

**MEM-205 Περιγραφική Στατιστική**  
Τμήμα Μαθηματικών και Εφ. Μαθηματικών, Πανεπιστήμιο Κρήτης

Κώστας Σμαραγδάκης (kesmarag@gmail.com)

4η εβδομάδα (διάλεξη θεωρίας)

### Μέτρα κεντρικής τάσης

- ▶ Μέση τιμή  $\bar{x}$
- ▶ Διάμεσος  $M$
- ▶ Γεωμετρικός μέσος  $G$
- ▶ Επικρατέστερη τιμή  $M_0$

### Μέτρα μεταβλητότητας

- ▶ Εύρος  $R$
- ▶ Ενδοτεταρτημορικό εύρος  $IQR = Q_3 - Q_1$

- ▶ Μέση τιμή δείγματος:  $\bar{x}$
- ▶ Μέση τιμή πληθυσμού:  $\mu$

Έστω  $x_1, x_2, \dots, x_N$  παρατηρήσεις που αντιστοιχούν σε ένα τυχαίο δείγμα ενός πληθυσμού.

Έχουμε ορίσει ως μέση τιμή των παρατηρήσεων του δείγματος την ποσότητα:

$$\bar{x} = 1/N \sum_{n=1}^N x_n$$

Αυτή η μέση τιμή εκφράζει μόνο το δείγμα και όχι τον πληθυσμό, αν και για μεγάλο  $N$  προσεγγίζει την αντίστοιχη μέση τιμή  $\mu$  του πληθυσμού.

## Μέση Τιμή του Πληθυσμού vs Μέση Τιμή του Δείγματος

Ανεξάρτητα των τιμών του δείγματος ισχύει η ανισότικη σχέση

$$\sum_{n=1}^N (X_n - \bar{X})^2 \leq \sum_{n=1}^N (X_n - \mu)^2$$

με ισότητα μόνο αν  $\bar{X} = \mu$ .

Έστω  $f(x) = \sum_{n=1}^N (x_n - x)^2$

$$f'(x) = -2 \sum_{n=1}^N (x_n - x)$$

$$f'(x^*) = 0 \Rightarrow -2 \sum_{n=1}^N (x_n - x^*) = 0 \Leftrightarrow \sum_{n=1}^N x_n - \sum_{n=1}^N x^* = 0 \Leftrightarrow \sum_{n=1}^N x_n - n x^* = 0 \Leftrightarrow x^* = \frac{\sum_{n=1}^N x_n}{n} = \bar{X}$$

$$f''(x) = 2 > 0 \Leftrightarrow f''(x^*) > 0 \quad \text{άρα στο } x^* \text{ είναι το ολικό ελάχιστο της } f$$

## Παράδειγμα

Έστω το πείραμα της ρίψης ενός αμερόληπτου ζαριού. (Για κάποιες επαναλήψεις του πειράματος)

$$\mu = \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 2 + \frac{1}{6} \cdot 3 + \frac{1}{6} \cdot 4 + \frac{1}{6} \cdot 5 + \frac{1}{6} \cdot 6 = 3.5$$

Ρίχνουμε το ζάρι 3 φορές και λαμβάνουμε τα αποτελέσματα: 3,2,6

Έχουμε  $\bar{x} = 3.66$

$$\sum_{i=1}^3 (x_i - \bar{x})^2 = 8.66 < 8.75 = \sum_{i=1}^3 (x_i - \mu)^2$$

$(3 - 3.66)^2 + (2 - 3.66)^2 + (6 - 3.66)^2$        $(3 - 3.5)^2 + (2 - 3.5)^2 + (6 - 3.66)^2$

$$\sigma^2 = \frac{8.75}{3}$$

## Διασπορά ή Διακύμανση (Variance)

### Διασπορά πληθυσμού

$$\frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})^2 = \sigma^2$$

Ορίζεται ως η μέση τιμή του συνόλου τιμών

$$\{(x - \mu)^2\}$$

$$\sigma^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2$$

για κάθε παρατήρηση  $x$  του πληθυσμού. Η διασπορά του πληθυσμού συμβολίζεται με  $\sigma^2$ .

### Διασπορά στατιστικού δείγματος

$$s^2 = \frac{1}{N-1} \sum_{n=1}^N (x_n - \bar{x})^2$$

Μπορούμε να γράψουμε ισοδύναμα:

$$s^2 = \frac{\sum_{n=1}^N x_n^2 - \frac{(\sum_{n=1}^N x_n)^2}{N}}{N-1}$$

Όσο το  $N$  αυξάνεται έχουμε  $s^2 \rightarrow \sigma^2$ .

## Διασπορά στατιστικού δείγματος

$$s^2 = \frac{1}{N-1} \sum_{n=1}^N (x_n - \bar{x})^2$$

Γιατί διαιρούμε με  $N - 1$  και όχι απλά με  $N$ ;

Έστω ότι συμφιζώμε το  $\mu$ . Τότε για  $x_1, \dots, x_N$  παρατηρήσεις  $\sigma^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2$

Εάν δεν συμφιζώμε το  $\mu$ . Τότε για  $x_1, \dots, x_N$  προσεγγίζουμε το  $\mu$  με το  $\bar{x}$

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n \Leftrightarrow \sum_{n=1}^N x_n = N\bar{x} \Leftrightarrow x_N = N\bar{x} - \sum_{n=1}^{N-1} x_n$$

$$x_1, \dots, x_{N-1} \rightarrow x_N$$

Σε αυτή την περίπτωση είναι σαν έχουμε  $N-1$  παρατηρήσεις

### Διασπορά ομαδοποιημένων δεδομένων

$$s^2 = \frac{1}{N-1} \sum_{j=1}^K f_j (m_j - \bar{x})^2$$

Μπορούμε να γράψουμε ισοδύναμα:

$$s^2 = \frac{\sum_{j=1}^K m_j^2 f_j - \frac{(\sum_{j=1}^K m_j f_j)^2}{N}}{N-1}$$



## Διασπορά ή Διακύμανση (Variance)

$$s^2 = \frac{\sum_{j=1}^K m_j^2 f_j - \frac{(\sum_{j=1}^K m_j f_j)^2}{N}}{N - 1}$$

### Άσκηση - Διασπορά ομαδοποιημένων δεδομένων

	<b>f</b>	$m_j$	$m_j^2$	$m_j f_j$	$m_j^2 f_j$	
<b>[0,2)</b>	3	1	1	3	3	} → $s^2 = ?$
<b>[2,4)</b>	4	3	9	12	12	
<b>[4,6)</b>	5	5	25	25	125	
<b>[6,8)</b>	2	7	49	14	98	
<b>[8,10)</b>	4	9	81	36	324	
<b>[10,12)</b>	2	11	121	22	242	
<b>Total</b>	$20 = N$			$\sum m_j f_j$	$\sum m_j^2 f_j$	

## Τυπική Απόκλιση (Standard Deviation)

$$\begin{array}{l} [x] \text{ } \mu \\ [\sigma^2] \text{ } \mu^2 \quad [\sigma] \text{ } \mu \end{array}$$

Αποτελεί το πιο συχνά χρησιμοποιούμενο μέτρο μεταβλητότητας. Ορίζεται ως η τετραγωνική ρίζα της διασποράς.

- ▶ Τυπική απόκλιση πληθυσμού:

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

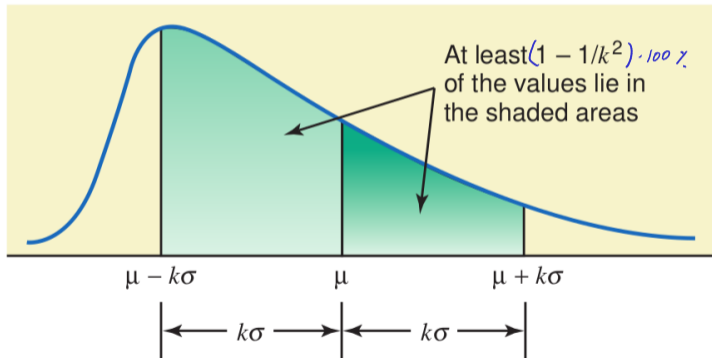
- ▶ Τυπική απόκλιση δείγματος:

$$s = \sqrt{s^2}$$

Η τυπική απόκλιση εκφράζεται στην ίδια μονάδα μέτρησης με τη μεταβλητή που αναφέρεται.

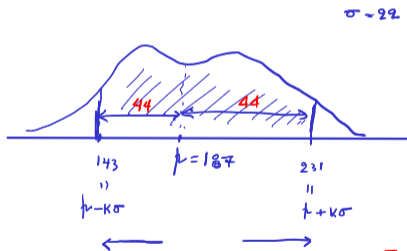
## Θεώρημα του Chebyshev

Για κάθε  $k > 1$ , τουλάχιστον  $(1 - 1/k^2)$  <sup>· 100%</sup> των παρατηρήσεων ανοίκουν στο διάστημα  $[\mu - k\sigma, \mu + k\sigma]$



## Άσκηση

Η μέση συστολική αρτηριακή πίεση 4000 γυναικών που υποβλήθηκαν σε εξέταση για υψηλή πίεση αίματος βρέθηκε να είναι 187 mm Hg με τυπική απόκλιση 22. Χρησιμοποιώντας το Θεώρημα του Chebyshev βρείτε το ελάχιστο ποσοστό των γυναικών αυτής της ομάδας με συστολική αρτηριακή πίεση μεταξύ 143 και 231 mm Hg.



$$\cancel{\mu - k\sigma} - \cancel{\mu + k\sigma} = 2k\sigma = \overset{22}{231 - 143} = 88 \Rightarrow \boxed{k=2}$$

Τουλάχιστον το  $(1 - \frac{1}{2^2}) \cdot 100\% = 75\%$  των

Παρατηρήσεων  $\in [143, 231]$ .

- ▶ Είναι το πηλίκο της τυπικής απόκλισης δια της μέσης τιμής. Συμβολίζεται ως CV:

$$CV = \frac{s}{\bar{x}}, \bar{x} > 0$$

- ▶ Είναι χρήσιμος για τη σύγκριση της ομοιογένειας δυο συσχετισμένων μεταβλητών με διαφορετικές μονάδες μέτρησης ή στο να συγκρίνουμε την ομοιογένεια μεταβλητών με ίδιες μονάδες μέτρησης αλλά με διαφορετικές μέσες τιμές.
- ▶ Επίσης χρησιμοποιείται για το χαρακτηρισμό ενός δείγματος ως ομοιογενές (CV  $\geq$  0.1) ή ~~α~~ομοιογενές (CV < 0.1) .

### Παράδειγμα

Έστω δείγματα με τις ημερήσιες μετρήσεις θερμοκρασίας 2 πολέων στη διάρκεια ενός έτους. Για την πόλη Α η μέση θερμοκρασία ήταν 20 βαθμούς °C και η τυπική απόκλιση 2, ενώ για την Β η μέση θερμοκρασία ήταν 15 βαθμούς °C και η τυπική απόκλιση 1.8

$$CV_A = \frac{s_A}{\bar{x}_A} = \frac{2}{20} = 0.1 < CV_B = \frac{s_B}{\bar{x}_B} = \frac{1.8}{15} = 0.12$$

### Παράδειγμα

Σε δυο γραπτές δοκιμασίες οι μαθητές μιας τάξης είχαν επιδόσεις που περιγράφονται παρακάτω:

δοκιμασία Α (κλίμακα 0-20): μέση τιμή 14, τυπική απόκλιση 1.4

δοκιμασία Β (κλίμακα 0-100): μέση τιμή 70, τυπική απόκλιση 3.5

$$CV_A = \frac{s_A}{\bar{x}_A} = \frac{1.4}{14} = 0.1 > CV_B = \frac{3.5}{70} = 0.05$$