

ПРАКТИЧНА РОБОТА № 1-3 ГРУПУВАННЯ ДАНИХ

1. Мета роботи та завдання:

Метою роботи є вивчення принципів обробки статистичних даних засобами Microsoft Excel і способів подання результатів у статистичних таблицях і на графіках.

Завданнями роботи є:

- Закріплення навичок формування та заповнення таблиць з використанням вбудованих формул;
- Оволодіння методикою групування даних і обчислення основних показників варіації;
- Надбання навичок побудови статистичних графіків.

2. Теоретичні відомості:

Для вивчення форми емпіричного розподілу проводять групування даних. Результати групування представляють у вигляді таблиць і графіків. Аналітичне групування даних призначене для аналізу кореляційного взаємозв'язку. Таке групування полягає в розбитті діапазону можливих значень на інтервали і підрахунку підсумків по кожній групі.

Можна використати такі правила групування:

- Межі інтервалів частіше за все круглі числа (наприклад, 10-20, 70-80, 120-150 і т.і.);
- Інтервали повинні мати однакову ширину (наприклад: <100, 100-120, 120-140, 140-160, > 160). При цьому ширина першого і останнього інтервалів приймається рівною іншим;
- Рекомендована кількість інтервалів частіше за все на порядок менше обсягу вибірки;
- Бажано уникати появи порожніх і малочисельних інтервалів.

Для кожного з інтервалів необхідно обчислити наступні показники:

- n_i - частота (кількість елементів вибірки, що потрапляють в даний інтервал);
- $n_i\%$ - відносна частота (частка числа елементів в даному інтервалі від обсягу вибірки);
- $k_i\%(w_i\%)$ - накопичена частота (сума частот поточного інтервалу і всіх попередніх).

Дослідження варіації в статистиці має велике значення, допомагає пізнати сутність досліджуваного явища. Вимірювання варіації, з'ясування її причини, виявлення впливу окремих факторів дає важливу інформацію, наприклад, про тривалість технологічного процесу, витрати на закупівлю обладнання, про фінансове становище підприємства і т.і.

Варіація - це відмінність у значеннях якої-небудь ознаки у різних одиниць даної сукупності в один і той же період або момент часу. Варіація виникає в результаті того, що індивідуальні значення ознаки складаються під сукупним впливом різноманітних факторів (умов), які по-різному поєднуються в кожному окремому випадку.

До основних показників варіації відносяться: розмах варіації, обсяг вибірки, медіана, мода, середнє, дисперсія і т.і. (Див. Табл.1.1).

Будь-яка випадкова величина має *функцію розподілу* - залежність густини ймовірності від значення випадкової величини. Для *нормального розподілу* (розподілу Гаусса) функція розподілу має такий вигляд:

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

де μ – середнє значення (математичне очікування), σ - стандартне відхилення. *Стандартним нормальним розподілом* називається нормальний розподіл з математичним очікуванням 0 і стандартним відхиленням 1.

Таблиця 1.1
Показники варіації

Назва	Позначення	Назва у зводній таблиці	Метод обчислення	Формула Excel
1	2	3	4	5
Розмах варіації	R	Інтервал	Різниця максимального і мінімального значень	МАКС(інтервал)-МІН(інтервал)
О'бєм вибірки	n	Рахунок	Кількість статистичних одиниць	РАХУНОК(інтервал)
Медіана	Me	Медіана	Центральне значення відсортованої вибірки	МЕДІАНА(інтервал)
Мода	Mo	Мода	Найбільш часто зустрічається значення	МОДА(інтервал)
Середнє	\bar{x}	Середнє	Середнє арифметичне	СРЗНАЧ(інтервал)
Середнє лінійне відхилення	δ, d	-	Середній модуль відхилення від середнього значення	СРВІДХ(інтервал)
Дисперсія	σ^2	Дисперсія	Середній квадрат відхилення від середнього значення	ДИСП(інтервал)
Середнє квадратичне відхилення	σ	-	середнє квадратичне відхилення від середнього значення	СТАНДВІДХ(інтервал)
Середнє квадратичне відхилення(незміщена оцінка)	σ	Стандартне відхилення	середнє квадратичне відхилення від середнього значення з поправкою на обсяг вибірки	СТАНДВІДХ(інтервал)- незміщена оцінка

Закінчення таблиці 1.1

1	2	3	4	5
Коефіцієнт осциляції	V_R	-	$V_R = \frac{R}{\bar{x}} * 100$	-
Лінійний коефіцієнт варіації	V_d	-	$V_d = \frac{d}{\bar{x}} * 100$	-
Коефіцієнт варіації	V_σ	-	$V_\sigma = \frac{\sigma}{\bar{x}} * 100$	-
Коефіцієнти асиметрії	A	асиметричність	$A_s = \frac{\sum (x_i - \bar{x})^3}{n\sigma^3}$	-
Коефіцієнт ексцесу	E	ексцес	$E_x = \frac{\sum (x_i - \bar{x})^4}{n\sigma^4} - 3$	-

Графік функції густини розподілу ймовірностей і інтегральної функції розподілу представлені на Рис. 1.1. і 1.2

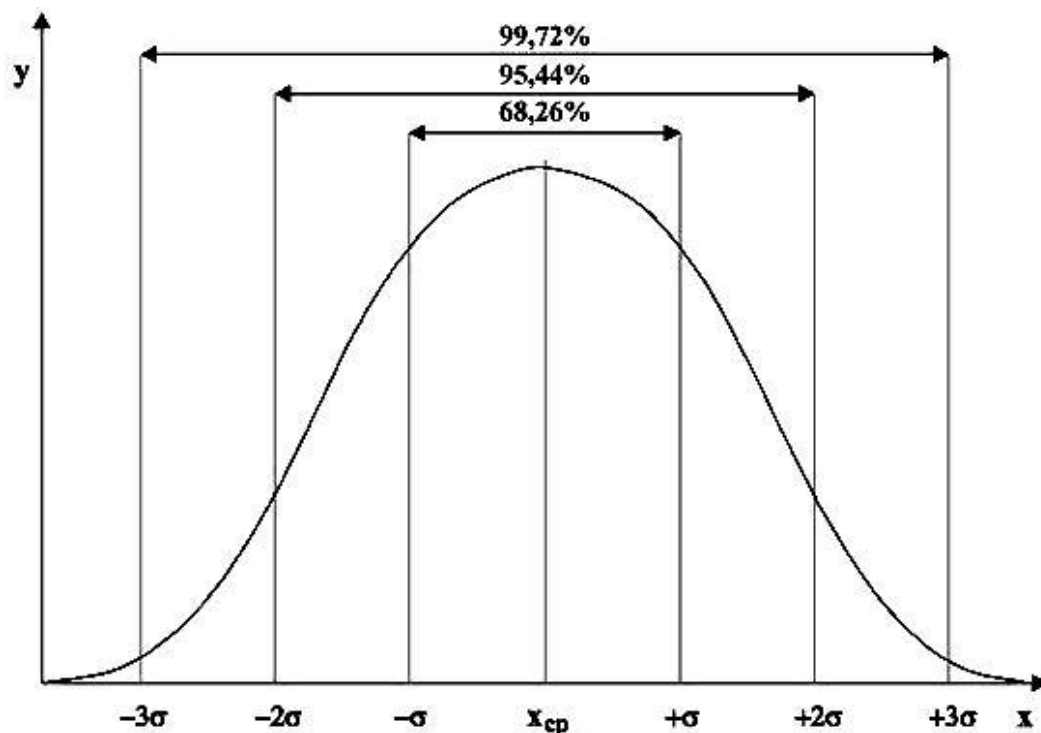


Рис.1.1 Графік густини розподілу.

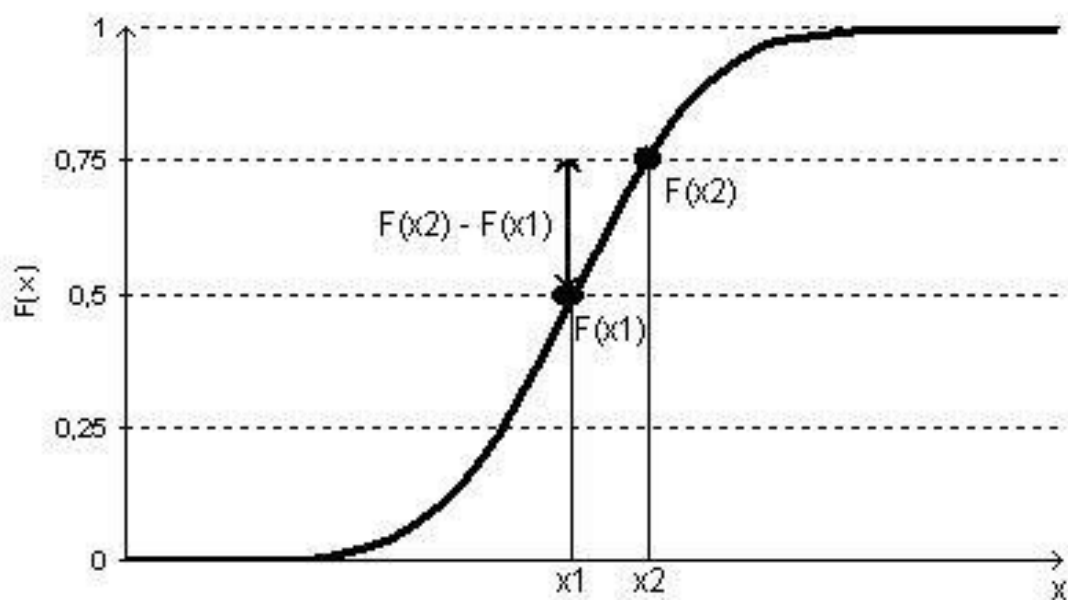


Рис.1.2 Функція розподілу.

Безперервна випадкова величина має рівномірний розподіл на відрізку $[a, b]$, якщо на цьому відрізку густина розподілу випадкової величини постійна, а поза ним дорівнює нулю.

$$f(x) = \begin{cases} 0, & x < a \\ C, & a \leq x \leq b, \\ 0, & x > b \end{cases} \quad \text{де } C = \frac{1}{b-a}.$$

Постійна величина C може бути визначена за умови рівності одиниці площі, обмеженої кривою розподілу. Густина рівномірного розподілу представлена на Рис.1.3

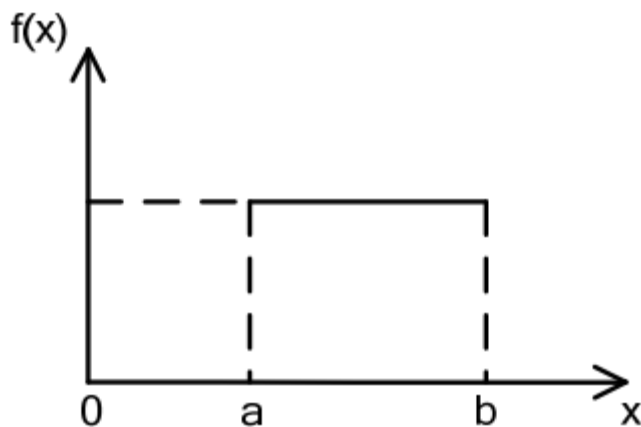


Рис.1.3 Густина рівномірного розподілу

Функція рівномірного розподілу $F(x)$ на відрізку $[a, b]$ дорівнює:

$$F(x) = \int_{-\infty}^x f(x) dx = \int_a^x \frac{1}{b-a} dx = \frac{x-a}{b-a}$$

3. Короткий опис програмного комплексу

Лабораторні роботи виконуються в пакеті Microsoft Excel, який представляє собою електронну таблицю. Документ Excel має розширення *.xls, *.xlsx і називається робочою книгою. Робоча книга складається з листів. Перемикатися між листами можна, використовуючи закладки (ярлички) в нижній частині вікна «Лист 1» і т.і.

Кожен лист представляє собою таблицю. Таблиця складається з стовпців і рядків. Кількість стовпців в листі - 256, рядків - 65536. Стовпці позначаються літерами латинського алфавіту (в звичайному режимі) від A до Z, потім йде AA-AZ, BB-BZ і т.і. Рядки позначаються звичайними арабськими числами.

На перетині стовпчиків та рядків знаходиться комірка. Кожна комірка має свою унікальну (в межах даного листа) адресу, яка складається з літери стовпчика (в звичайному режимі) і номера рядка. Адреса комірки використовується для роботи з даними (комірками) і формулами.

Статистичні розрахунки виконуються трьома способами: за допомогою формул, функцій і статистичної надбудови.

Формули вводяться в комірку шляхом набору з клавіатури. Формули починаються зі знака рівності (=), наприклад:

=A1 * \$B\$2

=СУММ(B2:B151)/150

Статистичні функції розташовують в формулі, вибравши в верхньому меню [Вставка → Функція → Категорія → Статистичні]. Довідкові матеріали з цих функцій можна отримати в довідковому керівництві Microsoft Excel ([Довідка → Довідка Microsoft Excel] або F1).

Приклад: обчислення середнього значення:

=СРЗНАЧ (B2: B151)

4. Методика виконання роботи

4.1. Генерація вихідних даних

При виконанні роботи також використовується статистична надбудова Microsoft Excel. Щоб активувати надбудову, необхідно вибрати *[Сервіс → Надбудови]* в меню і поставити галочку навпроти пункту **[Пакет аналізу]** (див. Рис. 1.4). Після цього буде досяжним пункт меню *[Сервіс → Аналіз даних]*.

Вихідні дані є вибіркою $\{x_1, x_2, \dots, x_n\}$, згенерованою за одним із законів розподілу в залежності від варіанту завдання (Табл. 1.2). Для генерації вихідних даних використовується функція *[Генерація випадкових чисел]* статистичної надбудови Microsoft Excel (див. Рис. 1.5). Після виклику функції, статистична надбудова запропонує вказати параметри генерації вибірки, як показано на Рис.1.6, 1.7.

Число змінних встановлюється рівним 1, обсяг вибірки n (число випадкових чисел). Параметри розподілу задаються відповідно до варіанту завдання (див. Табл.1.2). При генеруванні рівномірного розподілу параметрами є мінімальне і максимальне значення діапазону (*min* і *max*), які вводяться в вікні *[Параметри → Між]*. Для нормального розподілу вказують середнє значення μ і стандартне відхилення σ .

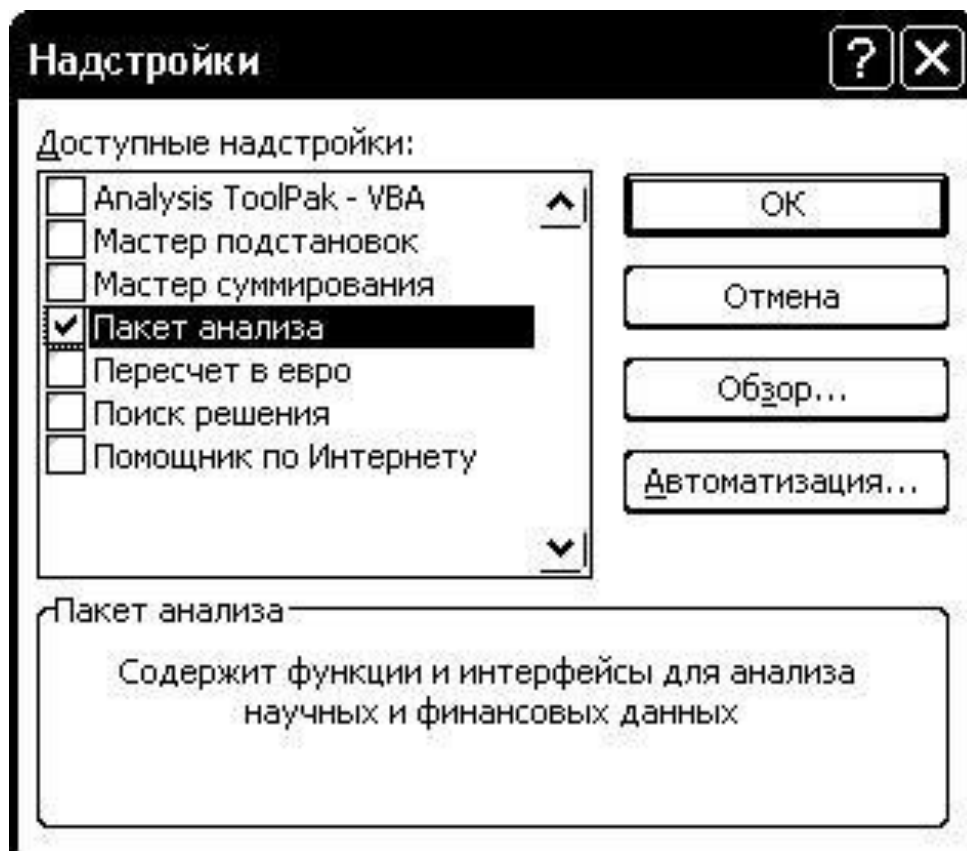


Рис1.4 Включення статистичної надбудови MS Excel

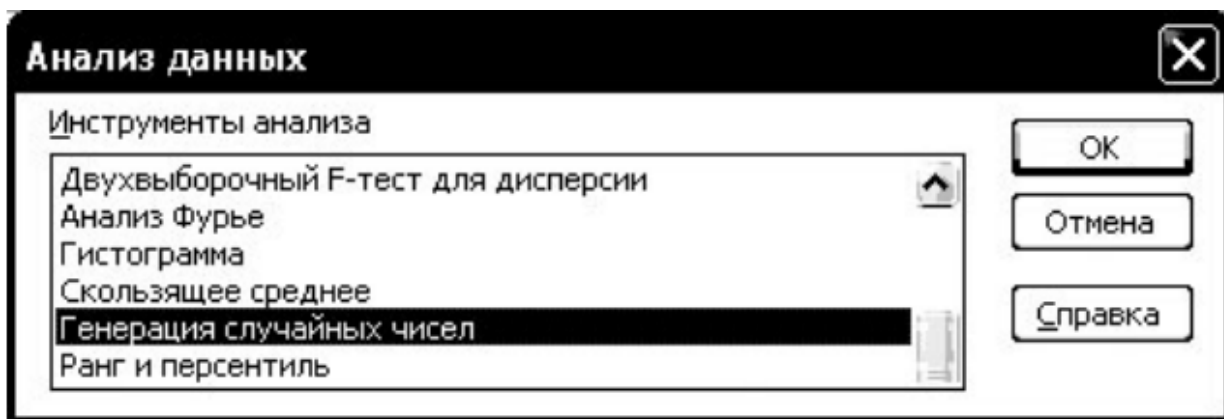


Рис. 1.5 Виблик функції «Генерація випадкових чисел»

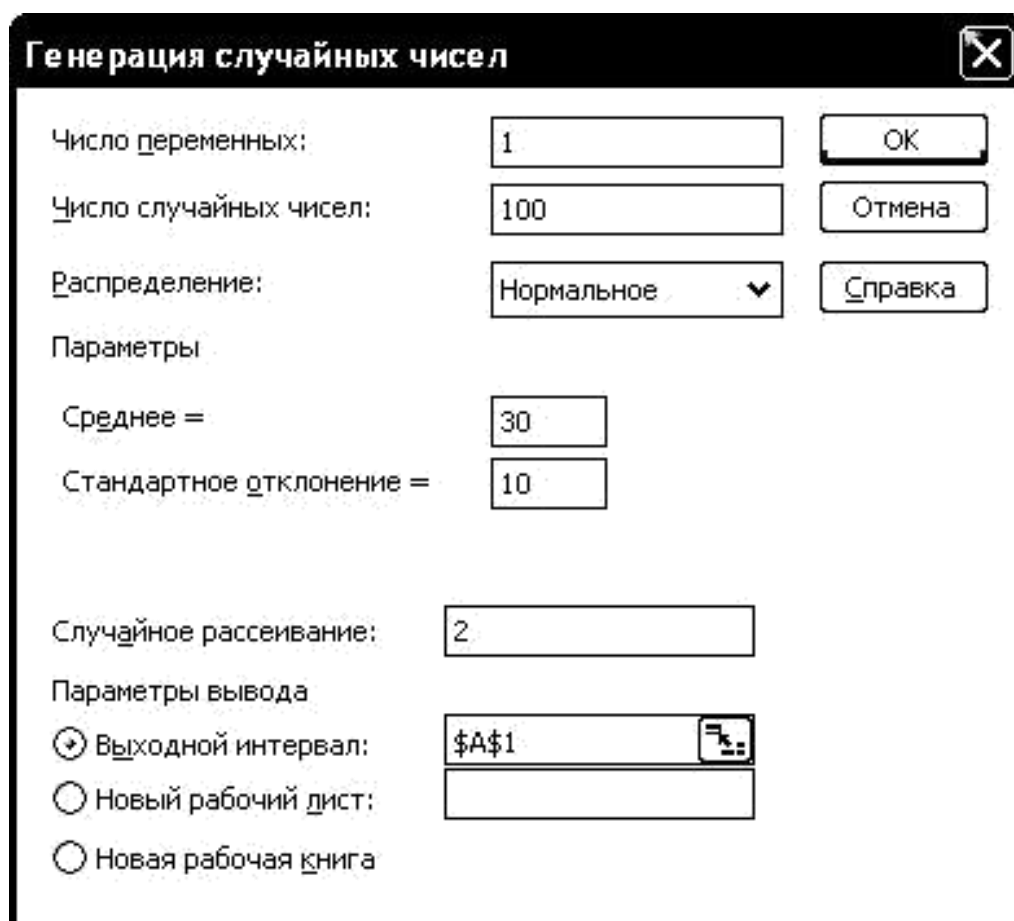


Рис. 1.6. Вибір параметрів генератора нормального розподілу

Випадкове розсіювання – параметр генератора псевдовипадкових чисел, що визначають початок послідовності.

Якщо задавати одно й те саме значення параметра, можна кожного разу отримати одну й ту саму послідовність.

Згенеровані випадкові числа необхідно округлити до цілих і розмістити в другій стовпчик таблиці. Для цього можна скористатися математичної функцією *ОКРУГЛ*, що має 2 аргументи: число, що округлюється, і значення десяткового розряду, до якого його потрібно округлити. Число розрядів дорівнює 0 в разі округлення до цілого.

Задаємо функцію округлення для однієї з комірок стовпчика «Округлені числа». Потім комірку з формулою копіюємо в буфер і вставляємо в усі інші комірки стовпця (див. Рис. 1.8). Надалі, в роботі використовуються тільки округлені значення.

Рис. 1.7. Вибір параметрів генератора рівномірного розподілу

	B2		fx =ОКРУГЛ(A:A;0)		
	A	B	C	D	E
1	Генератор	Рост, см		Рост, см	
2	186,207007	186			
3	180,8470412	181	Среднее		171,3666667
4	174,0076907	174	Стандартная ошибка		1,395531913
5	186,4932707	186	Медиана		173,5
6	175,9891049	176	Мода		175
7	170,5842769	171	Стандартное отклонение		17,09170554
8	183,2322153	183	Дисперсия выборки		292,1263982
9	166,4345225	166	Эксцесс		67,75392775
10	174,5641957	175	Асимметричность		-6,802877293
11	160,9356975	161	Интервал		189
12	185,1100192	185	Минимум		150
13	168,5142064	169	Максимум		189
14	155,6472976	156	Сумма		25705
15	189,3441572	189	Счет		150

Рис. 1.8. Генерація випадкових чисел і описова статистика

Варіанти завдання

Таблиця 1.2

№ варіанта	Випадкова величина	Параметри		n
		min	max	
	Рівномірний розподіл			
1	Зріст, см	155	190	150
2	Вага, кг	40	90	100
3	Витрати на споживання	2000	3500	170
4	Місячна зарплата, грн.	3000	7000	180
5	ВВП. трлн. дол.	2	5	180
	Нормальний розподіл	μ	σ	
6	Ціна автомобіля, тис. дол.	20	5	140
7	Число студентів у групі, чол.	31	10	190
8	Витрати енергії в умовних одиницях при проведенні експерименту	5000	500	170
9	Чисельність складових	20	4	150
10	Тривалість експерименту, хв.	240	60	140

4.2. Обчислення показників варіації

Для обчислення показників варіації застосовується функція Описова статистика статистичної надбудови Microsoft Excel (див. Рис. 1.9). У діалоговому вікні потрібно вибрати *Вхідний інтервал*, *Мітки в першому рядку*, *Вихідний інтервал* і *Підсумкова статистика*. Розрахуйте показники варіації, вказані в Табл 1.1

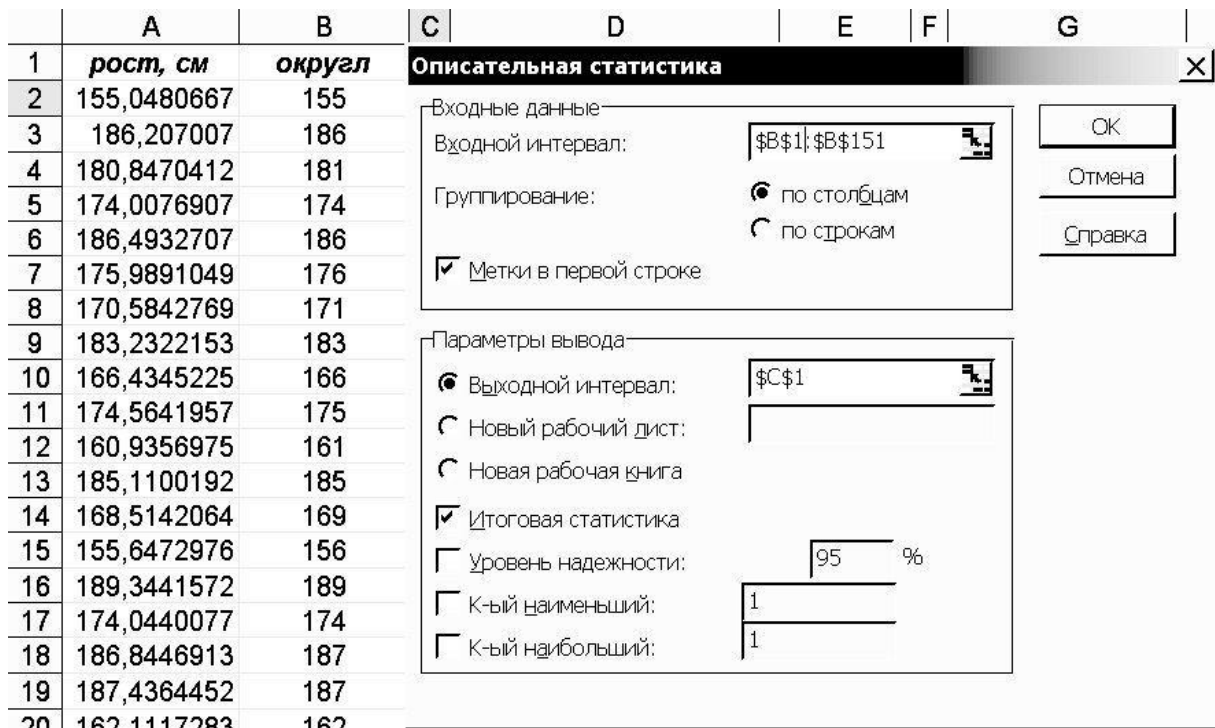


Рис. 1.9 Використання функції «Описова статистика»

При аналізі показників варіації можна використовувати такі правила:

- Вибірка вважається однорідною, якщо коефіцієнт варіації $V\sigma \leq 30\%$;

- Якщо коефіцієнти асиметрії та ексцесу близькі до нуля, то форму розподілу можна вважати близькою до нормальної. Критичні значення A і E обчислюють за формулами:

$$D(A) = \frac{6(n-1)}{(n+1)(n+3)}, \quad D(E) = \frac{24n(n-2)(n-3)}{(n+1)^2(n+3)(n+5)}$$

Критерій згоди: $|A| \leq 3\sqrt{D(A)}$; $|E| \leq 5\sqrt{D(E)}$

4.3. Групування за допомогою статистичної надбудови

Групування даних проводиться двома способами: за допомогою стандартних функцій Excel і статистичної надбудови. Спочатку слід створити таблицю нижніх меж інтервалів групування (Рис. 1.10).

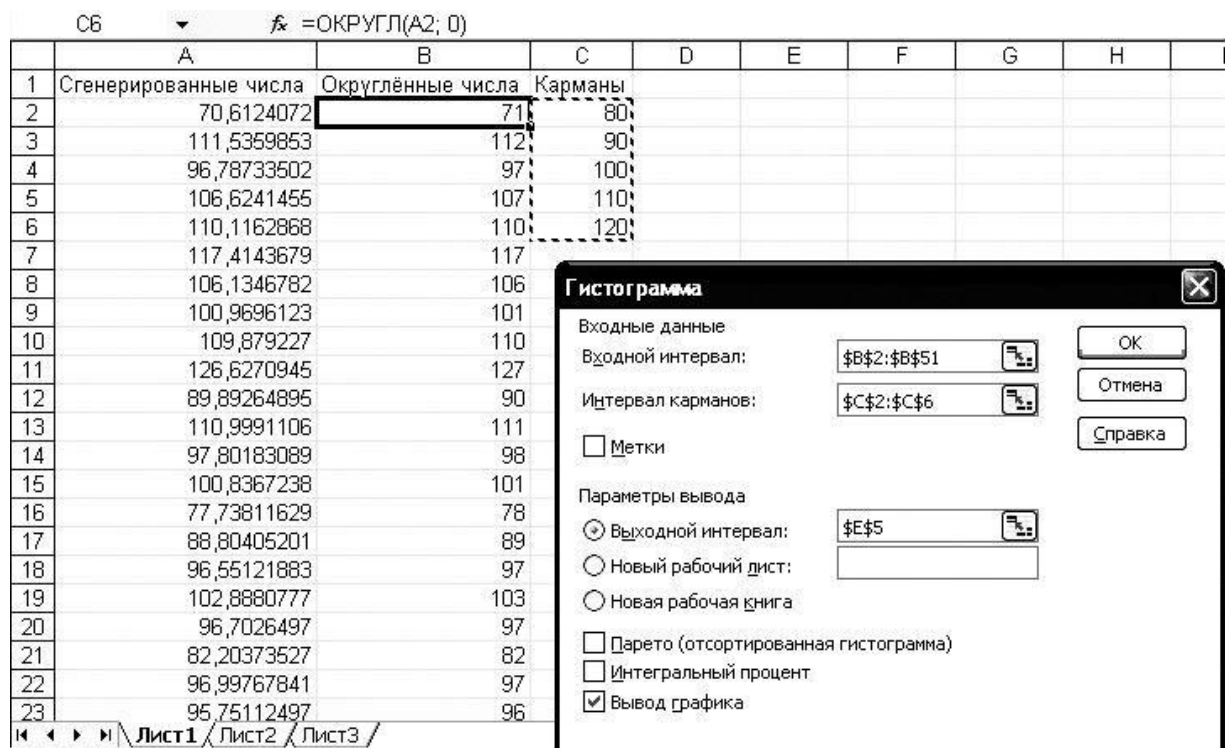


Рис.1.10 Встановлення функції Гістограма

Для групування даних за допомогою статистичної надбудови вибираємо меню [Сервіс → Аналіз даних → Гістограма]. Вказуємо наступні параметри:

Вхідні дані:

- Вхідний інтервал – вибірка вихідних даних;
- Інтервал кишень – нижні межі інтервалів групування.

Параметри виходу:

- Вхідний інтервал – розташування результатів групування на аркуші;
- Виведення графіка – побудова гістограми.
- Інтегральний відсоток – обчислення накопичених частот.

Результат роботи функції Гістограма представлений на Рис. 1.11.

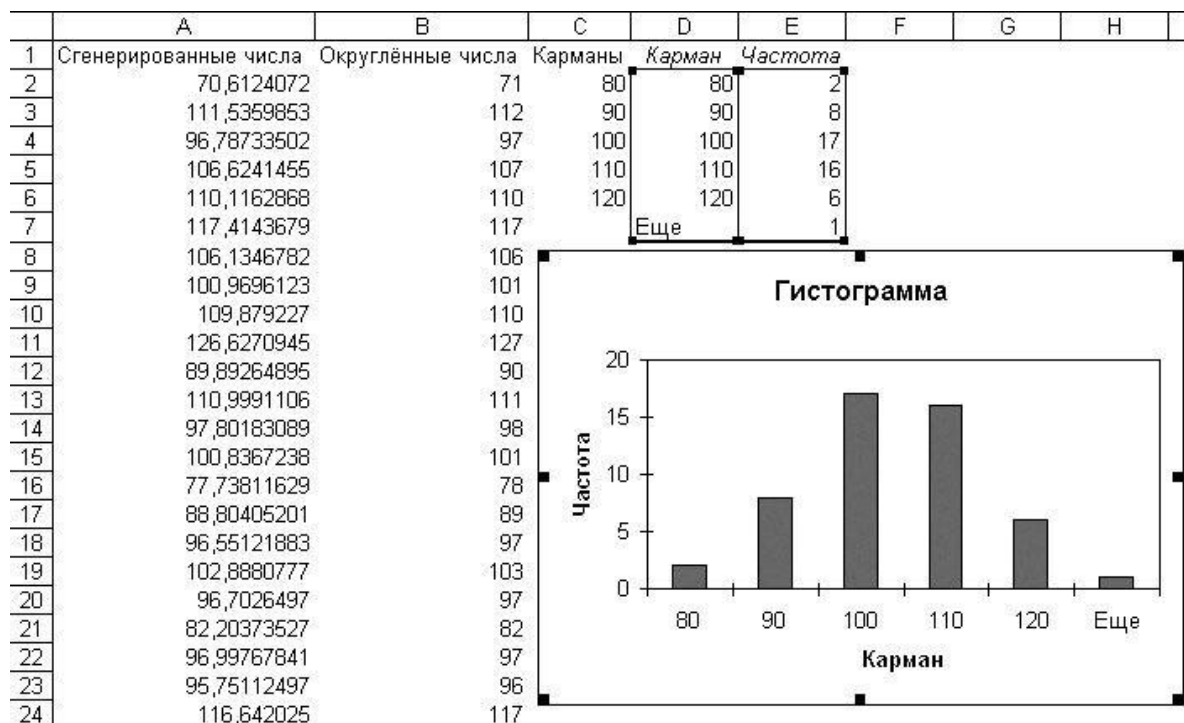


Рис. 1.11 Результат виклику функції «Гістограма»

Згенеровану таблицю необхідно доповнити відсутніми стовпчиками. Графік необхідно настроїти для коректного відображення. На Рис.1.12 показаний приклад рекомендованого оформлення таблиці та графіку.

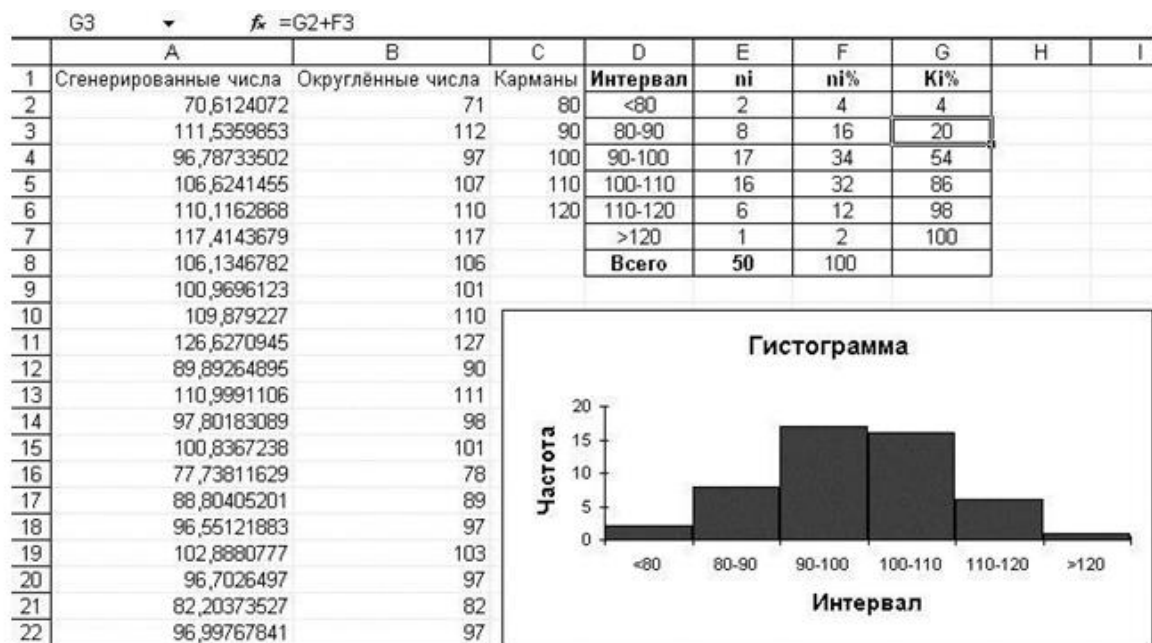


Рис.1.12 Налаштування результатів групування

Таблицю необхідно доповнити наступними стовпцями:

-Інтервал - підписи для стовпців гістограми у вигляді текстових міток, що описують межі інтервалів групування;

- Частота n_i ;
- Відносна частота $n_i, (\%)$;
- Накопичена частота $K_i, (\%)$.

У графі **Всього** виконують підрахунок суми частот.

Для обчислення відсутніх характеристик інтервалів слід використовувати формули. Наприклад, для обчислення накопиченої частоти для інтервалу 90..100 (комірка G3) використовується формула = G2 + F3. Обчислення відносної частоти для інтервалу 80..90 (комірка F3) виконується за допомогою формули: =100*E3/\$E\$8.

Після обчислень слід переконатися у відсутності грубих помилок. Наприклад, накопичена частота повинна дорівнювати 100%. Розташування стовпчиків гістограми має відповідати межах інтервалів групування даних. Для настройки графіка клацніть по стовпчику гістограми курсором і натисніть праву кнопку миші. Виберіть [*Формат рядів даних* → *Параметри*] і встановіть нульове значення параметрів *Перекриття* та *Ширина зазору*.

4.4. Групування за допомогою формул

Приклад групування за допомогою формул наводиться на рис. 1.13. Підрахунок частоти потрапляння в інтервал значень визначається як різниця кількості значень менше верхньої межі і менше нижньої межі інтервалу. Наприклад, частота для першого інтервалу (комірка G6) розрахована за допомогою функції *СЧЁТЕСЛИ*:

СЧЁТЕСЛИ(В:В;"<="&E6);*СЧЁТЕСЛИ*(В:В;"<="&D6)

Перша частина наведеної формули обчислює кількість комірок стовпчика В, значення яких менше або дорівнює вмісту комірки Е6 (верхня межа діапазону). Друга частина формули обчислює кількість комірок стовпчика В, значення яких менше або дорівнює вмісту комірки D6 (нижня межа діапазону). Таким чином, функція в цілому дає кількість комірок стовпчика В, значення яких потрапляють в інтервал між значеннями D6 і Е6.

G11	=СЧЁТЕСЛИ(B:B; "<="&E11)-СЧЁТЕСЛИ(B:B; "<="&D11)								
	A	B	C	D	E	F	G	H	I
1	70,61241	71							
2	111,536	112	80						
3	96,78734	97	90						
4	106,6241	107	100						
5	110,1163	110	110						
6	117,4144	117	120						
7	106,1347	106							
8	100,9696	101							
9	109,8792	110							
10	126,6271	127							
11	89,89265	90							
12	110,9991	111							
13	97,80183	98							

Сводная таблица					
Интервал			ni	ni, %	ki %
Нижняя граница	Верхняя граница	Середина			
0	70	35	0	0	0
70	80	75	3	3	3
80	90	85	13	13	16
90	100	95	37	37	53
100	110	105	32	32	85
110	120	115	13	13	98
120	130	125	2	2	100
Всего			100	100	-

Рис.1.13 Групування даних з використанням формул

За даною таблицею будується гістограма (рис. 1.14). для цього вибираємо в меню *[Вставка → Діаграма → Гістограма → Звичайна гістограма]*. Переходимо на закладку *Ряд* і натискаємо кнопку *Додати*. Натискаємо кнопку *Значення* і вказуємо діапазон значень частот. Натискаємо кнопку *Підпису осі X* і вказуємо діапазон міток для осі X. натискаємо кнопки *Далі → Готово*. Як міток можна вказати середини інтервалів групування.

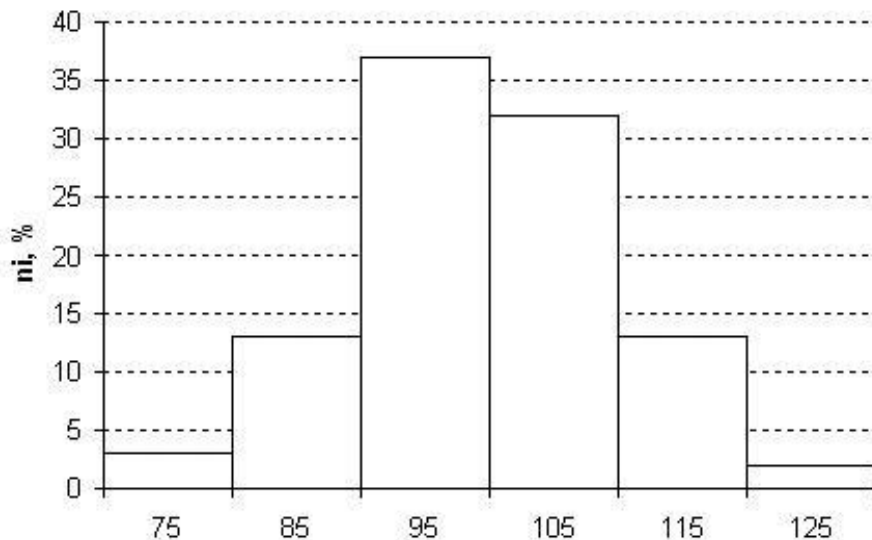


Рис. 1.14 Гістограма за результатом групування

4.5. Побудова графіків

Щоб додати полігон на графік гістограми, необхідно додати новий ряд даних. Для цього потрібно натиснути правою кнопкою миші по області графіка і в контекстному меню вибрати *[Вихідні дані → Ряд]* (Рис. 1.15).

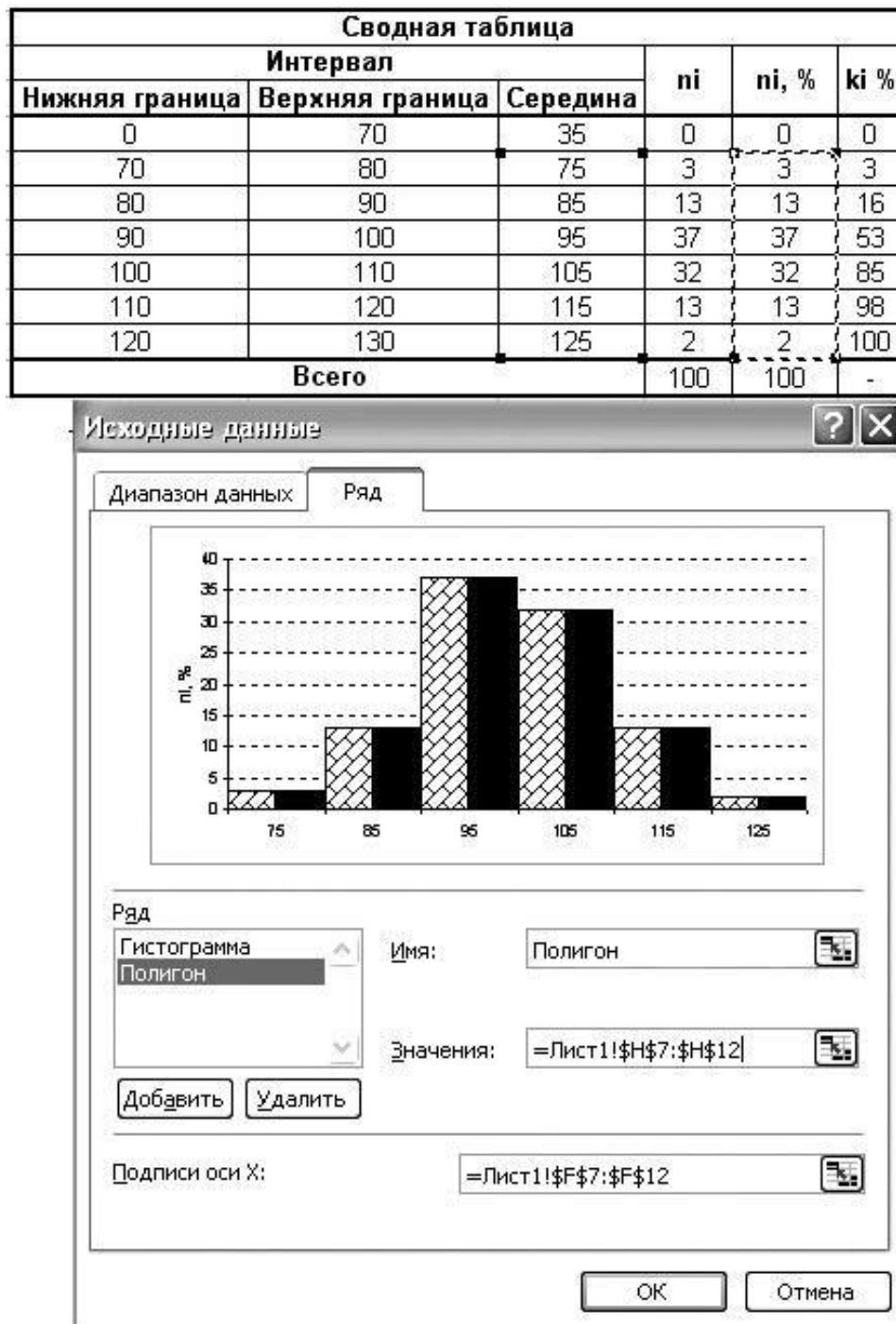


Рис.1.15 Додавання даних для побудови полігону.

Натискаємо кнопки *Додати* → *Значення* і вказуємо той же діапазон комірок, що і для гістограми. При роботі з рядами даних у вікні *Ім'я* можна ввести коментар, що пояснює зміст набору даних для графіка.

Накопичена частота будується на окремому графіку.

Обираємо в меню *[Вставка → Діаграма → Точкова діаграма → Точкова діаграма, на якій значення з'єднані відрізками]*. Додаємо ряд значень, вказавши діапазони по осі X і по осі Y. Накопичена частота (кумулята) – це лінійний графік накопичених частот, що представляє собою «інтеграл від гістограми». Приклад кумуляти показаний на Рис.1.16. На Рис.1.17 можна побачити остаточний варіант графіка гістограми та полігону.

Всі графіки повинні бути оформлені належним чином. Наприклад, мінімальне і максимальне значення по осі Y кумуляти повинні бути 0 і 100, відповідно (правий клік по осі Y, *[Формат осі → Шкала]*); лінія кумуляти повинна починатися від $x = 0$ (правий клік по осі X, *[Формат осі → Шкала]*, прибрати галочку *[Перетин з віссю Y (значень) в максимальному значенні]*. Підписи по осях повинні бути інформативними.

Положення осі Y встановлюється за допомогою пункту *[Вісь Y (значень) перетинає в значенні ..]*.

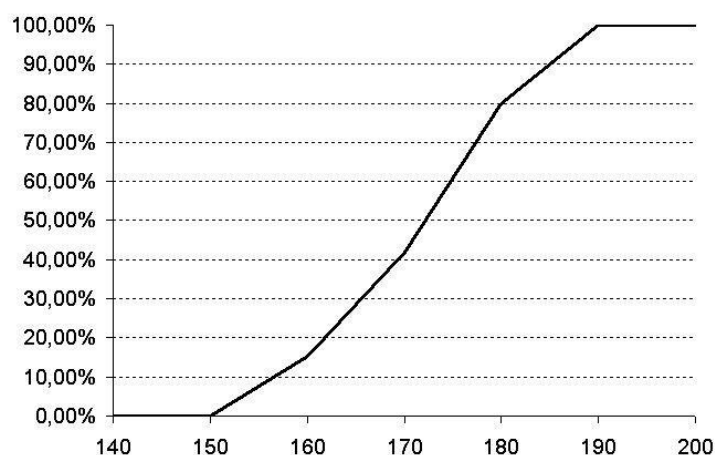


Рис.1.16 Приклад графіку накопиченої частоти

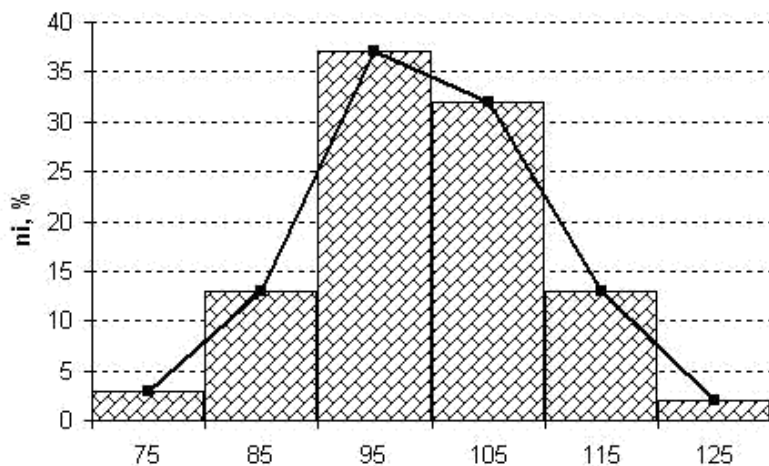


Рис. 1.17 Приклад гістограми та полігона.

4.6. Порівняння фактичного розподілу з теоретичним

Вихідна вибірка генерується за стандартним законом розподілу з параметрами відповідно до варіанта завдання (див. Розд. 4.1). Щоб порівняти фактичний і теоретичний розподіл, необхідно побудувати їх графіки. Для роботи з теоретичними розподілами використовуються готові статистичні функції, наприклад, *НОРМРАСП* і *НОРМОБР*. Опис статистичних функцій необхідно самостійно вивчити, викликавши *Довідку*.

На Рис. 1.18 показаний приклад графіка теоретичного і фактичного розподілів.

Щоб обчислити значення теоретичної ймовірності попадання випадкової величини в інтервал $[x_1, x_2)$, необхідно знайти різницю ймовірності попадання в інтервали $[0, x_2)$ і $[0, x_1)$ (див. Рис. 1.2). Наприклад, ймовірність попадання випадкової величини в інтервал $[100, 110)$ для нормального закону розподілу із середнім 100 і стандартним відхиленням 20, дорівнює:

$$= \text{НОРМРАСП}(110; 100; 20; \text{ІСТИНА}) - \text{НОРМРАСП}(100; 100; 20; \text{ІСТИНА}).$$

У разі рівномірного розподілу з нижньою межею 50 і верхньої 150, може бути використана наступна формула:

$$= (110-50) / (150-50) - (100-50) / (150-50) = (110-100) / (150-50) * 100\%.$$

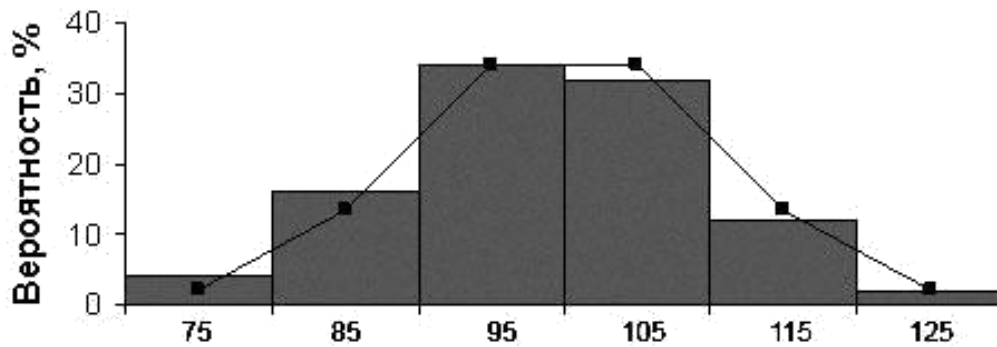


Рис.1.18. Емпіричний(стовпчики)та теоретичний(лінії) розподіл.

Вимоги до змісту та оформлення звіту

Звіт повинен бути продемонстрований як на паперовому носії, що містить показники варіації, результати угруповання, форми розподілу, статистичні графіки, так і в електронній формі, у вигляді файлу із заповненою таблицею та графіками.

Статистичні таблиці повинні містити всі елементи, необхідні для аналізу матеріалу: заголовки, назви, одиниці виміру, діапазони вимірювання. Основний заголовок повинен містити інформацію про об'єкт дослідження з прив'язкою до часу і місця.

Малюнки також повинні мати заголовок (підпис під малюнком). Позначення по осях повинні містити назви статистичних показників і одиниці вимірювання. Якщо на одному графіку зображено кілька кривих, необхідно вказати, що зображує кожна з них. Така інформація дається у вигляді легенди або в складі підпису під малюнком.

Звіт закінчується висновками щодо досліджуваної вибірки. У висновках не переказувати етапи проведених робіт, а коротко викладають результати. наприклад:

- Які дані проаналізовані;
- Чи є явні помилки і невідповідності в даних;
- Які закономірності в даних виявлені;
- Якими графіками і показниками це підтверджується.

Титульний аркуш звіту повинен містити всю інформацію, необхідну для однозначної ідентифікації авторів і роботи. Для цього на титульному аркуші вказують назву дисципліни, тему і номер роботи, варіант завдання, номер групи, прізвища та ініціали студентів, посаду, прізвище та ініціали викладача.

Порядок виконання роботи

1. Ознайомтеся з описом наступних функцій:
СТАНДОКЛ, СТАНДОТКЛП, НОРМРАСП.

2. Згенеруйте вихідні дані.
3. Обчисліть основні статистичні показники трьома способами:
 - За допомогою формул;
 - За допомогою статистичних функцій;
 - За допомогою статистичної надбудови.
4. Порівняйте результати розрахунків.
5. Зробіть висновок про однорідність вибірки.
6. Зробіть висновок про близькість до нормального розподілу.
7. Проведіть групування даних двома способами:
 - За допомогою стандартних функцій Excel;
 - За допомогою статистичної надбудови.
8. Побудуйте зведені таблиці і порівняйте результати групування, отримані двома способами.
9. Побудуйте гістограму, полігон і кумуляту двома способами:
 - За допомогою статистичної надбудови
 - «Вручну».
10. Побудуйте гістограми теоретичного та емпіричного розподілів на одному графіку.
11. Зробіть висновок про близькість емпіричного розподілу до теоретичного.
12. Оформіть звіт відповідно до вимог.

Контрольні питання

1. Як обчислюються основні статистичні показники?
2. Як проводиться групування?
3. Які існують статистичні графіки, як вони будуються і в чому їх призначення?
4. В чому різниця між теоретичним і емпіричним розподілом?
5. Які вимоги пред'являються до оформлення таблиць і графіків?
6. Що повинні містити висновки?