

Лекція 1

**Аналіз даних. Основні поняття.
Описова статистика.**

.

К.ф.-м.н. Щестюк Н.Ю.



IN AN
**INTERNET
MINUTE**

50,200
MOBILE APPS
DOWNLOADED¹



94
TWITTER
ACCOUNTS
CREATED²



2.4 MILLION
GOOGLE
SEARCHES³



1,389
UBER
RIDES⁴



30
IDENTITY
THEFTS⁵



2,083,333
MINUTES
USED ON
SKYPE CALLS⁶



142,361,111
EMAILS SENT
AND RECEIVED⁷



347,222
TWEETS⁸



120
LINKEDIN
ACCOUNTS
CREATED⁹



300
HOURS OF VIDEO
UPLOADED
ON YOUTUBE¹⁰



\$203,579
IN AMAZON
SALES¹¹



216,000
PHOTOS
POSTED TO
INSTAGRAM¹²



обмін даними

- **Uber** – найбільша в світі служба таксі, не володіє жодним авто
 - **Facebook** - найпопулярніший власник медіа, не створює жодного контенту
 - **Alibaba** - retailer, нічого не виробляє
 - **Airbnb** – сервіс короткострокової оренди житла, не володіє жодним помешканням
-
- сервіс для пошуку пари для шлюбу
 - сукупний дохід > 1 мільярда доларів
 - близько 4% шлюбів у США в 2012
 - більше 33 млн користувачів у понад 150 країнах



Основні компоненти аналізу даних



Dr

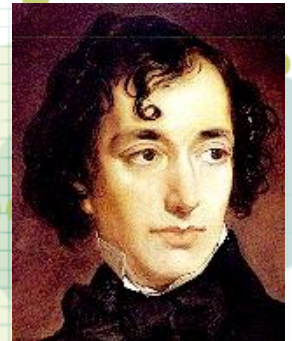
"There are three kinds of lies:
lies, damned lies, and statistics."

B. Disraeli

СТАТИСТИКА -

наука про збір, організацію,
аналіз та трактування даних.

(оцінити варіативність та
зменшити невизначеність)



За допомогою статистики вирішують наступні проблеми:

- оцінювання ризику вживання трансгенних продуктів харчування або нових вакцин
- передбачення кількості захворювань на грип за регіонами
- передбачення результатів наступних виборів
- управління самокерованими автомобілями
- голосові асистенти для смартфонів



Історія та нові напрямки розвитку статистики



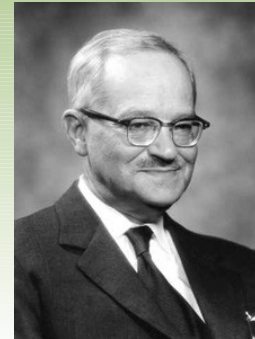
Карл Гаусс
(1777-1855)



Карл Пірсон
(1857-1936),
Метод моментів,
критерій «хі-
квадрат»



Рональд Фішер
(1890-1962)
Метод
максимальної
правдоподібності,
Fisher's exact test



Ежи Нейман
(1894-1977),
Біхевіористська
статистика



А.М. Колмогоров
(1903-1987),
Основи непар-
аметричної
статистики

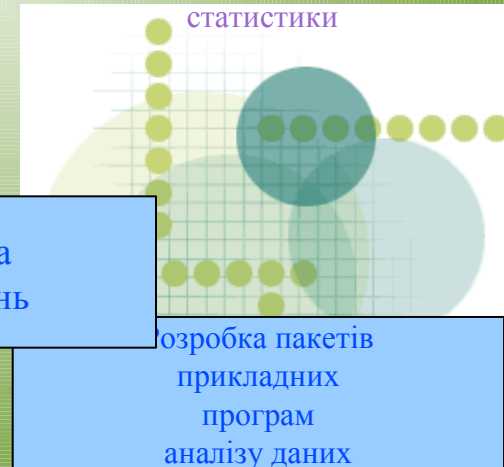
Статистика об'єктів
нечислової

прир

Математичні
методи планування
експерименту

Статистична
теорія рішень

розробка пакетів
прикладних
програм
аналізу даних



Вибіркове обстеження (Dalenius, 1974)

Потреба в статистичній інформації в предметній області

План вибіркового обстеження
(методи, обсяг вибірки)

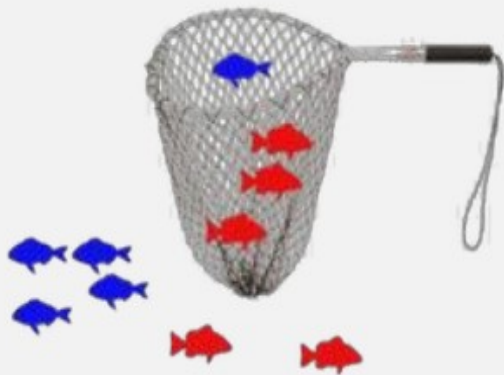
Збір даних

Обробка, очищення та аналіз даних
(побудова моделей, перевірка гіпотез)

Трактування, презентація та використання результатів



Генеральна сукупність



Генеральна сукупність - усі об'єкти, які хотів би вивчати дослідник при необмеженій кількості ресурсів.

Як сформувати вибірку



Простий випадковий вибір
Всі об'єкти мають однакову можливість бути вибраними. Випадковим чином обирається n об'єктів



Вибір з заміною
Після того, як об'єкт вибрано, він повертається і може бути обраний повторно



Вибір без заміни
Після того, як об'єкт вибрано, він вилучається і не може бути обраний повторно



Стратометричний вибір
Сукупність ділиться на гомогенні групи (населення за рівнем освіти чи віковою групою)

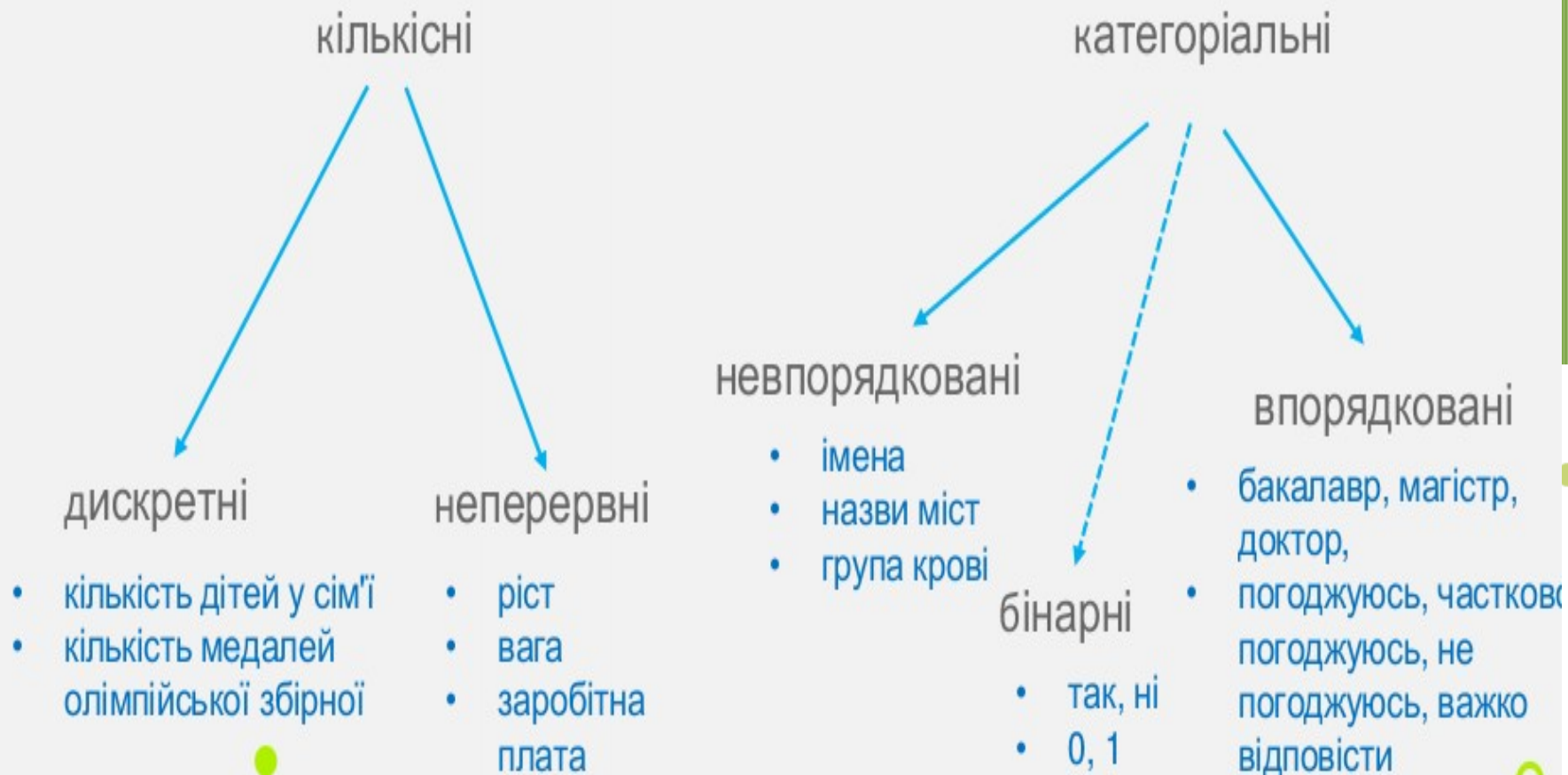


Кластерний вибір
Сукупність ділиться на кластери (місто на райони)



Систематичний вибір
Елементи сукупності впорядковуються і вибирається кожен k -ий елемент (елементи на конвейєрі з метою виявлення дефектів)

Типи даних



Основні поняття

Генеральна сукупність

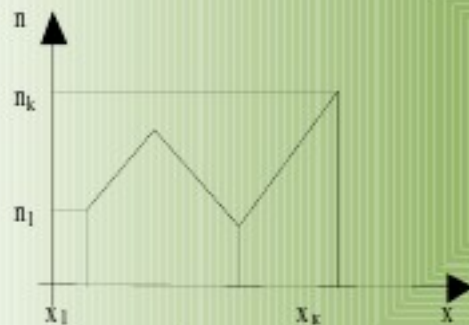
Вибірка
(n -об'єм)

Варіаційний ряд

Статистичний розподіл

x_i	x_1	x_2	x_3		x_{k-1}	x_k
n_i	n_1	n_2	n_3		n_{k-1}	n_k

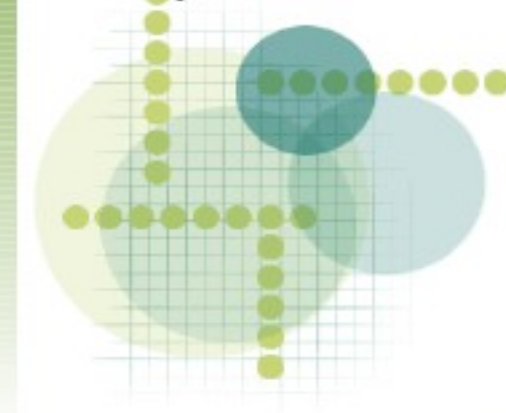
Емпірична функція розподілу



Полігон



Гістограма



Вибірковий метод

Генеральна сукупність — це множина всіх значень, яких може набувати дана випадкова величина.

$\{x_1, \dots, x_n\}$ — це **вибірка** спостережень за випадковою величиною ξ із генеральної сукупності (*набір з незалежних і однаково розподілених випадкових величин ("копій")*).

n — **обсяг (об'єм) вибірки**, кількість елементів вибірки

Впорядковуємо вибірку за зростанням:

$$x_1^* \leq x_2^* \leq \dots \leq x_n^*$$

Таким чином впорядковану вибірку називають **варіаційним рядом**.



Статистичний ряд

Статистичним рядом називається послідовність пар (z_i, n_i) , $i = \overline{1, l}$

Групований статистичний ряд (Frequency tables) – це сукупність пар (x_i^*, n_i^*) , $i = \overline{1, k}$, де x_i^* – середина i -того інтервалу,

n_i^* – частота попадання у i -й інтервал.

Одна з рекомендацій щодо вибору числа інтервалів групування – формула Стерджесса:

$$k = 1 + \log_2 n$$



ВІЗУАЛІЗАЦІЯ

Полігон та гістограма частот

Полігон частот – ламана з вершинами у точках

$$(x_i^*, n_i^*), \quad i = \overline{1, k}$$

Гістограма частот – це ступінчаста фігура, побудована з прямокутників, основами яких є інтервали групування, а висоти

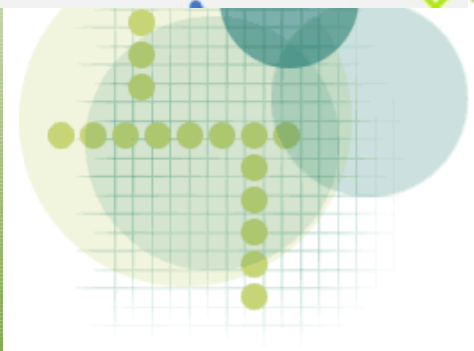
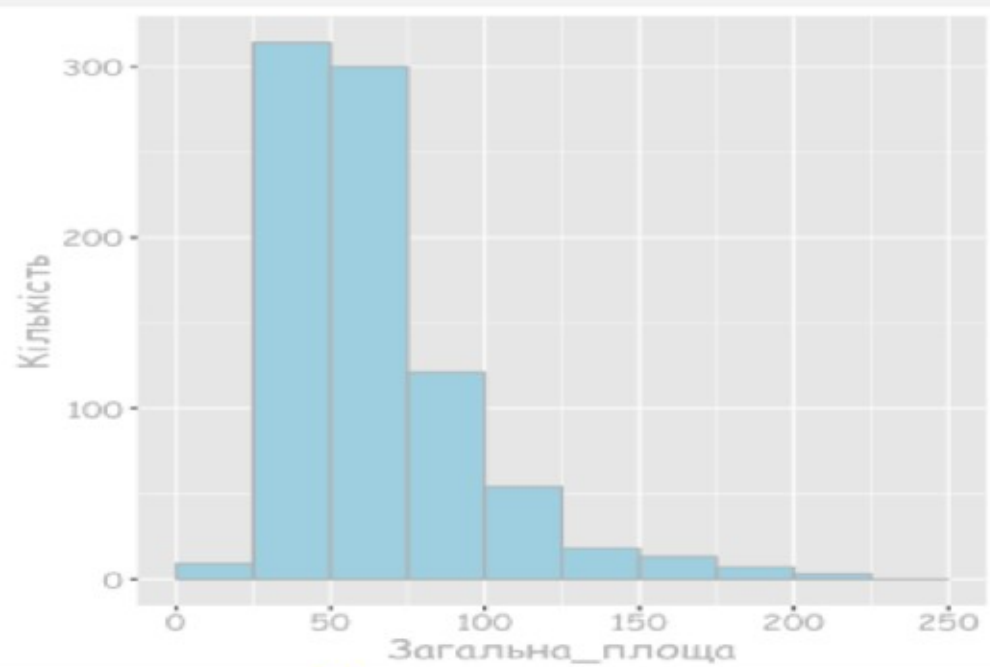
$$h_i^* = \frac{n_i^*}{\delta n},$$

де δ – це є довжина інтервалу групування.

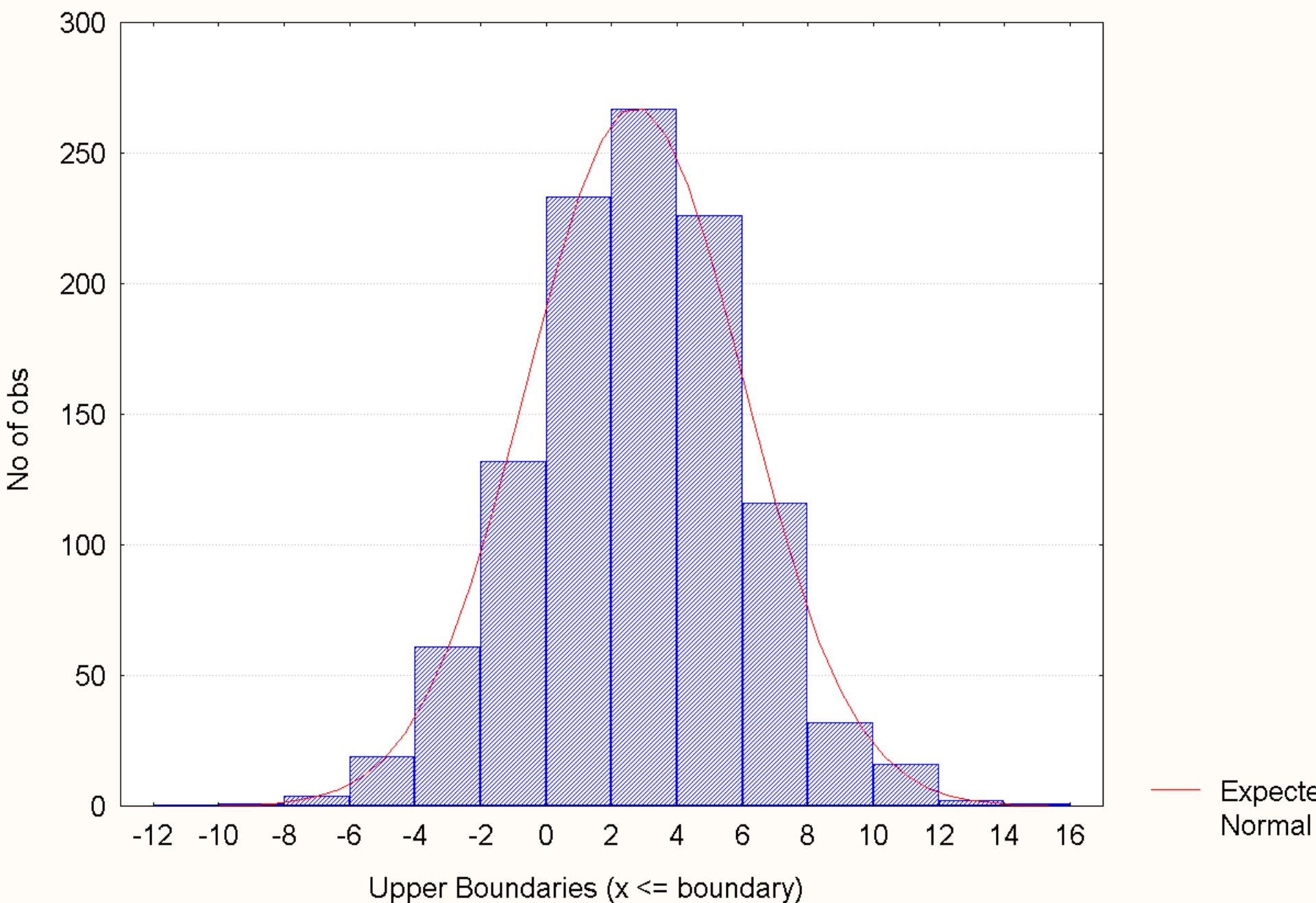
Площа гістограми дорівнює 1.



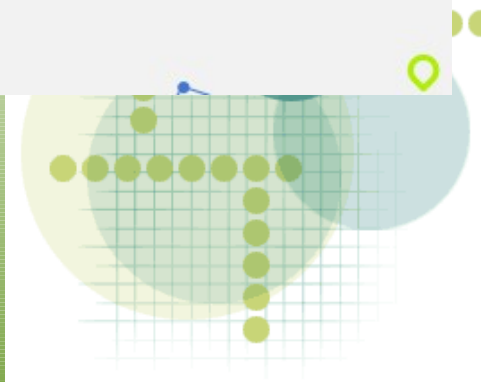
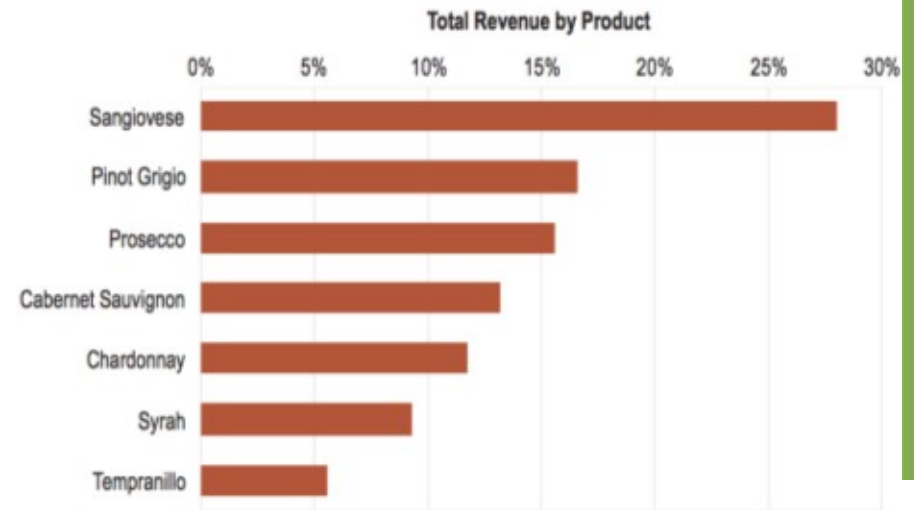
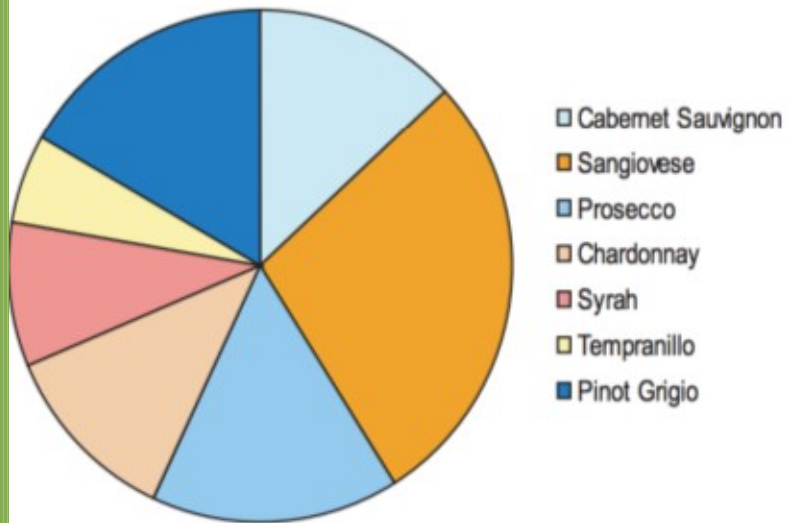
гістограма



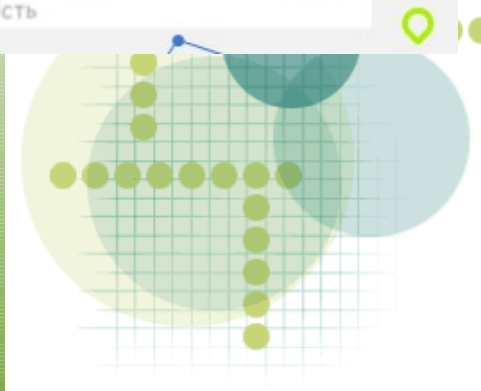
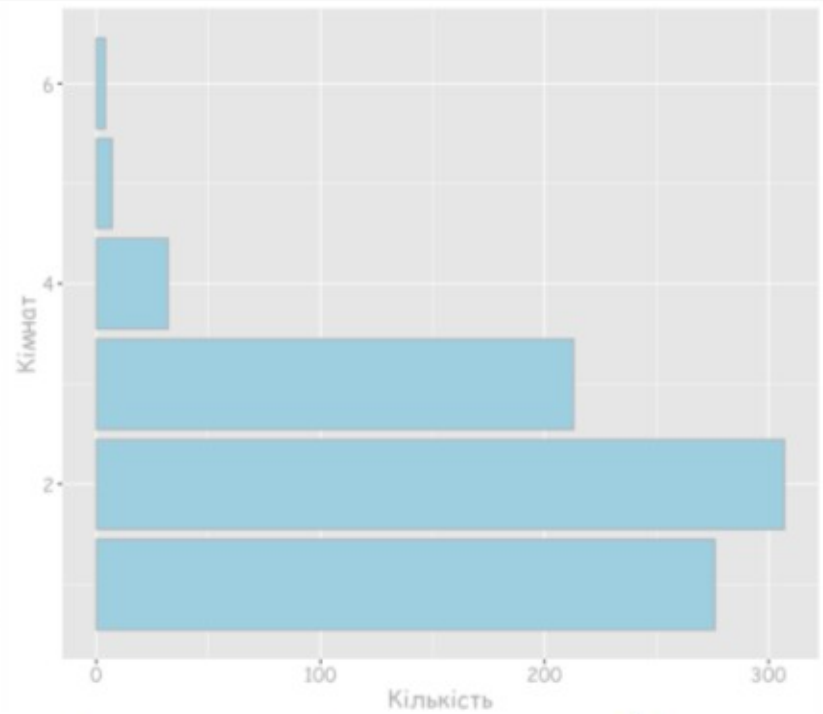
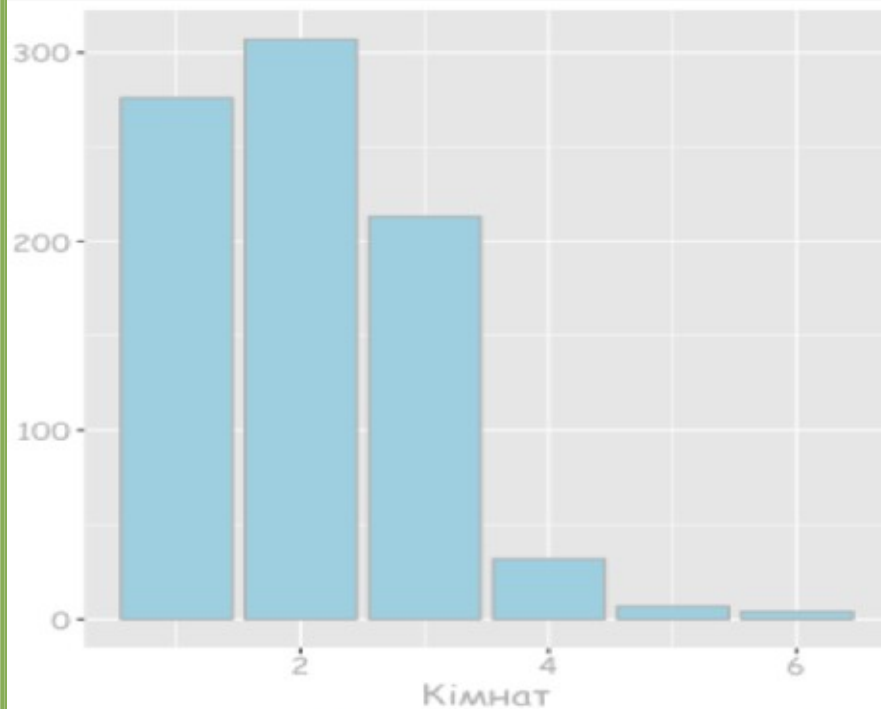
VAR3: =VNormal(Rnd(1); 2.678; 3.345);



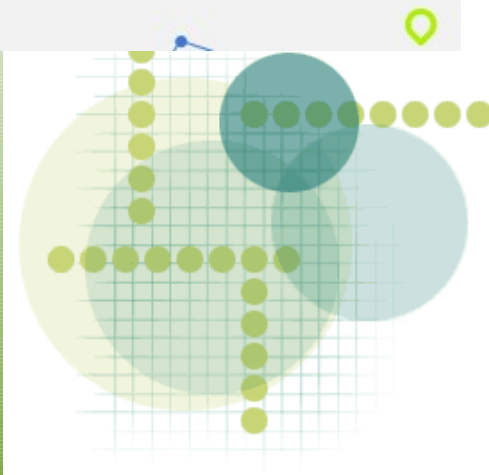
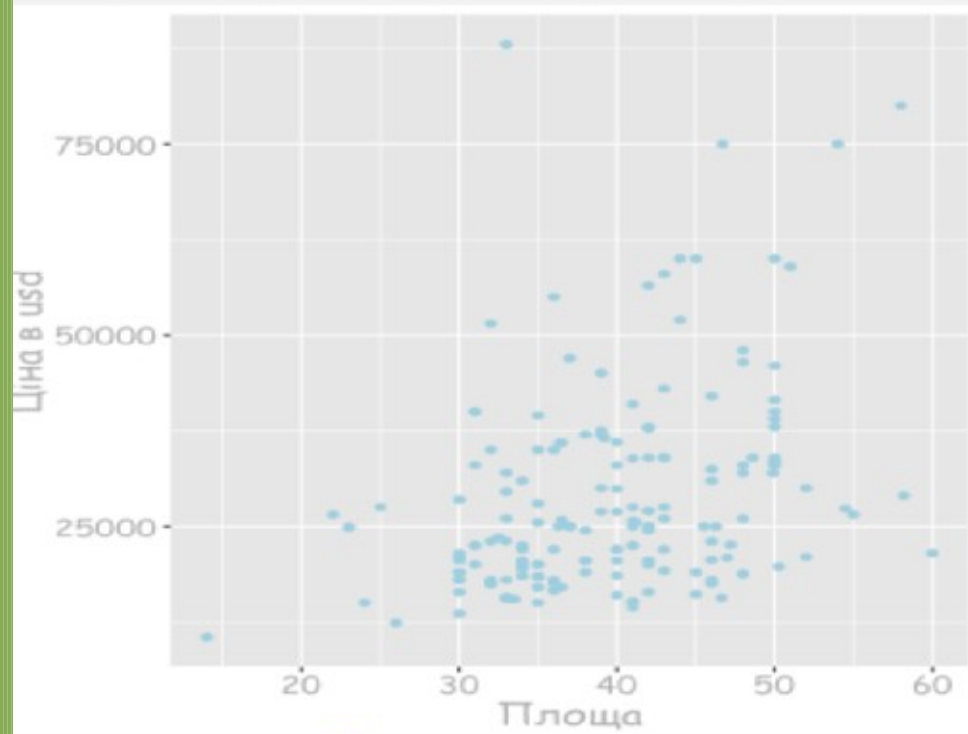
кругова діаграма



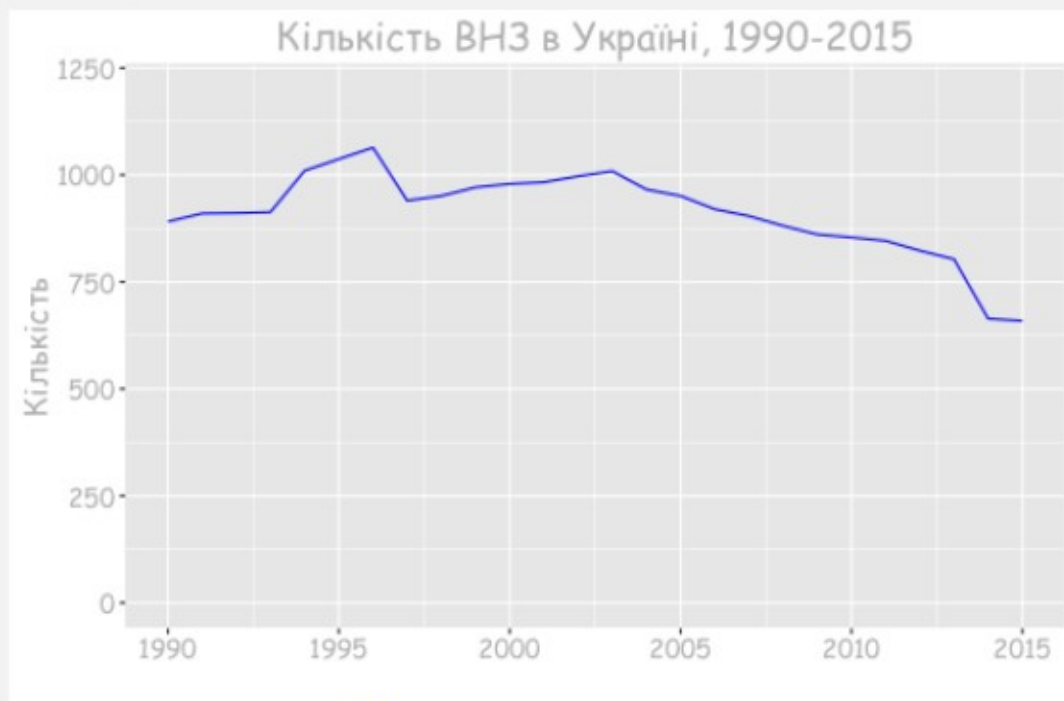
СТОВПЧИКОВА ДІАГРАМА



діаграма розсіювання



лінійний графік



Емпірична функція

Емпірична функція розподілу визначається через статистичний ряд співвідношенням

$$F_n^*(y) = \frac{1}{n} \sum_{i=1}^n I(x_i < y) \quad , \quad F_n^*(y) = \begin{cases} 0, & y \leq x_{(1)}; \\ \frac{n_1}{n}, & x_{(1)} < y \leq x_{(2)}; \\ \dots & \dots \\ \frac{n_1 + n_2 + \dots + n_k}{n}, & x_{(k)} < y \leq x_{(k+1)}; \\ \dots & \dots \\ 1, & y > x_{(n)}. \end{cases}$$

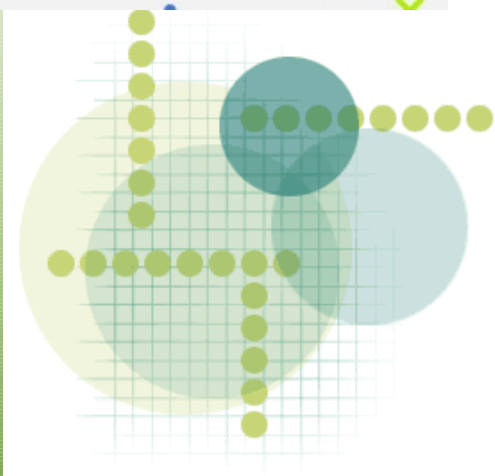
Кумулятивна крива – ламана з вершинами у точках

$$\left(x_i^* + \frac{\delta}{2}, \frac{1}{n} \sum_{j=1}^i n_j^* \right), i = 1, \dots, k$$



який тип діаграми краще застосовувати?

- Порівнювати значення: стовпчикова діаграма, лінійний графік, графік розсіювання
- Зрозуміти композицію: (виділити складові) - стовпчикова діаграма, кругова діаграма
- Оцінити розподіл даних: лінійний графік, графік розсіювання, стовпчикова діаграма, гістограма
- Зрозуміти тренд: лінійний графік, стовпчикова діаграма
- Зрозуміти відношення між даними: лінійний графік, графік розсіювання



Характеристики центральної тенденції

Середнє значення

$$\bar{x} = \frac{\sum x}{N}$$

Середня ціна
однокімнатної
квартири \$25880,
середня площа – 40 м²

Медіана

ділить вибірку навпіл

Медіана ціни
однокімнатної
квартири \$23000,
медіана площі 40 м²

Мода

найчастіше
трапляється

Характеристики центральної тенденції

Вибіркове середнє:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

**Вибіркове середнє для
згрупованої вибірки:
групування**

$$\bar{x} = \frac{1}{n} \sum_{i=1}^k n_i x_i^* \quad \text{де } k - \text{число інтервалів}$$

Мода (Md) – це елемент статистичного ряду з найбільшою частотою.

Медіана (Me)- це значення середнього елемента вибірки

Середнє геометричне -

$$\bar{x}_g = \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n}$$



середнє значення

3, 2, 5, 6, 7, 2, 3, 3

$$\bar{x} = \frac{3+2+5+6+7+2+3+3}{8} = \frac{31}{8} = 3.875$$

$$\bar{x} = \frac{\sum x}{N}$$

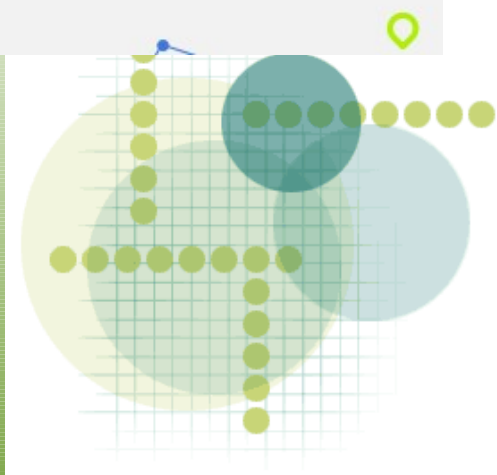
Місце серед членів ООН	Місце серед усіх територій	Країна	Загальна очікувана тривалість життя при народженні	Чоловіча очікувана тривалість життя при народженні	Жіноча очікувана тривалість життя при народженні
	1	 Макао	84,36	81,39	87,47
1	2	 Андорра	82,51	80,33	84,84
2	3	 Японія	82,12	78,8	85,62
3	4	 Сінгапур	81,98	79,37	84,78
4	5	 Сан-Маріно	81,97	78,53	85,72
	6	 Гонконг	81,86	79,16	84,79
5	7	 Австралія	81,63	79,25	84,14
6	8	 Канада	81,23	78,69	83,91
7	9	 Франція	80,98	77,79	84,33
8	10	 Швеція	80,86	78,59	83,26

медіана

3, 2, 5, 6,	7, 2, 4, 3
2, 2, 3, 3,	4, 5, 6, 7

$$\frac{3+4}{2} = 3.5$$

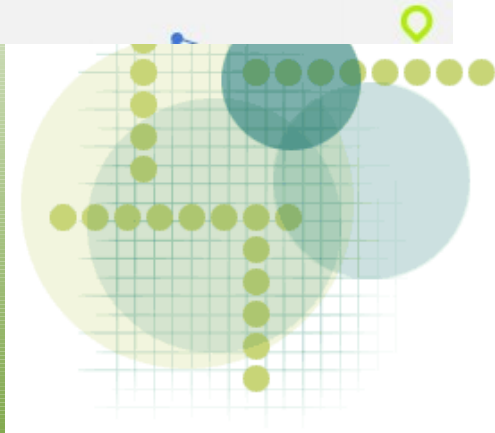
значення, яке ділить вибірку навпіл



мода

3, 2, 5, 6, 7, 2, 3, 3
2, 2, 3, 3, 3, 5, 6, 7

значення, яке найчастіше трапляється



Квантилі

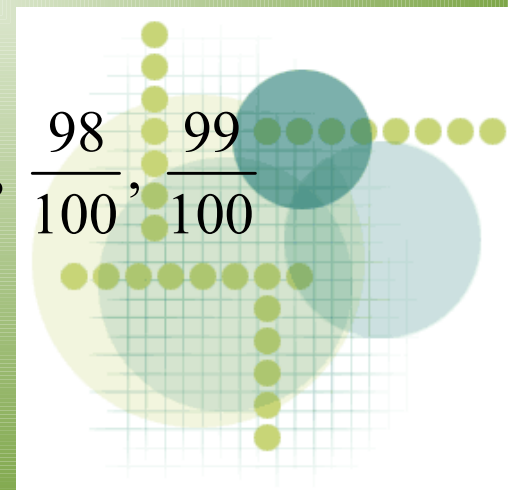
Вибірковий квантиль порядку p – це абсциса точки, яка лежить на кумулятивній кривій та має ординату p .

Медіана – це квантиль порядку $\frac{1}{2}$.

Квартилі – квантилі порядку $\frac{1}{4}$, $\frac{2}{4}$ та $\frac{3}{4}$.

Децилі – квантилі порядку $\frac{1}{10}$, $\frac{2}{10}$, ..., $\frac{9}{10}$.

Процентилі – квантилі порядку $\frac{1}{100}$, $\frac{2}{100}$, ..., $\frac{99}{100}$.



Характеристики розсіювання

Розмах (range) вибірки (ω), тобто різниця між найбільшим та найменшим елементами вибірки.

Дозволяє оцінити розкид елементів у вибірці.

Міжквартильний розмах (inter-quartile range або mid-spread) – це є різниця між нижнім і верхнім квантилем.

Дозволяє оцінити розкид 50% елементів вибірки і не враховує вплив викидів.



Характеристики розсіювання

Вибіркова дисперсія — $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$

Для групованих даних — $s^2 = \frac{1}{n} \sum_{i=1}^k n_i (x_i^* - \bar{x})^2$

Незміщена вибірка дисперсія

$$s_0^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Середньоквадратичне відхилення – це квадратний корінь з вибіркової дисперсії: $\sigma_B = \sqrt{s^2}$

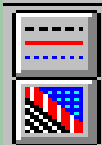


STATISTICA: Basic Statistics and Tables

File Edit View Insert Layouts Analysis
Graphs Options Window Help

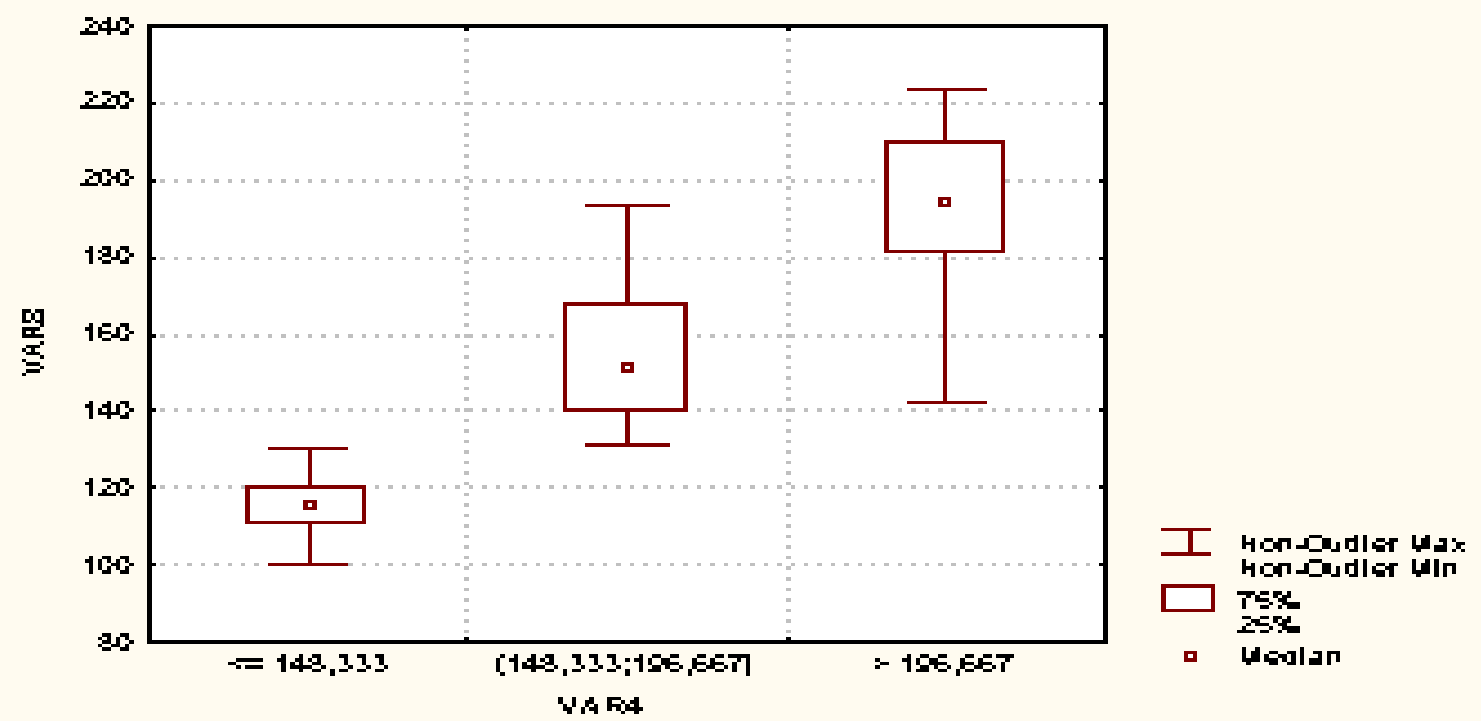
4,75; 34

Dyn.



Graph22: Box Plot

Box Plot (diggie.STA 20*310c)



Ready

Output:WINDOW

Sel:OFF

Weight:ON

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
	VAR1	VAR2	VAR3	VAR4	VAR5	VAR6	VAR7	VAR8	VAR9	VAR10	VAR11	VAR12	NEWVAR13	NEWVAR14	NEWVAR15	NEWVAR16	NEWVAR17	NEWVAR18
	.438	.534	2.215	.376	.410	1.030	.222	14.195	2.665	.330	4.261	3.052	-1.562	15.610	-.953	.771	3.709	.
	.291	.548	-2.984	.287	.641	1.531	.610	4.569	4.713	.181	-5.898	7.823	-11.534	-1.246	3.927	-.412	2.657	.
	.521	.353	3.409	.528	.826	-.084	.787	-6.207	4.217	.238	5.194	5.169	-8.495	23.806	-2.565	1.656	4.152	.
	.795	.906	.011	.360	.229	.011	.331	8.707	5.816	.416	6.034	4.816	2.255	1.048	3.045	2.488	5.996	2.
	.162	.234	4.709	.676	.524	.426	.312	1.193	3.921	.459	4.383	11.398	-19.419	32.949	-4.103	.238	1.958	.
	.457	.596	3.017	.545	.807	1.610	.613	-2.002	5.574	.223	7.810	1.882	-3.753	21.097	-.912	1.388	3.762	1.
	.276	.553	-1.475	.428	.753	3.147	.672	7.417	5.475	.206	4.683	7.172	-2.573	-6.059	3.012	.076	2.583	1.
	.314	.319	.718
	.729	.544	.712
	.479	.826	2.444
	.319	.671	.275
	.243	.425	1.873
	.421	.729	2.222
	.403	.720	4.978
	.443	.683	2.504
	.711	.565	8.893
	.293	.479	.324
	.146	.889	6.505
	.421	.582	4.453
	.319	.381	.363
	.547	.421	.988
	.463	.453	-1.581
	.323	.424	-1.224
	.288	.359	2.319
	.540	.603	6.488
	.480	.677	3.540
	.384	.319	-.713
	.128	.613	.876
	.468	.622	10.673
	.484	.434	-.866
	.638	.456	-.259
	.169	.751	1.150
	.297	.881	6.614
	.759	.753	4.169
	.473	.573	.237	.275	.846	.740	.507	-2.469	5.965	.188	-1.289	6.205	.894	2.124	.770	.493	3.607	.
	.784	.331	-1.267	.102	.896	-.690	.401	-.608	4.458	.071	2.192	6.596	-4.189	-5.949	-.370	.423	5.313	-2.
	.506	.872	5.129	.356	.302	4.797	.548	7.194	4.159	.154	1.935	1.493	-15.252	35.904	-2.259	.905	4.139	2.
	.806	.786	1.666	.105	.646	-.702	.517	4.684	3.303	.443	4.105	9.124	-.089	11.558	-.982	.779	5.793	-.1.
	.369	.667	2.473	.332	.742	-.885	.534	5.623	6.470	.277	2.008	4.446	-1.515	17.223	-.729	.164	3.106	1.
	.416	.598	3.109	.215	.872	.640	.444	16.439	3.822	.161	1.339	5.932	-5.537	21.375	-2.339	-.171	3.140	.
	.146	.512	2.867	.420	.613	-1.287	.481	-4.591	2.414	.176	11.479	4.184	-4.237	20.086	-1.559	-.652	1.852	1.
	.283	.578	6.768	.589	.805	3.919	.787	-1.176	4.960	.251	-2.800	8.261	-36.202	47.334	-4.643	.567	2.689	1.
	.633	.729	4.161	.632	.651	-.609	.343	-4.903	6.358	.596	-6.410	2.558	-8.796	29.139	-1.178	2.493	4.960	2.
	.419	.277	1.060	.448	.489	-3.124	.635	3.180	3.892	.109	-20.611	5.628	-.388	7.767	-.835	.914	3.587	.

2D Box Plots

Graph Type:

Box-Whiskers

Whiskers

Boxes

Columns

High-Low Close

Regular

Multiple

Variables:

Categories: none

Variables: none

OK

Cancel

Options...

BOX CATEGORIES

Variable: none

Integer Mode

Categories: 10

Boundaries: none

Codes: none

Multiple Subsets

Change Variable

FIT

Off

Linear

Custom: none

Non-Outlier Max

75%

Median

25%

Non-Outlier Min

Middle Point

Value: Median

Style: Mean

Median

Pooled Variance

Multiple Box Layout

Shifted

Overlaid

Connect Middle Points

Box

Value: Percentiles

Coefficient: 25

Whisker

Value: Non-Outlier Range

Coefficient: 1

Outliers

Outl. & Extremes

Coefficient: 1.5

Trim distrib. extremes: 0%

Як вибирати показник розсіювання?

1. Якщо за характеристику центральної тенденції вибрано медіану, то слід вибрати міжквартильний розмах.
2. Якщо за характеристику центральної тенденції вибрано вибіркове середнє, то слід вибрати дисперсію та середньоквадратичне відхилення.



Коефіцієнт асиметрії

Коефіцієнт асиметрії (skewness) обчислюється як:

$$a_s = \frac{m_3}{\sigma_B^3} \quad m_3 = \frac{1}{n} \sum n_i (x_i - \bar{x})^3$$

Якщо коефіцієнт асиметрії **додатний**:

- Розподіл має довгий правий хвіст
- Вибіркове середнє більше за медіану
- Медіана більша за моду

Якщо коефіцієнт асиметрії **від'ємний**:

- Розподіл має довгий лівий хвіст
- Вибіркове середнє менше за медіану
- Медіана менша за моду



Коефіцієнт ексцесу

Вибірковим коефіцієнтом ексцесу (kurtosis) називається величина

$$e_k = \frac{m_4}{\sigma_B^4} - 3 \qquad m_4 = \frac{1}{n} \sum n_i (x_i - \bar{x})^4$$

Якщо коефіцієнт ексцесу < 3 , то пік розподілу пологіший, ніж у нормального – розподіл **плосковершинний**.

Якщо коефіцієнт ексцесу > 3 , то пік розподілу крутіший, ніж у нормального – розподіл **гостровершинний**.



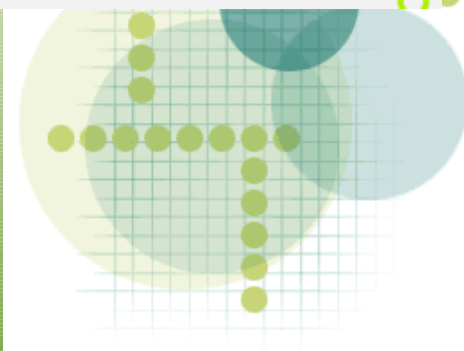
ІНТЕРПРЕТАЦІЯ РЕЗУЛЬТАТІВ парадокс Сімпсона

Факультет А

	Подало заяв	Прийнято	Відсоток прийнятих
Чоловіки	900	450	50%
Жінки	100	80	80%

Факультет Б

Чоловіки	100	10	10%
Жінки	900	180	20%



парадокс Сімпсона

Факультет А

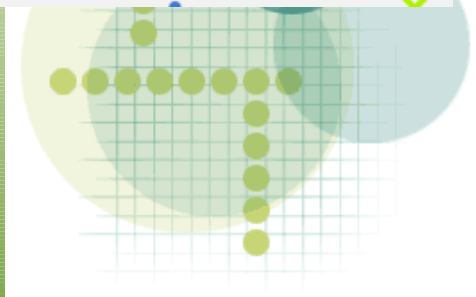
	Подало заяв	Прийнято	Відсоток прийнятих
Чоловіки	900	450	50%
Жінки	100	80	80%

Факультет Б

Чоловіки	100	10	10%
Жінки	900	180	20%

Обидва

Чоловіки	1000	460	46%
Жінки	1000	260	26%



Методика та філософія викладання

Теоретичний матеріал подається у вигляді
презентацій та викладень на дошці

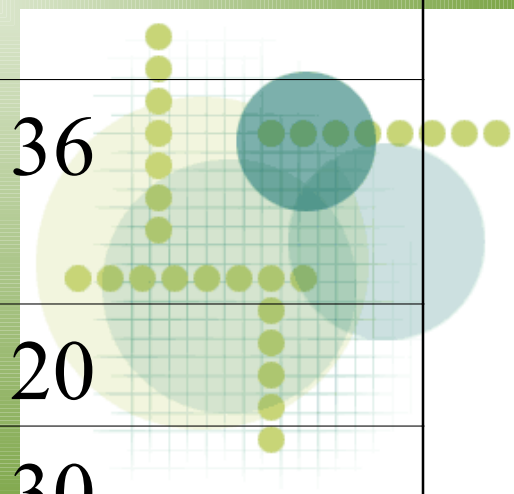
Засвоєння теоретичного матеріалу
відбувається за рахунок участі у
дискусіях, тестуванні, колоквіуму

Практичні навички формуються при
виконанні індивідуальних завдань



Умови визначення навчального рейтингу

№	Вид роботи	Кількість	Макс. кількість балів
1	Дискусія (тест) на лекціях	7	14
2	Участь у проекті	3	36
3	колоквіум	1	20
4	Ісит	1	30



Дякую за увагу!

