

**סטטיסטיקה תיאורית + לוחות שכיחות בדידים/רציפים**

בגדול מקצוע הסטטיסטיקה נחלק ל-2 תחומים עיקריים- סטטיסטיקה תיאורית וסטטיסטיקה היסקית;

בסטטיסטיקה היסקית משערים השערות, משווים בין קבוצות באוכלוסיה ועוד, אך גם מסטטיסטיקה תיאורית ניתן ללמוד הרבה על האוכלוסיה- ניתן ללמוד על המבנה שלה, על הנטיות המרכזיות שלה, על הפיזור שלה וכו';

בסטטיסטיקה תיאורית עוסקים מכמה זוויות בד"כ:

1. ניתוח נתונים גולמיים.
2. ניתוח טבלאות שכיחות בדידות.
3. ניתוח טבלאות שכיחות רציפות.

הניתוחים הללו מתבטאים הן במדדים (מרכז ופיזור) והן בגרפים מתאימים (גרף עוגה/גרף מקלות/היסטוגרמה);

**נתונים גולמיים** הינם פשוט רצפים של מספרים שמייצגים משהו; הם יכולים להיות טבעיים, שלמים, רציונליים, מרוכבים ועוד;

**טבלת שכיחות** הינה למעשה דרך לקחת קבוצה גדולה של תצפיות, ולרכז אותה בטבלה, כאשר יש לנו ערכים אחד (לרוב יקראו x), ושכיחויות מנגד (frequencies).

כאשר נצפה למגוון לא גבוה של תצפיות (לדוגמא מס' ילדים במשפחה, מס' חשבונות בנק, מס' מכוניות)- לרוב ירוכזו נתונים אלה בטבלאות שכיחות בדידות, וכאשר נצפה למגוון גבוה יותר של תצפיות (לדוגמא משכורת, גובה, משקל)- ירוכזו ע"פ רוב הנתונים הללו בטבלאות שכיחות רציפות.

**מדדי המרכז** המקובלים ע"פ רוב הם ממוצע (לרוב מסומן כ-  $\bar{x}$ ) שכיח (נקרא Mode ולרוב מסומן כ-Mo) וחציון (נקרא Median ולרוב מסומן כ-Me, Med או Md);

**ממוצע**- הוא סכום הערכים ברצף חלקי מספר הערכים; בכתיבה-  $\bar{x} = \frac{\sum x_i}{n}$ .

**שכיח**- הוא הערך שחוזר על עצמו הכי הרבה פעמים בהתפלגות; לעיתים מצביעים עליו ולעיתים הוא מתקבל ע"י נוסחה.

**חציון**- הוא הערך בהתפלגות, שמימינו ומשמאלו, ישנה בדיוק אותה כמות של תצפיות, כלומר ניתן גם להתייחס אליו לערך שלמעשה מחלק את האוכלוסיה ל-2 חלקים שווים; גם את החציון ניתן לחשב ע"פ נוסחה (נראה בהמשך).

מדדי הפיזור המקובלים ע"פ רוב הם סטיית תקן (לרוב תסומן ב-  $s$  או ב-  $\sigma$ ), שונות (סטיית תקן בריבוע; תסומן לרוב ב-  $s^2$  או ב-  $\sigma^2$  או ב-  $Var$  - מן המילה האנגלית Variance), תחום בין-רבעוני (תסומן-  $IQR$ ) ומקדם ההשתנות (לרוב יסומן כ-  $CV$ ).

נרחיב על כל המדדים בהמשך;

**טבלאות שכיחות בדידות:**

כאמור בטבלאות שכיחות נמצא ערכים השייכים להתפלגות מסוימת, ולצדם- השכיחויות של אותם ערכים; לרוב לא יהיו המון ערכים כאלה; ישנן כמה וכמה דרכים לחשב טבלאות שכיחות בדידות; אנחנו נציג דרך אחת, לא קשה במיוחד, שיעושה את העבודה'.

נניח שנבדקה קבוצה של  $n=200$  משפחות, ונתונה לנו התפלגות של מספר הילדים שם; להלן ההתפלגות שנתקבלה-

$f(x)$	$x$
25	0
30	1
40	2
65	3
28	4
8	5
4	6

נציג כעת דרך לחשב את המדדים שהזכרנו; יש להוסיף 3 טורים-

$x^2 \cdot f(x)$	$F(x)$	$x \cdot f(x)$	$f(x)$	$x$
0	25	0	25	0
30	55	30	30	1
160	95	80	40	2
585	160	195	65	3
448	188	112	28	4
200	196	40	8	5
144	200	24	4	6

ממוצע-  $\bar{x} = \frac{\sum x \cdot f}{n} = \frac{0 + 30 + \dots + 40 + 24}{200} = \frac{481}{200} = 2.405$  - בממוצע ישנם כ-2.4 ילדים למשפחה.

שכיח- הערך שחוזר על עצמו הכי הרבה (כלומר שה-  $f(x)$  עבורם הוא המקסימלי) -  $Mo=3$ .

חציון- הערך שמימנו ומשמאלו ישנה אותה כמות של תצפיות; כיוון שלפנינו כמות זוגית של תצפיות, החציון יהיה הממוצע של שתי התצפיות המרכזיות בהתפלגות, היינו ממוצע התצפיות במקום ה-

100 וה-101-

$$Med = \frac{x_{n/2} + x_{n/2+1}}{2} = \frac{x_{100} + x_{101}}{2} = \frac{3+3}{2} = 3$$

סטיית תקן - מדד המייצג את הפיזור באוכלוסיה, כלומר עד כמה האוכלוסיה הומוגנית או הטרוגנית; סטיית התקן היא מספר טהור, 0 או יותר, שככל שהערך גבוה יותר כך סטיית התקן גבוהה יותר; כשלעצמו אין למספר חשיבות מיוחדת, אך כאשר משווים בין אוכלוסיות, אפשר ללמוד עליהן המון מתוך ההשוואה בין סטיות התקן שלהן.

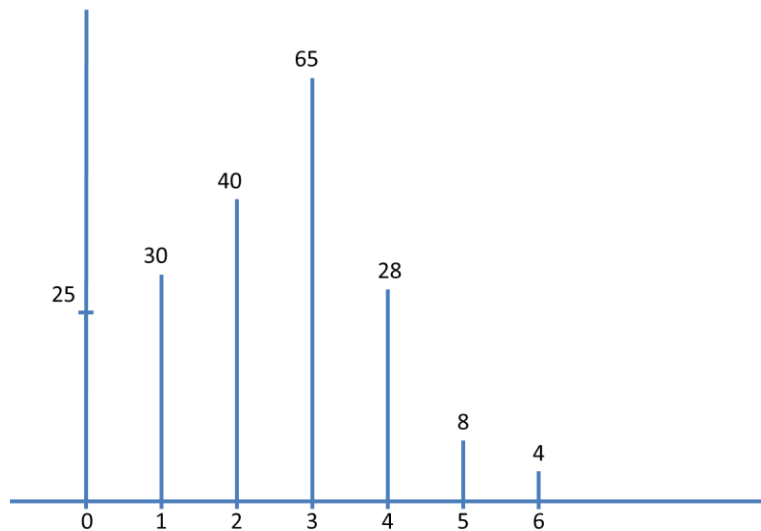
ישנן מגוון שיטות לחשב את סטיית התקן; אנחנו נחשב את הסטייה באמצעות שימוש בנוסחת עבודה, שמקלה מאד על החישוב (יחסית לנוסחה ע"פ ההגדרה); נשים לב שהנוסחה מותאמת לטבלאות שכיחות, וכדי שתתאים לנתונים גולמיים יש לבצע שינויים קלים;

$$\sigma = \sqrt{\frac{\sum x_i^2 f(x)}{n} - \bar{x}^2} = \sqrt{\frac{1,567}{200} - 2.405^2} = 1.4321$$

$$\sigma^2 = \frac{1,567}{200} - 2.405^2 = 2.051$$

סטיית התקן בריבוע - 2.051

הצגה גרפית - שיטת ההצגה המקובלת במקרה של משתנה כמותי בדיד היא גרף מקלות (Bar Chart) -



**טבלאות שכיחות רציפות:**

גם כאן ישנם ערכים השייכים להתפלגות מסוימת ולצדם השכיחויות שלהם; לרוב מדובר בערכים רבים ומגוונים, ולכן יש צורך לקבץ אותם לקבוצות; לעיתים הקבוצות מקובצות עם גבולות מדומים ולכן יש להפוך את הקבוצות לגבולות אמיתיים; ישנן שיטות שונות למעבר לגבולות אמיתיים; אנחנו נתמקד בטבלאות כבר בגבולות האמיתיים שלהם;

נניח שנבדקה אותה קבוצה של  $n=200$  משפחות, ונתונה לנו התפלגות השכר שלהם (באלפים); להלן ההתפלגות שנתקבלה-

f(x)	x
15	0-3
18	3-5
55	5-8
50	8-12
37	12-16
20	16-20
5	20-25

כאן יש לבנות טבלה מעט יותר גדולה;

$x^2 \cdot f(x)$	צפיפות	רוחב מחלקה	F(x)	$x \cdot f(x)$	אמצע קטע	f(x)	x
33.75	5	3	15	23	1.5	15	0-3
288	9	2	33	72	4	18	3-5
2,323.75	18.33	3	88	357.5	6.5	55	5-8
5,000	12.50	4	138	500	10	50	8-12
7,252	9.25	4	175	518	14	37	12-16
6,480	5	4	195	360	18	20	16-20
2,531.25	1	5	200	113	22.5	5	20-25

ממוצע- מחושב באותה דרך כמו בטבלאות בדידות, רק עם אמצעי המחלקות-

$$\bar{x} = \frac{\sum \hat{x} \cdot f}{n} = \frac{1,942.5}{200} = 9.7125$$

שכיח- אמצע הקטע של המחלקה עם הצפיפות הגבוהה ביותר-  $Mo = 6.5$ .

חציון- ישנן כמה גרסאות לנוסחת החציון; כאן תוצג נוסחה אחת, יחסית נוחה לשימוש ולעין; את המחלקה החציונית נבחר לפי מציאת הערך  $n/2$  (במקרה שלנו-100) ומציאת המחלקה הרלבנטית בה הוא עובר (במקרה שלנו- 8-12); אחרי שמצאנו את המחלקה הרלבנטית- מחלקת החציון- **נשים לב שלמעשה החציון לא יכול להיות מחוץ לגבולות המחלקה הזו- מחלקת החציון.**

$$Med = L_0 + \frac{n/2 - F_{x(m-1)}}{f_{x(m)}} \cdot (L_1 - L_0) = 8 + \frac{100 - 88}{50} \cdot (12 - 8) = 8.96$$

ס.תקן- מחושב באותה דרך כמו בטבלאות בדידות, רק עם אמצעי המחלקות-

$$\sigma = \sqrt{\frac{\sum x_i^2 f(x)}{n} - \bar{x}^2} = \sqrt{\frac{23,908.75}{200} - 9.7125^2} = 5.021$$

שונות- סטיית התקן בריבוע-  $\sigma^2 = \frac{23,908.75}{200} - 9.7125^2 = 25.211$

תחום בין רבעוני- ההפרש בין האחוזון ה-75 לאחוזון ה-25 בהתפלגות, כלומר בין הרבעון השלישי לרבעון הראשון; נחשב כ"א מהם בנפרד ולבסוף נחסיר זה מזה-

$$Q_1 = L_0 + \frac{n/4 - F_{q_1-1}}{f_{q_1}} \cdot (L_1 - L_0) = 5 + \frac{50 - 33}{55} \cdot (8 - 5) = 5.927 \text{ רבעון 1-}$$

$$Q_3 = L_0 + \frac{3n/4 - F_{q_3-1}}{f_{q_3}} \cdot (L_1 - L_0) = 12 + \frac{150 - 138}{37} \cdot (16 - 12) = 13.297 \text{ רבעון 3-}$$

$$IQR = Q_3 - Q_1 = 13.297 - 5.927 = 7.369 \text{ ההפרש-}$$

$$CV = \frac{s}{\bar{x}} = \frac{5.021}{9.7125} = 0.5169 \text{ נקרא גם CV-}$$

מקדם ההשתנות-מדד פיזור שמשמשים בו לעיתים; נקרא גם CV-

תצוגה גרפית- עבור משתנה רציף, ההצגה הגרפית המתאימה היא היסטוגרמה, אשר לוקחת למעשה בחשבון גם את רוחב המחלקה (ציר ה-X) וגם את הצפיפות שלה (ציר ה-Y)-

