

Χαρακτηριστικό, Τυχαίο Δείγμα και παρατηρήσεις

X : χαρακτηριστικό (δέντρα)

Έχω δείγμα μεγέθους $n (= 100)$

$X_1, X_2, X_3, \dots, X_{100}$

τυχαίες μεταβλητές
(κάθε φορά που κάνω δει-
ματοληψία έχω αλλη-
τη τιμή)

Έχω πάρει ένα συγκεκριμένο
δείγμα

$x_1 = 0$ $x_2 = 5$ $x_3 = 1$ $x_4 = 0$ $x_5 = 0$ ----- $x_{100} = 1$

(έχω διακριτό χαρακτηριστικό)

	0	1	2	3	4	5
συχνότητα	46	32	8	11	2	1

↓
από τις x_1, \dots, x_{100} οι 46 ισούται με 0

Επιλογή «μοντέλου» για την καλύτερη περιγραφή του χαρακτηριστικού

$X \sim \text{Poisson}(\theta)$
 ↓ αποδοθεί → μοντέλο

παράμετρος: αν τυχερά προσδιορίσω "ξέρω" τα πάντα για το χαρακτηριστικό

$P(X=x) = e^{-\theta} \frac{\theta^x}{x!}$ $x=0,1,2,\dots$
 ↓ συγκεκριμένη τιμή
 ↓ χαρακτηριστικό
 ↓ πιθανότητα

$$x! = 1 \cdot 2 \cdot 3 \cdot \dots \cdot x$$

$$0! = 1$$

$X =$ γεγονότα x παραγοντ στον χώρο ή στον χρόνο

$$\sum_{x=0}^{\infty} P(X=x) = \sum_{x=0}^{\infty} e^{-\theta} \frac{\theta^x}{x!} = 1 \quad \theta > 0$$

Εκτιμητές των παραμέτρων- εκτίμηση από το συγκεκριμένο δείγμα

Το θ είναι η μέση τιμή: (πληθυσμός)

$$E(X) = \sum_{x=0}^{\infty} x P(X=x) \rightarrow \text{από το θεωρητικό μοντέλο}$$

↓ expected value $E(X) = \theta$

$$\bar{X} = \sum_{i=1}^n X_i \frac{1}{n} \rightarrow \text{από τις παρατηρήσεις του δείγματος}$$

↓ δειγματικός μέσος

$$\theta = \bar{X} \quad (\text{εκτιμητής: συνάρτηση των παρατηρήσεων})$$

↓
εκτίμω την άγνωστη παράμετρο από την αντίστοιχη δειγματική της τιμή.

κάθε φορά \bar{x} είναι διαφορετικό (έχει άλλη τιμή για κάθε δείγμα)

και έτσι παίρνω μια εκτίμηση για το θ . Μπορώ να ^{βάλω} την συγκεκριμένη τιμή στο μοντέλο και να το προσδιορίσω πάλι.

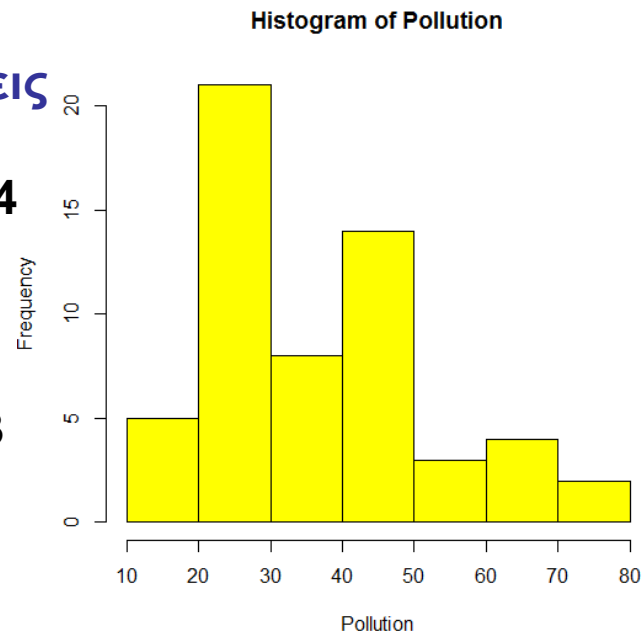
Δείγμα από τη συγκέντρωση ενός συγκεκριμένου ρύπου (σε mgr/cm³) σε δείγματα αέρος από 57 πόλεις

68 63 42 27 30 36 28 32 79 27 22 23 24 25 24

65 43 25 74 51 36 42 28 31 28 25 45 12 57 51

12 32 49 38 42 27 31 50 38 21 16 24 69 47 23

22 43 27 49 48 23 12 19 46 30 49 49



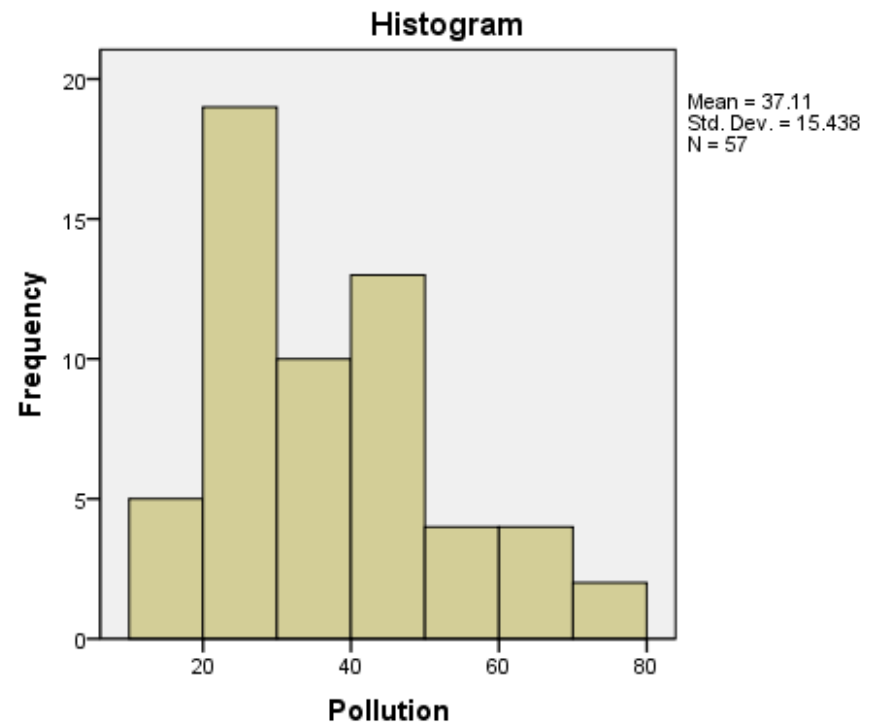
$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n} = \frac{68 + 63 + 42 + \dots + 30 + 49 + 49}{57} = \frac{2099}{57} = 36.82 \text{ mgr/cm}^3$$

$$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n - 1} = \frac{(68 - 36.82)^2 + (63 - 36.82)^2 + \dots + (49 - 36.82)^2}{56}$$

$$= \frac{971.90 + 685.15 + \dots + 148.24}{56} = \frac{14366.25}{56} = 256.54 = (16.02)^2$$

$$S = 16.02 \text{ mgr/cm}^3$$

Όρια κλάσης L_i-U_i	Κεντρική τιμή x_i	Συχνότητα f_i	Αθροιστική Συχνότητα F_i
[10-20)	15	5	5
[20-30)	25	19	24
[30-40)	35	10	34
[40-50)	45	13	47
[50-60)	55	4	51
[60-70)	65	4	55
[70-80]	75	2	57
		57	



Cases weighted by Frequency

$$\bar{X} = \frac{\sum_{i=1}^k x_i f_i}{\sum_{i=1}^k f_i} = \frac{15 \times 5 + 25 \times 19 + \dots + 65 \times 4 + 75 \times 2}{5 + 19 + \dots + 4 + 2} = \frac{2115}{57} = 37.11 \text{ mgr/cm}^3$$

$$S^2 = \frac{\sum_{i=1}^k (X_i - \bar{X})^2 f_i}{n - 1} = \frac{(15 - 37.11)^2 \times 5 + (25 - 37.11)^2 \times 19 + \dots + (75 - 37.11)^2 \times 2}{56} = \frac{2443.21 + 2784.21 + \dots + 2872.02}{56} = \frac{13347.37}{56} = 238.34 = (15.43)^2$$

$$S = 15.43 \text{ mgr/cm}^3$$

Statistics

Pollution

N	Valid	57
	Missing	0
Mean		37.11
Mode		25
Std. Deviation		15.438
Variance		238.346
Percentiles	25	24.79 ^a
	50	34.66
	75	47.65

a. Percentiles are calculated from grouped data.

Pollution

Όρια κλάσης $L_i - U_i$	Frequency	Percent	Αθροιστική Συχνότητα F_i	Cumulative Percent
[10-20)	15	5	8.8	8.8
[20-30)	25	19	33.3	42.1
[30-40)	35	10	17.5	59.6
[40-50)	45	13	22.8	82.5
[50-60)	55	4	7.0	89.5
[60-70)	65	4	7.0	96.5
[70-80)	75	2	3.5	100.0
Total	57	100.0		

$$\delta = Q_2 = L_3 + \frac{\frac{n}{2} - F_2}{f_3} c = 30 + \frac{28.5 - 24}{10} 10$$

$$= 30 + 4.5 = 34.5 \text{ mgr/cm}^3$$

Διάμεσος (Median) δ

• ενός πεπερασμένου συνόλου τιμών μιας μεταβλητής ή παρατηρήσεων είναι εκείνη η τιμή που χωρίζει το σύνολο σε δύο ίσα μέρη
 Δηλ. 50% παρατηρήσεων $\geq \delta$
 ————— $\leq \delta$

• Διατάσσουμε τις παρατηρήσεις κατά αύξουσα σειρά μεγέθους. Έστω το πλήθος των παρατηρήσεων n είναι περιττό τότε η διάμεσος είναι η μεσαία τιμή. Έστω το n είναι άρτιο $\delta =$ κριτήριο των δύο μεσίων παρατηρήσεων

• Σε ομαδοποιημένες παρατηρήσεις η διάμεσος δίνεται από τον τύπο

$$\delta = L_i + \frac{\frac{n}{2} - F_{i-1}}{f_i} \cdot c$$

L_i : το κάτω αλγεθινό όριο της κλάσης που βρίσκεται η διάμεσος
 n : πλήθος παρατηρήσεων
 f_i : η διάμεσος

$F_{i-1} = \sum_{j=1}^{i-1} f_j$ το άθροισμα των συχνοτήτων μέχρι των προηγούμενων κλάσεων

c : εύρος της κλάσης

με i τέτοιο ώστε $F_{i-1} < n/2 \leq F_i$

$$\frac{n}{2} = 28.5 \quad F_2 = 24 < 28.5 \leq 34 = F_3$$

Όρια κλάσης $L_i - U_i$	Κεντρική τιμή x_i	Συχνότητα f_i	Αθροιστική Συχνότητα F_i
[10-20)	15	5	5
[20-30)	25	19	24
[30-40)	35	10	34
[40-50)	45	13	47
[50-60)	55	4	51
[60-70)	65	4	55
[70-80]	75	2	57
		57	

• Για ομαδοποιημένες παρατηρήσεις

$$p_k = L_i + \left(\frac{\frac{kn}{100} - F_{i-1}}{f_i} \right) c$$

με i τέτοιο ώστε $F_{i-1} < k \cdot n / 100 \leq F_i$

• 1^ο τεταρτημόριο

$$Q_1 = L_i + \left(\frac{\frac{n}{4} - F_{i-1}}{f_i} \right) c$$

• 2^ο τεταρτημόριο (Μέσος)

$$\delta = L_i + \left(\frac{\frac{n}{2} - F_{i-1}}{f_i} \right) c$$

• 3^ο τεταρτημόριο

$$Q_3 = L_i + \left(\frac{\frac{3n}{4} - F_{i-1}}{f_i} \right) c$$

$$\frac{n}{4} = 14.25$$

$$F_1 = 5 < 14.25 \leq 24 = F_2$$

$$\frac{3n}{4} = 42.75$$

$$F_3 = 34 < 42.75 \leq 47 = F_4$$

$$Q_1 = L_2 + \frac{\frac{n}{4} - F_1}{f_2} c = 20 + \frac{14.25 - 5}{19} 10 = 20 + 4.49 = 24.49 \text{ mgr/cm}^3$$

$$Q_3 = L_4 + \frac{\frac{3n}{4} - F_3}{f_4} c = 40 + \frac{42.75 - 34}{13} 10 = 40 + 6.73 = 46.73 \text{ mgr/cm}^3$$

• Σε ομαδοποιημένες παρατηρήσεις η κορυφή δίνεται από τον τύπο

$$M_0 = L_i + \frac{\Delta_1}{\Delta_1 + \Delta_2} \cdot c$$

$$\Delta_1 = f_i - f_{i-1}$$

$$\Delta_2 = f_i - f_{i+1}$$

$$\max f_i = f_2 = 19$$

$$\Delta_1 = f_2 - f_1 = 19 - 5 = 14 \quad \Delta_2 = f_2 - f_3 = 19 - 10 = 9$$

$$M = L_2 + \frac{\Delta_1}{\Delta_1 + \Delta_2} c = 20 + \frac{14}{14 + 9} 10 = 20 + 6.09 = 26.09 \text{ mgr/cm}^3$$

με i τέτοιο ώστε $f_i > f_j$ για κάθε $j \neq i$

ΕΦΑΡΜΟΣΜΕΝΑ ΜΑΘΗΜΑΤΙΚΑ

Περιγραφική Στατιστική

Εναλλακτικός τύπος υπολογισμού της διασποράς

Από τον δειγματικό μέσο \bar{X} προκύπτει ότι $\bar{X} = \frac{\sum_{i=1}^n X_i}{n} \Rightarrow \sum_{i=1}^n X_i = n\bar{X}$ (1)

Αναπτύσσοντας το τετράγωνο, στον τύπο της δειγματικής διασποράς παίρνουμε:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i^2 - 2\bar{X}X_i + \bar{X}^2) \stackrel{*}{=}$$

$$\begin{aligned} \text{Γαυτό το άθροισμα είναι το εἶναι το εἶναι:} \quad & \sum_{i=1}^n (X_i^2 - 2\bar{X}X_i + \bar{X}^2) = \\ & = (X_1^2 - 2\bar{X}X_1 + \bar{X}^2) + (X_2^2 - 2\bar{X}X_2 + \bar{X}^2) + \dots + (X_n^2 - 2\bar{X}X_n + \bar{X}^2) \end{aligned}$$

$$\stackrel{*}{=} \frac{1}{n-1} \left\{ \sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + n\bar{X}^2 \right\} \stackrel{(1)}{=} \frac{1}{n-1} \left\{ \sum_{i=1}^n X_i^2 - 2\bar{X}n\bar{X} + n\bar{X}^2 \right\} \Rightarrow$$

$$S^2 = \frac{1}{n-1} \left\{ \sum_{i=1}^n X_i^2 - n\bar{X}^2 \right\} \stackrel{(1)}{=} \frac{1}{n-1} \left\{ \sum_{i=1}^n X_i^2 - \frac{\left(\sum_{i=1}^n X_i\right)^2}{n} \right\}$$

το S ονομάζεται δειγματική τυπική απόκλιση και το μετράμε στις ίδιες μονάδες με το χαρακτηριστικό

Στο διάστημα $(\bar{X} - 3S, \bar{X} + 3S)$ περιλαμβάνονται, σχεδόν βέβαια, όλες οι παρατηρήσεις του δείγματος

Εναλλακτικός τύπος υπολογισμού της διασποράς

κεντρικά τιμή ομαδοποιημένες παρατηρήσεις

$$S^2 = \frac{1}{n-1} \sum_{i=1}^k (x_i - \bar{x})^2 f_i \quad \text{διασπορά} \quad k: \text{πλήθος κλάσεων}$$

x_i : αντιπρόσωπος για τις f_i το πλήθος παρατηρήσεων που i κλάση

$$= \frac{1}{n-1} \sum_{i=1}^k (x_i^2 - 2\bar{x}x_i + \bar{x}^2) f_i$$

$$= \frac{1}{n-1} \left\{ \sum_{i=1}^k x_i^2 f_i - 2\bar{x} \sum_{i=1}^k x_i f_i + \bar{x}^2 \sum_{i=1}^k f_i \right\} \quad \bar{x} = \frac{1}{n} \sum_{i=1}^k x_i f_i$$

$$= \frac{1}{n-1} \left\{ \sum_{i=1}^k x_i^2 f_i - 2n\bar{x}^2 + n\bar{x}^2 \right\} \quad n\bar{x} = \sum_{i=1}^k x_i f_i$$

$$= \frac{1}{n-1} \left\{ \sum_{i=1}^k x_i^2 f_i - n\bar{x}^2 \right\} = \frac{1}{n-1} \left\{ \sum_{i=1}^k x_i^2 f_i - n \left(\frac{\sum x_i f_i}{n} \right)^2 \right\}$$

$$= \frac{1}{n-1} \left\{ \sum_{i=1}^k x_i^2 f_i - \frac{(\sum x_i f_i)^2}{n} \right\}$$

Γραμμικός Μετασχηματισμός των παρατηρήσεων

Έστω $y_i = \alpha x_i + \beta$ $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ $S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

\hookrightarrow αλλαγή θέσης \hookrightarrow κλίμακας

$$\bar{y} = \frac{\sum_{i=1}^n y_i}{n} = \frac{\sum_{i=1}^n (\alpha x_i + \beta)}{n} = \frac{\sum_{i=1}^n \alpha x_i}{n} + \frac{\sum_{i=1}^n \beta}{n} = \alpha \frac{\sum_{i=1}^n x_i}{n} + \frac{n\beta}{n}$$
$$= \alpha \bar{x} + \beta$$

$$S_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n-1} \sum_{i=1}^n (\alpha x_i + \beta - (\alpha \bar{x} + \beta))^2$$

$$= \frac{1}{n-1} \sum_{i=1}^n (\alpha x_i - \alpha \bar{x})^2 = \frac{1}{n-1} \sum_{i=1}^n [\alpha (x_i - \bar{x})]^2$$

$$= \alpha^2 \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \Rightarrow S_y^2 = \alpha^2 \cdot S_x^2 \Rightarrow S_y = |\alpha| S_x$$

Τα ποσοτικά σημεία και η κορυφή μετασχηματίζονται, «όπως αριβώς» και οι παρατηρήσεις προσοχή για $\alpha < 0$