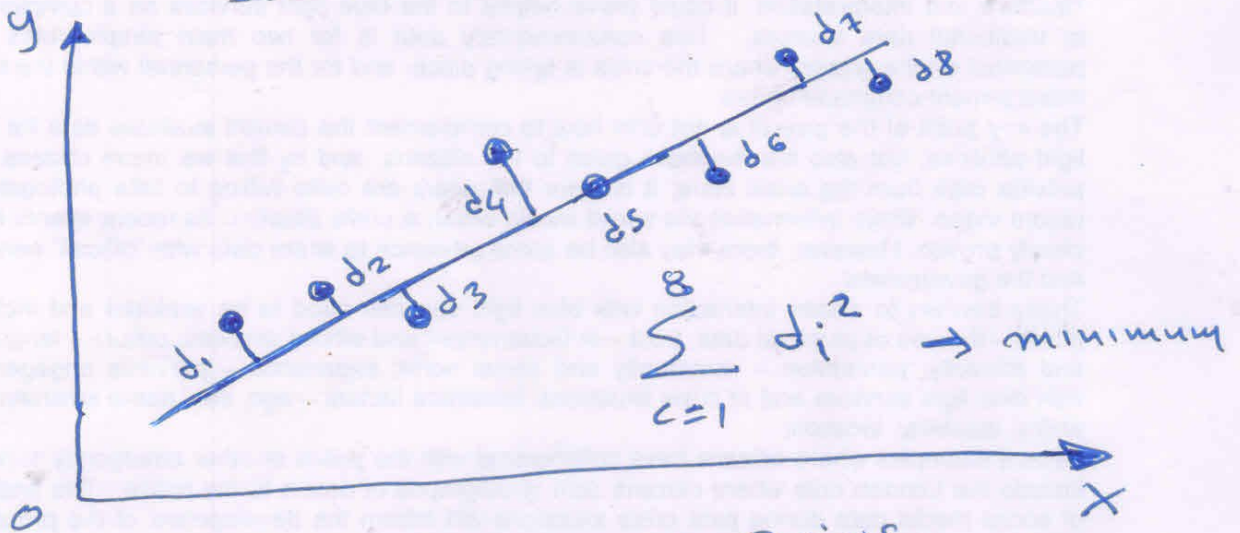
(A)

## ΕΥΘΕΙΑ ΕΞΑΧΙΣΤΩΝ ΤΕΤΡΑΓΩΝΩΝ

Έστω συντεταγμένες  $(x_i, y_i)$  ( $i=1, 2, 3, \dots, n$ ) που φέρνεται ως σχετιζόμενα γεγονότα τους  $y$  ο γραμμικός τύπος και κατά συνέπεια σχετιζόμενα με ευθεία. Ανάμεσα τους ευθεία ευθεία  $n$  στοιχεία χάρη στην επίλυση των εξισώσεων:

Το άθροισμα των τετραγώνων των απόστασεων  $d_i$  ( $i=1, 2, 3, \dots, n$ ) των παραπάνω σημείων / σημείων από αυτή την ευθεία είναι ελάχιστο

### Γεωμετρική Ερμηνεία



Εάν  $n$  εξισώσεις αυτές της ευθείας έχει τη μορφή

$$y = b_0 + b_1 x$$

Το πρόβλημα είναι να βρεθεί στην ευθεία των τιμών των συντελεστών  $b_0$  και  $b_1$

Απλά δεδομένα  $\forall$  ζεύγη  $(x_i, y_i)$

(5)

Οι τιμές των  $b_0$  και  $b_1$  προκύπτουν από τη λύση του συστήματος ( $v = N$ )

$$\sum_{i=1}^N y_i = N b_0 + b_1 \sum_{i=1}^N x_i$$

$$\sum_{i=1}^N x_i y_i = b_0 \sum_{i=1}^N x_i + b_1 \sum_{i=1}^N x_i^2$$

άπο άπου προκύπτει ότι (παράλειπεται οι βήματα στο άθροισμα για ευκολία)

$$b_1 = \frac{\sum x_i y_i - N \bar{X} \bar{Y}}{\sum x_i^2 - N \bar{X}^2} \text{ και στη συνέχεια}$$

$$b_0 = \bar{Y} - b_1 \bar{X} \text{ όπου } \bar{X} \text{ και } \bar{Y} \text{ οι}$$

μέσες όροι των  $x_i$  και  $y_i$  αντίστοιχα

ΣΤΙΣ παραπάνω εκφράσεις το  $b_1$  ονομάζεται συντελεστής παλινδρόμησης ή κλίση της ευθείας

και εκφράζει τη μεταβολή του  $y$  όταν το  $x$  μεταβληθεί κατά μία μονάδα.

- $b_1 > 0$  θετική εξάρτηση
- $b_1 < 0$  αρνητική εξάρτηση

Για κάθε τιμή του  $x_i$  υφίσταται ένα  
 τιμές  $y$ , η τιμή  $y_i$  που γαί:  $y \in \mathbb{R}$   $x_i$  σχμ-  
 φαι:  $y$  του  $f$  τύπος  $(x_i, y_i)$  που χρησιμο-  
 ποιήθηκε στον παραπάνω υπολογισμό με  
 η τιμή  $y_i^{\wedge}$  που προκύπτει αντικαθιστώντας  
 την τιμή του  $x_i$  στην εξίσωση της ευθείας  
 ελαχίστων τετραγώνων -  $\Theta$  είναι δηλαδή

$$y_i^{\wedge} = b_0 + b_1 x_i$$

Η διαφορά  $d_i = y_i - y_i^{\wedge}$  δεν είναι παρά το  
 σφάλμα που θεωρούμε προφουέρως.

ΠΑΡΑΔΟΧΕΣ

- Το σφάλμα  $d_i$  αποτελεί για τυχαία γετα-  
 βλήτη  $y$  φυσική γεγη τιμή.
- Η διανομή του σφάλματος είναι normal  
 για όλα τα  $x$
- Οι τιμές του σφάλματος είναι ανεξάρτητες  
 μεταξύ τους
- Το σφάλμα χαρακτηρίζεται από υπόνομή  
 υατδνομή

Η ευθεία ελαχίστων τετραγώνων χρησιμοποι-  
 είται για πρόβλεψη πάρφουέρως

Εάν το πλήθος των  $f_{ij}$  σημείων  $(x_i, y_j)$  είναι πάρα πολύ μεγάλο, προχωρούμε ως γνωστόν στην ομαδοποίησή τους. Έστω

$f_{ij}$  το πλήθος των  $f_{ij}$  σημείων  $(x_i, y_j)$

$f_{xi}$  το πλήθος των  $f_{ij}$  σημείων  $(x_i, \alpha)$  που έχουν ως πρώτο στοιχείο το  $x_i$  ανεξάρτητα από το  $y_j$

$f_{yj}$  το πλήθος των  $f_{ij}$  σημείων  $(\alpha, y_j)$  που έχουν ως δεύτερο στοιχείο το  $y_j$  ανεξάρτητα από το  $x_i$ .  $\Theta \alpha$  είναι το  $\alpha$

$$b_1 = \frac{N \sum \sum f_{ij} x_i y_j - (\sum f_{xi} x_i)(\sum f_{yj} y_j)}{N \sum f_{xi} x_i^2 - (\sum f_{xi} x_i)^2}$$

και

$$b_0 = \bar{y} - b_1 \bar{x} \text{ όπου } \bar{x} \text{ και } \bar{y} \text{ οι μέσοι όροι}$$

των  $x_i$  και  $y_j$  αντίστοιχα

Η προσάρτηση των  σημείων  $(x_i, y_j)$  είναι άπο φυθία γλφει να γίνει και σε άλλες

υφάνει οπως πάρβωτη, υλφρβωτη ή

εκθετική υφάνει.

## Προσαρμογή παραβολών σε δεδομένα

(8)

Έστω σύνολο σημείων  $M_i (x_i, y_i)$  ( $i=1, 2, \dots, N$ ).  
Ζητείται παραβολή  $y \in \mathbb{R}$  ως εξής

$$y = b_0 + b_1 x + b_2 x^2$$

η οποία να προσαρμόζεται στα δεδομένα και  
οδηγείται στο γινώμενο συνολικό σφάλμα

Εδώ οι συντελεστές  $b_1$  και  $b_2$  ονομάζονται  
συντελεστές παραμόρφωσης ενώ  $b_0$  είναι ένας  
σταθερός όρος.

Εδώ οι παραπάνω συντελεστές προκύπτουν  
από τη λύση του συστήματος

$$\sum y_i = N b_0 + b_1 \sum x_i + b_2 \sum x_i^2$$

$$\sum x_i y_i = b_0 \sum x_i + b_1 \sum x_i^2 + b_2 \sum x_i^3$$

$$\sum x_i^2 y_i = b_0 \sum x_i^2 + b_1 \sum x_i^3 + b_2 \sum x_i^4$$

Ποιο είναι το καλύτερο σχήμα κατεύθυνσης? Υπάρχουν  
υψηλά κριτήρια αλλά γενικά αρκεί για αυτήν την παρ-  
τήρηση του διόρθωματος. διότι αν εδω εδω αν  
ηδηω σουφε πως τα επιπλέον σχηματίζουν ευθεία η  
παραβολή εφαρμόζεται αναλόγως.



## Μέσο τετραγωνικό σφάλμα

9

Αποτελεί τον αριθμητικό μέσο των τετραγώνων των σφαλμάτων  $d_i$ .

Για ευθεία βεταζινόγραμμα δεδομένα έχουμε

$$\begin{aligned}\sigma^2 &= \frac{1}{N} \sum d_i^2 = \frac{1}{N} \sum (y_i - \hat{y}_i)^2 = \\ &= \frac{1}{N} \left( \sum y_i^2 - b_0 \sum y_i - b_1 \sum x_i y_i \right)\end{aligned}$$

Για ταξινόμημένα δεδομένα η εξίσωση είναι

$$\begin{aligned}\sigma^2 &= \frac{1}{N} \left( \sum_j f_{y_j} y_j^2 - b_0 \sum_j f_{y_j} y_j - \right. \\ &\quad \left. - b_1 \sum_i \sum_j f_{ij} x_i y_j \right)\end{aligned}$$

Για των παραβόλη έχουμε

$$\begin{aligned}\sigma^2 &= \frac{1}{N} \left( \sum y_i^2 - b_0 \sum y_i - b_1 \sum x_i y_i - \right. \\ &\quad \left. - b_2 \sum x_i^2 y_i \right)\end{aligned}$$

Η τετραγωνική ριζή  $\sigma$  του μέσου τετραγωνικού σφάλματος ονομάζεται ωλικό σφάλμα εκτίμησης.

ΣΥΣΧΕΤΙΣΗ Μέτρα των έντασης του βαθμού 10  
εξαρτήσεως ανάμεσα σε μεταβλητές.

Συνδιακύμανση (Covariance)

Μέτρο της πραγγυικής συνρροφείας μεταξύ  
δύο ποσοτικών μεταβλητών  $X$  και  $Y$ .

Έστω μεταβλητή  $X$  με σφισφάτα  $x_1, x_2, \dots, x_N$   
και μεταβλητή  $Y$  με σφισφάτα  $y_1, y_2, \dots, y_N$

Ορίζουμε (για ταξινόγητα δεδομένα)

$$\text{Cov}(X, Y) = \frac{1}{N} \sum (x_i - \bar{X})(y_i - \bar{Y}) =$$
$$\frac{1}{N} \sum x_i y_i - \bar{X}\bar{Y}$$

όπου  $\bar{X}$  και  $\bar{Y}$  οι μέσοι όροι των  $X$  και  $Y$ .

Εάν τα δεδομένα είναι ταξινόγητα θα έχουμε

$$\text{Cov}(X, Y) = \frac{1}{N} \sum_i \sum_j f_{ij} (x_i - \bar{X})(y_j - \bar{Y}) =$$
$$\frac{1}{N} \sum_c \sum_j f_{ij} x_i y_j - \bar{X}\bar{Y}$$

Εάν τα δεδομένα σχετίζονται με πραγγυικό τρόπο,  
ο συντελεστής  $b_1$  της ευθείας ελαχίστων τετρα-  
γωνων πραφείας ως  $b_1 = \frac{\text{Cov}(X, Y)}{\sigma_X^2}$

Οπου  $\sigma_X^2 = \text{cov}(X, X)$  η διακύμανση - (II)  
 δηλ η μεταβλητότητα του  $X$ . Από τον  
 ορισμό είναι

$$\sigma_X^2 = \frac{1}{N} \sum X_i^2 - (\bar{X})^2$$

Επειδή είναι  $\sigma_X^2 > 0$ , το πρόσημο του συντελεστή  
 β1 εξαρτάται από το πρόσημο του  $\text{cov}(X, Y)$ .  
 Λαμβάνοντας υπόψη ότι το β2 είναι η κλίση  
 της ευθείας ελαχίστων τετραγώνων, έχουμε:

- Εάν είναι  $\text{cov}(X, Y) > 0$  οι μεταβλητές μετα-  
 βάζονται αναλόγως (η αύξηση της γιά, οδηγεί σε αύξηση της άλλης).
- Εάν είναι  $\text{cov}(X, Y) < 0$  οι μεταβλητές  
 μεταβάζονται αντίστροφα αναλόγως  
 (η αύξηση της γιά) οδηγεί σε μείωση της  
 άλλης).
- Εάν είναι  $\text{cov}(X, Y) = 0$  οι μεταβλητές  
 είναι ανεξάρτητες μεταβλητές.

### ΙΔΙΟΤΗΤΕΣ

$$\text{cov}(X, Y) = \text{cov}(Y, X), \text{cov}(X, -Y) = -\text{cov}(X, Y)$$

$$\text{cov}(X+Y, Z) = \text{cov}(X, Z) + \text{cov}(Y, Z)$$

$$Z = a + bX, W = c + dY \Rightarrow \text{cov}(Z, W) = bd \text{cov}(X, Y)$$

Η χρησιμότητα της συνδιακύμανσης  $\text{COV}(X, Y)$  είναι η περιορισμένη, ελπίδα συνόψι-  
 εται από τις γωνίες των  $X$  και  $Y$ , ενώ αν  
 είναι  $\text{COV}(X, Y) = 0$  ποτε ο ένας από τους δύο  
 συσχετίζεται πως οι μεταβλητές είναι άσχετες

ΣΥΝΤΕΛΕΣΤΗΣ ΣΥΣΧΕΤΙΣΗΣ

$$r = \frac{\text{COV}(X, Y)}{\sigma_X \sigma_Y} = \frac{\text{COV}(X, Y)}{\sqrt{\sigma_X^2 \sigma_Y^2}} = \frac{\text{COV}(X, Y)}{\sqrt{\text{COV}(X, X) \text{COV}(Y, Y)}}$$

Για αταβινόμενα δεδομένα έχουμε

$$r = \frac{(1/N) \sum x_i y_i - \bar{X} \bar{Y}}{\sqrt{\frac{\sum x_i^2}{N} - \bar{X}^2} \sqrt{\frac{1}{N} \sum y_i^2 - \bar{Y}^2}} = \frac{N \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{N \sum x_i^2 - (\sum x_i)^2} \sqrt{N \sum y_i^2 - (\sum y_i)^2}} = \frac{\sum x_i y_i - N \bar{X} \bar{Y}}{N \sigma_X \sigma_Y}$$

τύπος του Pearson

Για τα ταξινομημένα δεδομένα έχουμε

(13)

$$r = \frac{N \sum \sum f_{ij} x_i y_j - \sum f_{x_i} x_i \sum f_{y_j} y_j}{\sqrt{N \sum f_{x_i} x_i^2 - (\sum f_{x_i} x_i)^2} \sqrt{N \sum f_{y_j} y_j^2 - (\sum f_{y_j} y_j)^2}}$$

I διαλύστες συντελεστή συσχέτισης

- Είναι αδιάστατο μέγεθος
- Η αίρεση ρίχνει στο διάστημα  $(-1, 1)$
- Για  $r = +1$  ή  $r = -1$  έχουμε τέλεια  $[+]$  ή  $[-]$  συσχέτιση (γραμμική πορεία θετική ή αρνητική)
- Εάν  $r = 0 \Rightarrow$  οι μεταβλητές είναι αβυσχέτιστες

# ΙΒΑΛΙΚΕΣ ΕΝΝΟΙΕΣ ΠΙΘΑΝΟΤΗΤΩΝ

①

Συνοχά : οριστός, υποσυνολο, κούθολιό,  
ένωση, ζούμ, διαφορά, καρτεσιάνο  
γινούενο  
κενό συνοχό

ιδιότητες ένωσης και τογής

Προσεταιριστική

Επιμεριστική

Υεταθετική, ουβέτερο στοιχείο

Σενα συνοχά

Κανόνες de Morgan

$$\overline{A \cup B} = \bar{A} \cap \bar{B}$$

$$\overline{A \cap B} = \bar{A} \cup \bar{B}$$

ΜΑΘΗΜΑ

6

Στατιστική Τυχαίοτητα, Χόος

Πείρατα τύχης

ορίζεται ως γία διαδιδάβιδ η οποία αν  
κδι εα αν αλαγβάνεται γε τον ίδιο τρόπο και  
κδιω αλο τις ίδιες κεριβώς συνθήκες δεν  
γπορούφε να προβλεγούγε το αποτέλεσμα  
λου προυνατεί αλο για ελανδαμυ

πχ Ριγμ Ζαριού

(διαγχαμ Bohr - Einstein)

Δειγματοχώρος  $\Omega$  ενός πειράματος (2)

Τυχής ονομάζεται το σύνολο όλων των δυνατών αποτελεσμάτων του ΠΧ

$$Ζ\alpha\pi = \{1, 2, 3, 4, 5, 6\}$$

$$Κερα = \{κ, γ\} \quad \begin{array}{l} \text{Κορώνα} \\ \text{Γράμματα} \end{array}$$

Κάθε δυνατό αποτέλεσμα

ενός πειράματος λέγεται από γεγονός

ή ενδεχομένο

Ενώ

Γεγονός ή ενδεχομένο ονομάζεται κάθε υποσύνολο  $A$  του δειγματοχώρου  $\Omega$

ΠΧ Ζαρι :

Αλλά ενδεχομένα : 1, 2, 3, 4, 5, 6

Γεγονότα : ΠΧ

$$A = \{1, 4\} \quad B = \{3, 5, 6\}$$

Δειγματοχώρος 2 νομισμάτων

$$\Omega = \{κκ, κγ, γκ, γγ\} \quad \begin{array}{l} \text{δίδυμη} \\ \text{κεραίατων} \end{array}$$

$$\underline{\Omega} = \{κκ, κγ, γγ\} \quad \begin{array}{l} \text{χωρίς δίδυμη} \\ \text{κεραίατων} \end{array}$$

οποτε  $κγ = γκ$

## ΟΡΙΣΜΟΙ

(3)

Αν σε μία επανάληψη ενός πειράματος προκύψει το αποτέλεσμα  $\omega \in \Omega$

τότε

- Αν είναι  $\omega \in A$  λέμε πως πραγματοποιείται το γεγονός  $A$
- Αν είναι  $\omega \notin A$  λέμε πως δεν πραγματοποιείται το γεγονός  $A$
- Αν είναι  $\omega \in A$  και  $\omega \in B$  ή ισοδύναμα εάν  $\omega \in (A \cap B)$  λέμε πως πραγματοποιούνται ταυτόχρονα τα γεγονότα  $A$  και  $B$ .
- Αν είναι  $\omega \in A$  ή  $\omega \in B$  ή ισοδύναμα εάν  $\omega \in (A \cup B)$  λέμε πως πραγματοποιείται τουλάχιστον ένα από τα  $A$  και  $B$
- Αν είναι  $\omega \in A$  και  $\omega \notin B$  ή ισοδύναμα  $\omega \in (A - B)$  λέμε πως πραγματοποιείται το  $A$  αλλά όχι το  $B$
- Αν είναι  $\omega \notin A$  και  $\omega \notin B$  ή ισοδύναμα  $\omega \in \bar{A} \cap \bar{B}$  ή  $\omega \in \overline{A \cup B}$  λέμε πως δεν πραγματοποιείται κανένα από τα  $A$  και  $B$



4 Ανάλυση των εννοιών της συνάρτησης

$f: A \rightarrow B$  η εικόνα ορισμού  
σύνολο τιμών

$y = f(x)$ , εικόνα,  $x$  προς  $y$ , επί

### ΠΙΘΑΝΟΤΗΤΑ

Εάν σε πειχτούμε ένα κέρμα  $N$  φορές και καταγράψουμε το πλήθος  $K$  του αποτελέσματος

"Γράμματα"

ο λόγος  $f = \frac{K}{N}$  ονομάζεται  
Γραμμική Συχνότητα

Η επί τοις 100 γραμμική συχνότητα ορίζεται ως

$$f_i \% \Rightarrow f_i \times 100$$

για το γεγονός  $\neq i$ .

Παράδειγμα Εάν σε πειχτούμε ένα κέρμα

100 φορές και τις 63 φορές έχουμε

"Γράμματα".

$$\text{Άρα } f_g = \frac{63}{100} = 0.63 \Rightarrow f_g \% = 63$$

Επειδή το άλλο αποτέλεσμα  $\Theta$  είναι γόνο

$$\text{"υπογράμματα"} \Rightarrow f_x = \frac{37}{100} = 0.37.$$

# Παράδειγμα

5

Εάν οι πιθανότητες για 50 φάρμακα να  
εξοφλή

- 1 17 φάρμακα  $\Rightarrow f_1 = \frac{17}{50} = 0.34$
- 2 8 φάρμακα  $\Rightarrow f_2 = \frac{8}{50} = 0.16$
- 3 7 φάρμακα  $\Rightarrow f_3 = \frac{7}{50} = 0.14$
- 4 10 φάρμακα  $\Rightarrow f_4 = \frac{10}{50} = 0.20$
- 5 3 φάρμακα  $\Rightarrow f_5 = \frac{3}{50} = 0.06$
- 6 5 φάρμακα  $\Rightarrow f_6 = \frac{6}{50} = 0.12$

Παρατηρούμε ότι

$$f_1 + f_2 + f_3 + f_4 + f_5 + f_6 = \frac{17 + 8 + 7 + 10 + 3 + 5}{50} = \frac{50}{50} = 1$$

Αρα  $F = \sum_{i \in \Omega} f_i = 100\%$  (Βεβαιότητα)

Είναι βεβαιότητα να αποτελέσει ο φάρμακος  
1 2 3 4 5 6.

	<u>Πιψεις</u>	<u>Σχετική Γυχνότητα</u>
Ας ενοστρεφούσε στο	20	0.652
λεπτά του κέρδιου	40	0.574
και εσω λω μ σχετική	60	0.510
γυχνότητα για <u>δράφατα</u>	80	0.489
είναι	100	0.401
	120	0.509 .....

Αυτή η σχετική συχνότητα ονομάζεται

(6)

"πιθανότητα να πραγματοποιηθεί"

Ορισμός ορίζεται ως πιθανότητα σε έναν

δείγματο χώρο  $\Omega$  για συνάρτηση που αντιστοιχεί σε κάθε γεγονός  $A \in \mathcal{A}$  έναν θετικό αριθμό

$P(A)$  τέτοιον ώστε

$$1) \quad P(A) \geq 0$$

$$2) \quad P(\Omega) = 1$$

Για δύο γεγονότα  $\{A_1, A_2\}$  ένα μετά το άλλο ανά 2

δηλαδή τέτοια ώστε  $A_1 \cap A_2 = \emptyset$  είναι

$$P(A_1 \cup A_2) = P(A_1) + P(A_2)$$

Γενικά

$$P(A_1 \cup A_2 \cup \dots \cup A_n) = P(A_1) + P(A_2) + \dots + P(A_n)$$

Δηλαδή

Η πιθανότητα εμφάνισης ενός ενδεχόμενου  $A$  είναι ίση με το άθροισμα των πιθανοτήτων των αλληλ αποκλειστικών ενδεχομένων από τα οποία αποτελείται.

Επομένως

$$P(A) = \sum_{\omega \in A} P(\omega)$$

Εξ ορισμού

(7)

$$\sum_{\omega \in \Omega} P(\omega) = 1$$

Εάν τα στοιχεία ενός πεπερασμένου δειγματοχώρου  $\Omega$  ψε πληθος στοιχείων  $n(\Omega)$  έχουν ίση πιθανότητα (ισοπιθάνος δειγματοχώρος) τότε η πιθανότητα ενός γεγονότος  $A$  του  $\Omega$  ή ο πληθος  $n(A)$  είναι

$$P(A) = \frac{n(A)}{n(\Omega)}$$

Νόμοι πιθανότητων

1)  $P(\bar{A}) = 1 - P(A)$

2)  $P(\emptyset) = 0$

3)  $A \subseteq B \Rightarrow P(A) \leq P(B)$

4)  $P(A) \leq 1$

5)  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

6)  $P(A - B) = P(A) - P(A \cap B)$

---

Η υδρανοφή συχνοτήτων δεν αποτελεί παρά έναν πίνακα στον οποίο αναφέρονται οι τιμές της τυχαίας μεταβλητής και οι αντίστοιχες συχνοτήτες εμφάνισης.

<u>ΤΥΧ</u>	Πλήθος οιμογενειών
Αριθμός παιδιών ανά οικογένεια	
0	50
1	250
2	360
3	220
4	80
5	40
	<hr/>
	900

Η σχετική συχνότητα

Εκχ 250/900 για οικογένειες με 1 παιδί) νοείται ως πιθανότητα. Ο αντίστοιχος πίνακας είναι γνωστός ως πίνακας πιθανοτήτων.

Οι τυχαίες μεταβλητές οπωσδήποτε αναφέρονται σε προηγούμενο γάμμα (α γλαφεί να είναι είτε διακριτές όταν λαμβάνουν τιμές από ένα πεπε-  
ρισμένο σύνολο ή αλλιώς διακριτές γκαβύ ταν)  
Τιμές είτε συνεχείς όταν λαμβάνουν οποιαδήποτε  
τιμή σε ένα διάστημα ή σύνολο διαστημάτων.

## ΑΓΩΓΕΣ ή ΔΙΑΦΗΤΕΙ ΥΕΤΑΒΛΗΤΕΣ

(2)

ΜΕΣΗ (mean) ή αναγενομένη (expected) τιμή  
ή γαθονυατιμή εδνιδά ή προοδονία (expectation)

$$E(X) = \mu = \sum x P(x)$$

οπου  $P(x)$  ή πιθανότητα εμφάνιςι της τιμής  $x$ .

Διασάγανση (variance)

$$\sigma^2 = \sum [x^2 \cdot P(x)] - \mu^2$$

Τυπιή απόκλιση  $\sigma = \sqrt{\sigma^2}$

Κατανομές πιθανότητας

4) Διωνυμιή κατανομή (Binomial distribution)  
ή διαεργαία Bernoulli

$$X \sim B(n, p)$$

- Χρησιφαποιείται σε σειρά πδνογοιότυπων πειραγάτων οπω είναι  $n$  πηγμ ενός νογιόγματος.
- Για καθε πειραφα υπάρχουν γόνα δύο αποτελέγγατα: επιτυχία (γε πιθανότητα  $p$ ) ή αποτυχία (γε πιθανότητα  $q$ ). Αρα  $p + q = 1$
- Τα πειραφατα είναι ανεξάρτητα γεταβί τους.
- Η τυχαία γεταβλημή  $X$  γερα το πλίοθος των επιτυχιών και κποφι να παρρι τιφές 0 εως  $n$ .

3  
Η πιθανότητα να έχουμε ακριβώς  $x$   
επιτυχίες σε  $n$  επαναλήψεις του πειράματος  
δίδεται από την έκφραση

$$P(X=x) = \binom{n}{x} p^x q^{n-x} \quad x=0,1,2,\dots,n$$

όπου  $n$  το πλήθος των επαναλήψεων και

$$\binom{n}{x} = \frac{n!}{x!(n-x)!}$$

Παράμετροι

Μέση τιμή  $E(X) = np$

Διακύμανση  $Var(X) = \sigma^2 = npq$

Τυπική απόκλιση  $\sigma = \sqrt{npq}$

ΠΑΡΑΔΕΙΓΜΑ

Το 5% των οδύων φορτηγών στην Αγγλία  
είναι γυναικες

Επιλέγουμε στην τύχη 10 οδύους.

Εφόσον το αποτέλεσμα είναι Άνδρας / Γυναίκα  
 $n$  κτηνοτρόφι είναι διωνυμική

1) Ποια είναι η πιθανότητα ώστε 2 οδύοι  
να είναι γυναίκες?

2) Ποια είναι η πιθανότητα κάποιος οδύος να  
την είναι γυναίκα?





- Η πιθανότητα να συμβεί ένα γεγονός είναι  $n$  ίδια για δύο διαδοχικά ισού πλάτους.
- Η πράξη ατονοίηση  $n$  οχι ενός γεγονότος σε ένα διάστημα είναι ανεξάρτητη από την πράξη ατονοίηση των  $n$  οχι σε άλλο άλλο διάστημα.

Ιδιότητες κατανομής Poisson

Μεση τιμή  $E(X) = \lambda = np$

Διακύμανση  $Var(X) = \sigma^2 = \lambda$

Τυπική απόκλιση  $\sigma = \sqrt{\lambda}$

ΠΑΡΑΔΕΙΓΜΑ

Ταξιδιώτες φτάνουν τυχαία και ανεξάρτητα στο σταθμό ελεγχου ενός αεροδρομίου. Ο υφιστάμενος αξιωματικός είναι  $\lambda = 10$  ταξιδιώτες λ.ε.π.

- Ποια είναι η πιθανότητα να φτάσει κανένας ταξιδιώτης κατά τη διάρκεια ενός λ.ε.π.?

$$P(X=0) = 10^0 \frac{e^{-10}}{0!} = 0.000045399 = 0.00454\%$$

- Ποια είναι η πιθανότητα να φτάσουν το πολύ τρεις ταξιδιώτες?

$$P(X \leq 3) = P(X=0) + P(X=1) + P(X=2) + P(X=3) = 10^0 \frac{e^{-10}}{1!} + 10^1 \frac{e^{-10}}{1!} + 10^2 \frac{e^{-10}}{2!} + 10^3 \frac{e^{-10}}{3!} = 1.03394\%$$

## Κανονική κατανομή ή κατανομή Gauss

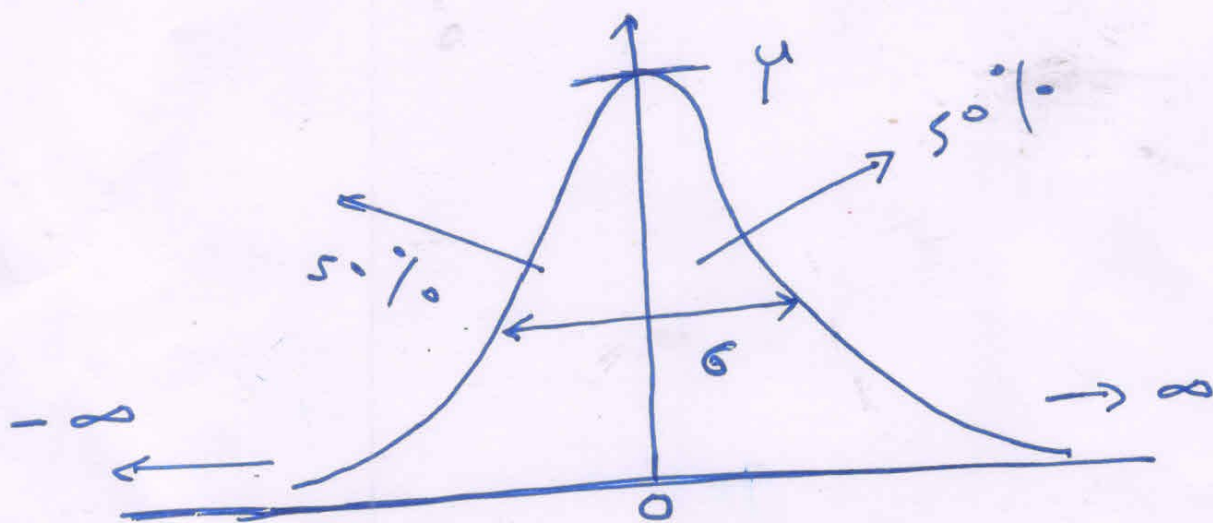
(6)

Αποτελεί μία από τις πιο σημαντικές κατανομές συνεχούς μεταβλητής και ορίζεται ως

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[ -\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2 \right]$$

όπου  $\mu$  ο αριθμητικός μέσος όρος και  $\sigma$  η τυπική απόκλιση.

Η γραφική της παράσταση έχει τη μορφή



- Ο αριθμητικός μέσος όρος αποτελεί το υψηλότερο σημείο της καμπύλης και αποτελεί επίσης την διαμέσο και μέτρο συχνοτήτων. Μπορεί να λάβει θετική, αρνητική ή μηδενική τιμή.
- Η καμπύλη είναι ψωδωμοειδής γι' αυτό λέγεται κορυφή και είναι συμμετρική ως προς την ευθεία του αριθμητικού μέσου  $x = \mu$ .

- Η κατανομή τείνει ασυμπτωτικά στο f όταν  $x \rightarrow \infty$  και  $x \rightarrow -\infty$  (8)
- Η τυμική απόκλιση καθορίζει το πλάτος της κατανομής ενώ το συνολικό εμβαδόν της είναι 100%  $\forall \epsilon$  της συνάρτησης.
- Το 80% των εμβαδών αντιστοιχεί σε διαστήματα  $\pm \sigma$ , το 95,99 σε διαστήματα  $\pm 2\sigma$  ενώ το 99,72 σε διαστήματα  $\pm 3\sigma$ .
- Η πιθανότητα να έχει η μεταβλητή  $X$  τιμή σε ένα διάστημα είναι ίση με το αντίστοιχο εμβαδόν της κατανομής.

Γράψουμε  $X \sim N(\mu, \sigma)$

Η πιθανότητα για ένα διάστημα  $[a, b]$  υπολογίζεται ως

$$P(a \leq X \leq b) = \int_a^b f(x) dx$$

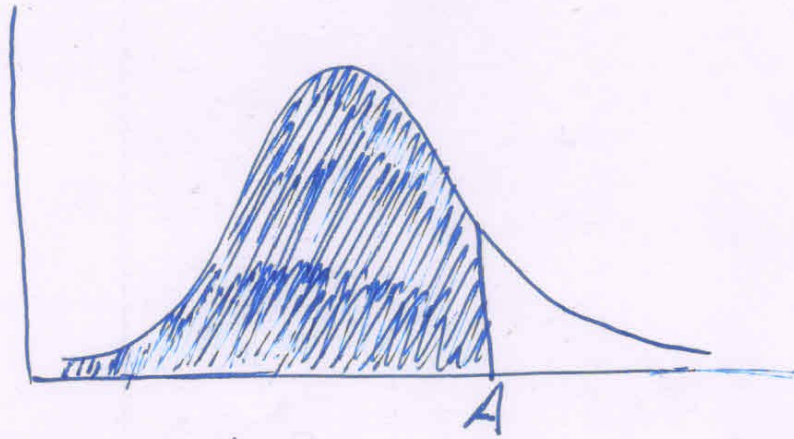
Προκειμένου να απλοποιήσουμε την διαδικασία υπολογισμού του ολοκληρώματος ορίζουμε τη μεταβλητή

$$Z = \frac{X - \mu}{\sigma}$$

με την αντίστοιχη κατανομή που προκύπτει να χαρακτηρίζεται από όλες τις ιδιότητες της αρχικής κατανομής αλλά να διαθέτει μηδενικό αριθμητικό μέσο όρο και μοναδικά τυμική απόκλιση ( $\sigma = 1$ )

Γράψουμε  $Z \sim N(0, 1)$

Το Βαθμίο πληθυντικά αυτής της νέας κατανομής (9) είναι πως οι τιμές της γίνουν να βρεθούν σε λιγότερη και δεν αλλοτείνει ο υπολογισμός του ολοκληρωτός.



$$P(X \leq A) = \text{το συνολικό εμβαδόν}$$

### Δειγματοληπτική κατανομή αριθμητικού μέγεθ

Έστω πληθυσμός  $N$  και δείγμα μέγεθος  $n$  που χρησιμοποιείται για την διεξαγωγή κάποιων ερευνών. Για να πάρουμε ακριβή αποτελέσματα ενδέχεται να χρειαστεί να πάρουμε πολλές φορές θεωρητικά πολλά δείγματα από τον ίδιο πληθυσμό με το ίδιο μέγεθος  $n$ . Το κάθε δείγμα έχει διαφορετικό αριθμητικό μέγεθος  $\bar{x}$  και η κατανομή πιθανότητας όλων των πιθανών τιμών του δειγματικού μέγεθους ονομάζεται δειγματική κατανομή.

Παράγωγοι δειγματοληπτικής  
κατανόησης αριθμητικού γέγον

Αναμενόμενη τιμή του  $\bar{X}$ :  $E(\bar{X}) = \mu$

( $n$  γέγον τιμή του πληθυσμού από τον οποίο  
ελήφθη το δείγμα).

Τυπική απόκλιση του  $\bar{X}$

Εάν είναι  $n/N > 0.05 \Rightarrow \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$

Εάν  $N \rightarrow \infty$  ή  $N$  πεπερασμένο τε

$n/N \leq 0.05$  τότε  $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$

όπου:  $\sigma$  η τυπική απόκλιση του πληθυσμού  
 $n$  το κοινό γέγον των επαναλαμβανόμενων  
δείγματων και  
 $N$  το γέγον του πληθυσμού

Όταν το  $n$  είναι μεγάλο, η δειγματοληπτική  
κατανόηση του αριθμητικού γέγον  $\bar{X}$  μπορεί  
να προσεγγιστεί από την κανονική κατανομή  
(central limit theorem).

Μεγάλο δείγμα  $\iff n > 30$

---

# Επιθυμητές ιδιότητες γεμάτων για στατιστική εκτίμηση

2

- Ανεροληψία: Ο γετός ορος της δείγματος πληθυσμιακής κατανομής των αριθμητικών γεμάτων που προσεχεται από τον ίδιο πληθυσμό είναι ο ίδιος με τον αριθμητικό γετό του πληθυσμού.
- Αποτελεσματικότητα: Σχετίζεται με το μέγεθος του τυπικού σφάλματος.
- Συνέπεια: Ο εκτιμητής θεωρείται συνεπής εάν καθώς αυξάνει το μέγεθος του δείγματος, τότε το προς γετόμενη γετόμενη γίνεται ίσο με την τιμή του αντίστοιχου γετού του πληθυσμού.
- Ικανότητα: Σχετίζεται με τον όγκο των πληροφοριών που επηρεάζονται στη στατιστική εκτίμηση.

## Τύποι Στατιστικής Εκτίμησης

### Εκτίμηση σε σημείο (point estimation)

Εκτιμούντε την τιμή μιας άγνωστης παράμετρος του πληθυσμού από μια τυχαία τιμή της αντίστοιχης παράμετρος της δείγματος πληθυσμιακής κατανομής.

Συγκεκριμένοι εκτιμητές είναι ο αριθμητικός γετός, η τυπική απόκλιση και το ποσοστό του δείγματος.

Η τιμή  $|\bar{X} - \mu|$  ονομάζεται δεδειγματοληπτική (3)  
υπό σφάλμα και όσο πιο γύρω είναι η τιμή του  
σφάλματος αυτού, τόσο καλύτερη είναι η εκτίμηση

### Εκτίμηση σε διάστημα (interval estimation)

Η αληθινή τιμή μιας παραμέτρου του πληθυσμού  
αναμένεται να βρεθεί μέσα σε ένα διάστημα  
τιμών το οποίο ονομάζεται διάστημα εμπιστοσύ-  
νης (confidence interval)  $1 - \alpha$  το οποίο  
υποδηλώνει την πιθανότητα να περιλάβει ένα  
διάστημα η προεπιλεγμένη τιμή της παραμέτρου του  
πληθυσμού.

Η πιθανότητα σφάλματος είναι  $\alpha$  ( $0 < \alpha < 1$ ) και  
ονομάζεται επίπεδο σημαντικότητας.

Εάν είναι  $\alpha = 0.05$ , ο συντελεστής εμπιστοσύνης

είναι  $1 - \alpha = 0.95$  και το διάστημα εμπιστοσύνης  
ονομάζεται διάστημα εμπιστοσύνης 95%.

### ΔΙΑΣΤΗΜΑ ΕΜΠΙΣΤΟΣΥΝΗΣ ΓΙΑ ΤΗ ΜΕΣΗ ΤΙΜΗ ΠΛΗΘΥΣΜΟΥ

Στον υπολογισμό του λάτβαν και να σημειωθεί:

- Η γορφή της κατανομής του πληθυσμού (κανονική ή κατά προσέγγιση κανονική)
- Το μέγεθος του δείγματος (μικρό ή μεγάλο)
- Η γνώση ή όχι της τυπικής απόκλισης του πληθυσμού.

- Το εάν ο πληθυσμός  $N$  είναι πεντασθενός ή απείρος.

(4)

- Μεγάλο δείγμα ( $n > 30$ ), γνωστό  $\sigma$ :  
πεντασθενός πληθυσμός,  $n/N > 0.05$

$$\bar{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$$

- : κατώτατο όριο εφελκυστικής

+ : ανώτατο όριο εφελκυστικής

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} \text{ τυπικό σφάλμα του μέσου}$$

$$E = Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \text{ μέγιστο σφάλμα εκτίμησης}$$

Όταν  $N \rightarrow \infty$  ή  $N$  πεντασθενός + ε  
 $n/N < 0.05$

το διάστημα είναι

$$\bar{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

- Μεγάλο δείγμα ( $n > 30$ ),  $\sigma$  άγνωστο

$$\bar{X} \pm Z_{\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}, \quad \hat{\sigma} = s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

Για  $N \rightarrow \infty$  ή  $N$  πεντασθενός και

$$n/N < 0.05:$$



Το διαστήμα είναι

5

$$\bar{X} \pm Z_{\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}}$$

Επειδή ΔΕΝ γνωρίζουμε την τυπική απόκλιση του πληθυσμού, χρησιμοποιούμε ως την τυπική απόκλιση του δείγματος.

Εδώ γίνεται για εκτιμώμενο τυπικό σφάλμα του μέσου  $\bar{X}$  και εκτιμώμενο μέγεθος σφάλμα επιτήρησης.

- Μικρό δείγμα ( $n \leq 30$ )  
Κανονική ή κατά προσέγγιση κανονική κατανομή,  $\sigma$  άγνωστο

### Κατανομή Student's t (W.S. Gosset)

- 1) Είναι συστηματική σφαλ και  $n$  κανονική κατανομή
- 2) Είναι διασπορά για κατά διαφορετικό μέγεθος δείγματος Επειδή διαφέρουν οι Βαθμοί Ελευθερίας
- 3) Όσο πιο μεγάλο είναι το μέγεθος του δείγματος τόσο πιο υψηλή είναι η ακρίβεια της πρόβλεψης παράστασης.
- 4) Σε σχέση με την κανονική κατανομή είναι πιο χαλαρή αλλά πιο ακριβής
- 5) Όταν οι Βαθμοί Ελευθερίας αυξάνονται προσεγγίζονται την κανονική κατανομή.

## Διαστήματα εφαιγοσυνών

(5)

$$\bar{X} \pm t_{\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}} \quad \text{όπου } \hat{\sigma} = s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

Το μέγεθος του δείγματος που πρέπει να επιλέξετε για να έχετε τα μέγιστα επιθυμητά βάρη επιτημών  $E$  σε δεδομένο  $\alpha$  και γνωστή τυπική απόκλιση είναι

$$n = \frac{(Z_{\alpha/2})^2 \sigma^2}{E^2}$$

Αυτός ο τύπος χρησιμοποιείται και όταν δεν είναι γνωστή η τυπική απόκλιση αλλά έχουμε επιλεγεί κάποια τιμή.

Σε όλα τα παραπάνω τα εύρος του διαστήματος εφαιγοσυνών εξαρτάται από:

- Το μέγεθος του δείγματος.
- τη διασπορά του πληθυσμού
- το συντελεστή εφαιγοσυνών.

Γίνεται τόσο μικρότερο όσο

- μεγαλύτερο είναι το μέγεθος του δείγματος
- μικρότερη είναι η διασπορά
- μικρότερος είναι ο συντελεστής εφαιγοσυνών.

## Παράδειγμα Α (Μεγάλο δείγμα, $\sigma$ γνωστό) (7)

Εταιρεία λωλοθούς υαλοκαθαρτήτων  
Θέλει να επιβεβαιώσει την μέση διάρκεια ζωής τους.

Δεδομένα Τυπική απόκλιση  $\sigma = 6$  μήνες

$N \rightarrow \infty$  (κανονική κατανομή)

Μέση διάρκεια ζωής  $\bar{x} = 21$  (μήνες)

$$\alpha = 0.05 \rightarrow 1 - \alpha = 0.95$$

(Μεγάλο δείγμα,  $\sigma$  γνωστό).

ΛΥΣΗ  $Z_{\alpha/2} = Z_{0.025} = 1.96$  (από πίνακα).

$$E = Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = 1.96 \frac{6}{\sqrt{100}} = 1.176.$$

Αρα

$$\left. \begin{aligned} \bar{x}_{\min} &= 21 - 1.176 = 19.824 \\ \bar{x}_{\max} &= 21 + 1.176 = 22.276 \end{aligned} \right\} \text{ με βεβαιότητα } 95\%.$$

## Παράδειγμα Β (Μεγάλο δείγμα, $\sigma$ άγνωστο)

Μεγεθος πληθυσμού  $N = 700$  ομογενείς

Μεγεθος δείγματος  $n = 50$  ομογενείς (Μεγάλο)

Κατανομή ετήσιου εισοδήματος  $\mu$   $\in$  στοχοποιημένη

Κανονική κατανομή,  $1 - \alpha = 0.9 \Rightarrow \alpha = 0.1$

Από το δείγμα προκύπτει:

$$\text{ΜΕΓΑ ΕΙΣΟΔΗΤΑ } \bar{X} = 11800 \text{ €}, \hat{\sigma} = 950 \text{ €}$$

Ποιο είναι το διάστημα εμπιστοσύνης?

ΛΥΣΗ είναι  $n/N = \frac{50}{700} = 0.071 > 0.05$   
η γεράλα

Αρα

$$\bar{X} \pm Z_{\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}, \quad \hat{\sigma} = S = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

$$Z_{\alpha/2} = Z_{0.05} = 1.64 \quad (\alpha \text{ no nivdymfi})$$

$$E = Z_{\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}} = 1.64 \frac{950}{\sqrt{50}} = 220.3$$

Αρα

$$\bar{x} - Z_{\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}} = 11800 - 220.3 \sqrt{\frac{700-50}{700-1}}$$

$$= 11800 - 220.3 \times 0.9693 =$$

$$= 11587.5$$

ομοίως για το (+):

$$\bar{x} + Z_{\alpha/2} = 11800 + 220.3 \times 0.9693 = 12.0125$$

Αρα με βεβαιότητα 90%, το μέσο ετήσιο εισόδημα των οικογενειών είναι από

$$\underline{\underline{11587.5 \text{ €} \text{ ως } 12.012.5}}$$

# Διαφορά γέγων τιμών

9

Δύο δείγματα είναι εξαρτημένα ή εξισωμένα κατά φύση, όταν το ένα επιλεγεται τυχαία και το άλλο επιλέγεται ισότιμο γέγων πρώτο.

Πχ για ένα δείγμα 10 τυχαίων γαθώνων καταγράψετε:

- Βαθμό γαθώνων σταθμής
- IQ
- Άνοιξις
- Ποσότητα ελιφίνου γούδα

Στο δεύτερο δείγμα επιλεγούτε 10 γαθούς για τους οποίους ισχύουν οι ίδιες τιμές IQ, άνοιξις και ποσότητας ελιφίνου γούδα και καταγράψετε το βαθμό της σταθμής

## ΠΙΝΑΚΑΣ ΔΕΔΟΜΕΝΩΝ

Άνοιξις	IQ	Μορφή	Βαθ1	Βαθ2
			$X_{11}$	$X_{21}$
			$X_{12}$	$X_{22}$
			$X_{13}$	$X_{23}$

Ανεξάρτητα δείγματα

Επιλεγούναί τε τυχαία δείγματα

$$d_i = X_{1i} - X_{2i}$$

Ζευγαρωτά δεδομένα

ο αριθμητικός μέσος των διαφορών είναι

$$\bar{d} = \frac{1}{n} \sum_i d_i$$

Η τυπική απόκλιση είναι

$$S_d = \sqrt{\frac{\sum d_i^2 - \frac{(\sum d_i)^2}{n}}{n-1}}$$

Διαστήματα εμπιστοσύνης

$$\bar{d} \pm t_{\alpha/2} \frac{S_d}{\sqrt{n}}$$

Υπολογισμός τυπικού σφάλματος

Ανεξάρτητα δείγματα

Μεσος ορος  $\mu_1 - \mu_2$

Τυπικό σφάλμα  $\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$

Μεγάλα δείγματα ( $n_1, n_2 > 30$ )  $\sigma_1^2, \sigma_2^2$  γνωστά

$$(\bar{X}_1 - \bar{X}_2) \pm Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

E  
Υπολογισμός εμπιστοσύνης σφάλματος