

# Лекція 2

## Вивідна статистика.

Збіжність емпіричних  
характеристик до теоретичних.

•

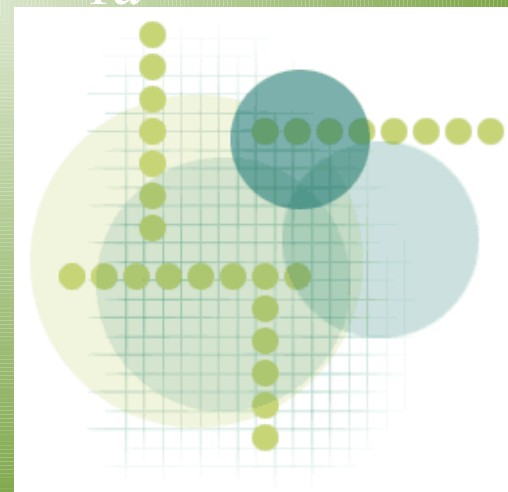
К.ф.-м.н. Щестюк Н.Ю.



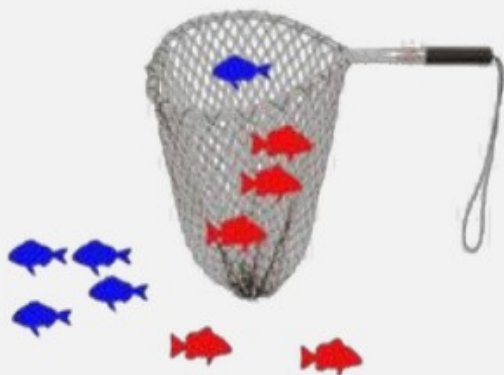
**Описова статистика** цікавиться виключно властивостями спостережуваних даних.

**Вивідна статистика** дозволяє робити висновки про генеральну сукупність на основі вибірки. Впевненість у цих висновках можна представити чисельно.

Розуміння термінів "генеральна сукупність" та "вибірка" є надзвичайно важливим для розуміння вивідної статистики.



# Вибірковий метод. 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)



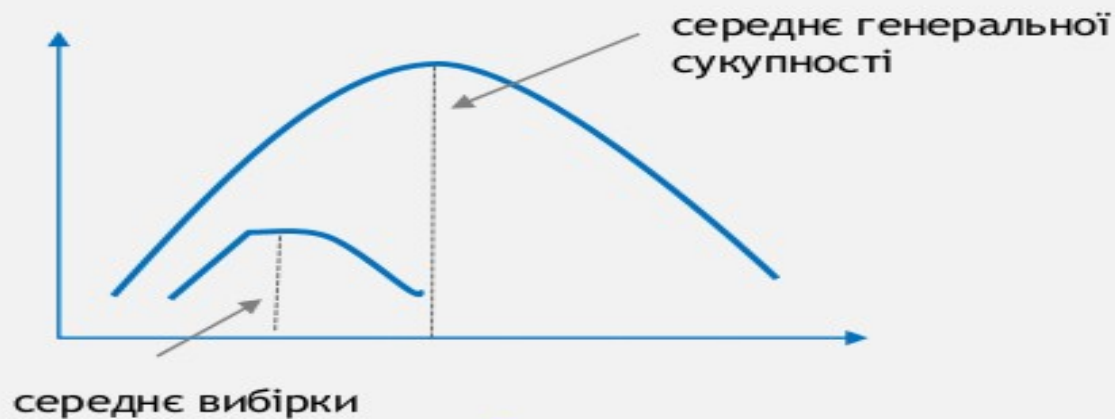
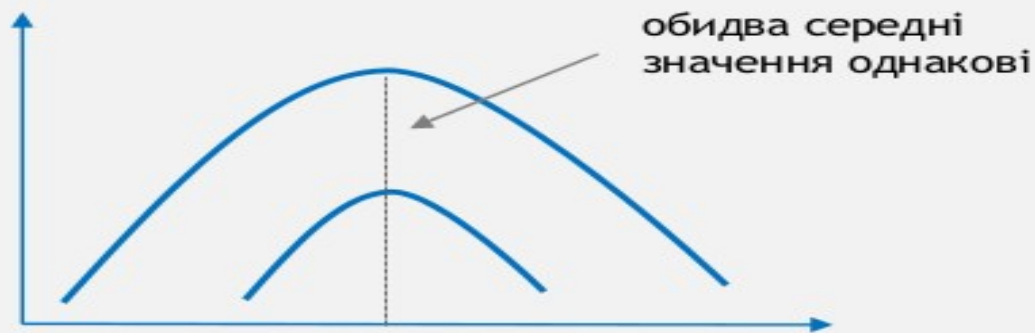
## Систематичний вибір

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

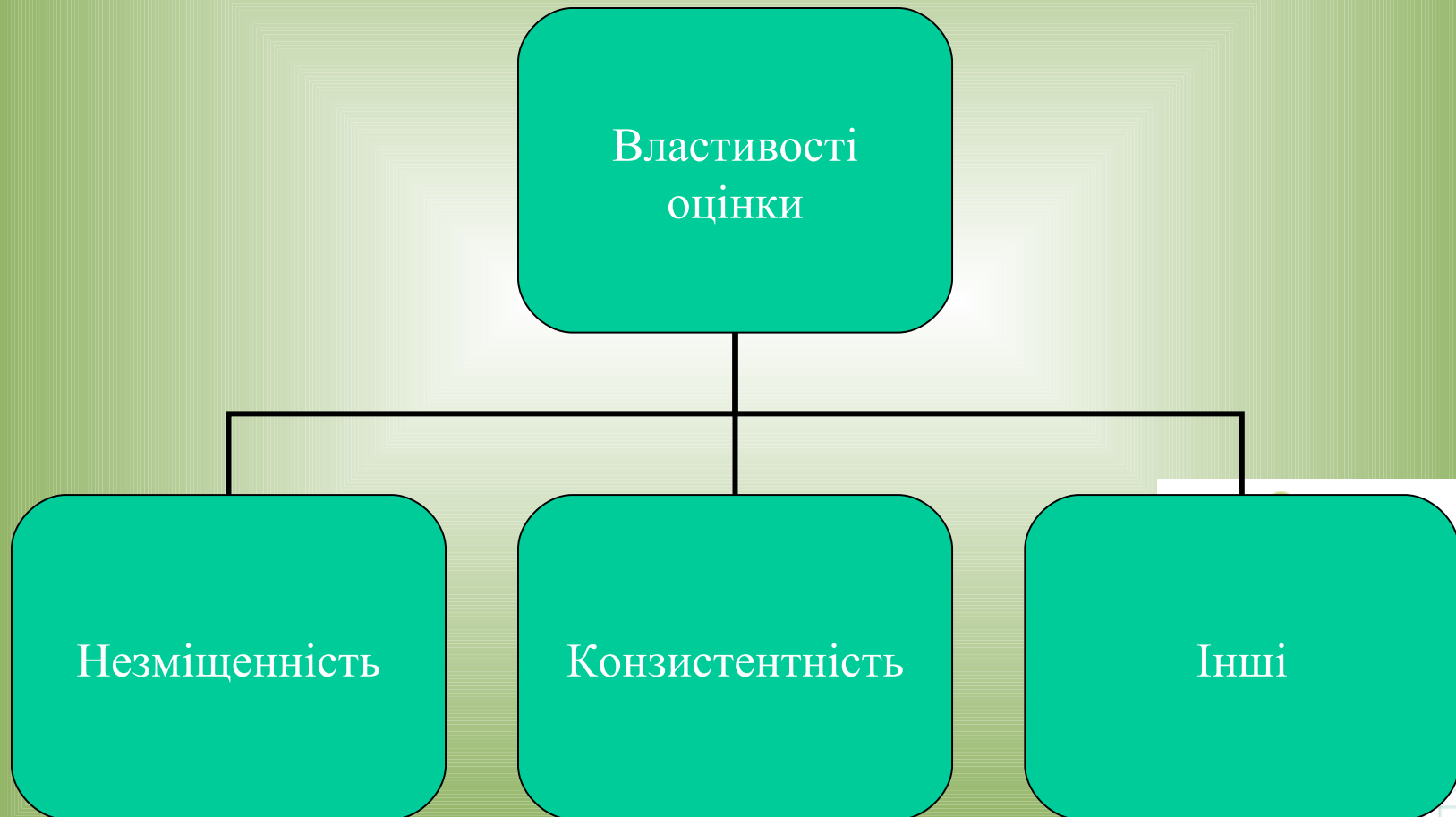
# Емпіричні характеристики вибірки: середнє, дисперсія, емпірична функція розподілу

	вибірка	генеральна сукупність
розмір	$n$	$N$
середнє значення	$\bar{x} = \frac{\sum x}{n}$	$\mu = \frac{\sum x}{n}$
дисперсія	$s^2 = \frac{\sum (x - \bar{x})^2}{n - 1}$	$\sigma^2 = \frac{\sum (x - \bar{\mu})^2}{N}$
середньоквадратичне відхилення	$s = \sqrt{s^2}$	$\sigma = \sqrt{\sigma^2}$
пропорція	$\bar{p} = \frac{n \text{ успіхів}}{n \text{ випробувань}}$	$p = \frac{N \text{ успіхів}}{N \text{ випробувань}}$

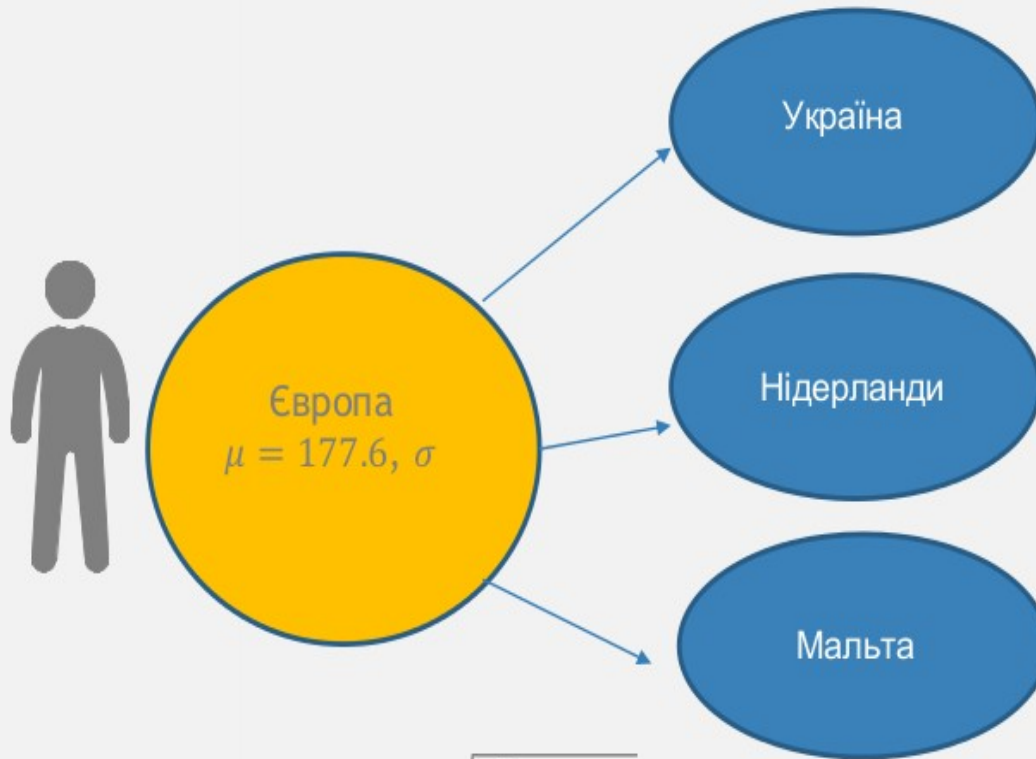
# Середнє вибірки і середнє генеральної сукупності (похибка репрезентативності)



# Властивості якості емпіричних оцінок



# Середнє сукупності вибірок



$$\mu = \frac{x_1 + x_2 + \dots + x_N}{N} \quad \sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N}}$$

вибірка 1  
 $\bar{x} = 176.5, s$

вибірка 2  
 $\bar{x} = 183.8, s$

вибірка 3  
 $\bar{x} = 169.9, s$

середнє ( $\bar{x}$ )  $\approx \mu$   
 $sd(\bar{x}) < \sigma$



# ЗВЧ у формі Чебишова та Бернуллі

$$P\left(\left|\frac{\sum_{j=1}^n z_j}{n} - \frac{\sum_{j=1}^n Ez_j}{n}\right| < \varepsilon\right) \geq 1 - \frac{D(z_j)}{n \varepsilon^2} \rightarrow 1$$

$$P\left\{\left|\frac{m}{n} - p\right| > \varepsilon\right\} \leq \frac{pq}{n \varepsilon^2}$$



# Припущення щодо вибірки

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

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

$$Ex_j = \mu$$



# Незміщенність та конзистентність для вибіркового середнього

$$1) \quad E(\bar{x}) = E\left(\frac{\sum_{j=1}^n x_j}{n}\right) = \frac{1}{n} \sum_{j=1}^n E x_j = \frac{1}{n} \sum_{j=1}^n \mu = \mu$$

$$2) \quad P\left(\left|\frac{\sum_{j=1}^n x_j}{n} - \frac{\sum_{j=1}^n E x_j}{n}\right| < \varepsilon\right) \geq 1 - \frac{D(x_j)}{n \varepsilon^2} \longrightarrow 1$$



# Незміщенність (зміщеність) для вибіркової дисперсії

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

$$\begin{aligned} 1) \quad E(s^2) &= E(\overline{x^2} - (\bar{x})^2) = E(\overline{x^2}) - E(\bar{x})^2 = \\ &= E(\overline{x^2}) - (D(\bar{x}) + (E\bar{x})^2) = \end{aligned}$$

$$= E(x_j^2) - \left( D\left( \frac{\sum x_j}{n} \right) - (Ex)^2 \right) =$$

$$= D(x) - \frac{n D(x)}{n^2} = D(x) \left( 1 - \frac{1}{n} \right) = \frac{n-1}{n} D(x)$$



# Конзистентність для вибіркової дисперсії

Нехай  $z_j = s_j^2 = x_j^2 - (\bar{x})^2$  ЗБЧ:  $P\left(\left|\frac{\sum_{j=1}^n z_j}{n} - \frac{\sum_{j=1}^n Ez_j}{n}\right| < \varepsilon\right) \geq 1 - \frac{D(z_j)}{n \varepsilon^2} \rightarrow 1$

$$\frac{\sum_{j=1}^n z_j}{n} = \frac{\sum_{j=1}^n (x_j^2 - (\bar{x})^2)}{n} = \frac{\sum_{j=1}^n x_j^2}{n} - \frac{\sum_{j=1}^n (\bar{x})^2}{n} = \overline{x^2} - (\bar{x})^2$$

$$\frac{\sum_{j=1}^n Ez_j}{n} = \frac{\sum_{j=1}^n Es_j^2}{n} = \frac{\sum_{j=1}^n E(x_j^2 - (\bar{x})^2)}{n} = \frac{nE(x_j^2 - (\bar{x})^2)}{n} = \frac{n-1}{n} D(x_j)$$



# Незміщенність для емпіричної функції розподілу

$$F_n^*(y) = \frac{1}{n} \sum_{i=1}^n I(x_i < y)$$

$$\begin{aligned} 1) \quad E(F_n^*(y)) &= \frac{\sum_{j=1}^n E(x_j < y)}{n} = \frac{1}{n} n E(x_j < y) = \\ &= 1P(x_j < y) + 0P(x_j \geq y) = P(x_j < y) = F(y) \end{aligned}$$



# Конзистентність для емпіричної функції розподілу

$$F_n^*(y) = \frac{1}{n} \sum_{i=1}^n I(x_i < y)$$

Нехай  $z_j = I(x_j < y)$

$$\text{ЗБЧ: } P \left( \left| \frac{\sum_{j=1}^n z_j}{n} - \frac{\sum_{j=1}^n E z_j}{n} \right| < \varepsilon \right) \geq 1 - \frac{D(z_j)}{n \varepsilon^2} \rightarrow 1$$

$$\frac{\sum_{j=1}^n z_j}{n} = \frac{\sum_{j=1}^n I(x_j < y)}{n} = F^*(y)$$

$$\frac{\sum_{j=1}^n E z_j}{n} = \frac{\sum_{j=1}^n E I(x_j < y)}{n} = \frac{n E I(x_j < y)}{n} = P(x_j < y) = F(y)$$



# ЦГТ і її застосування

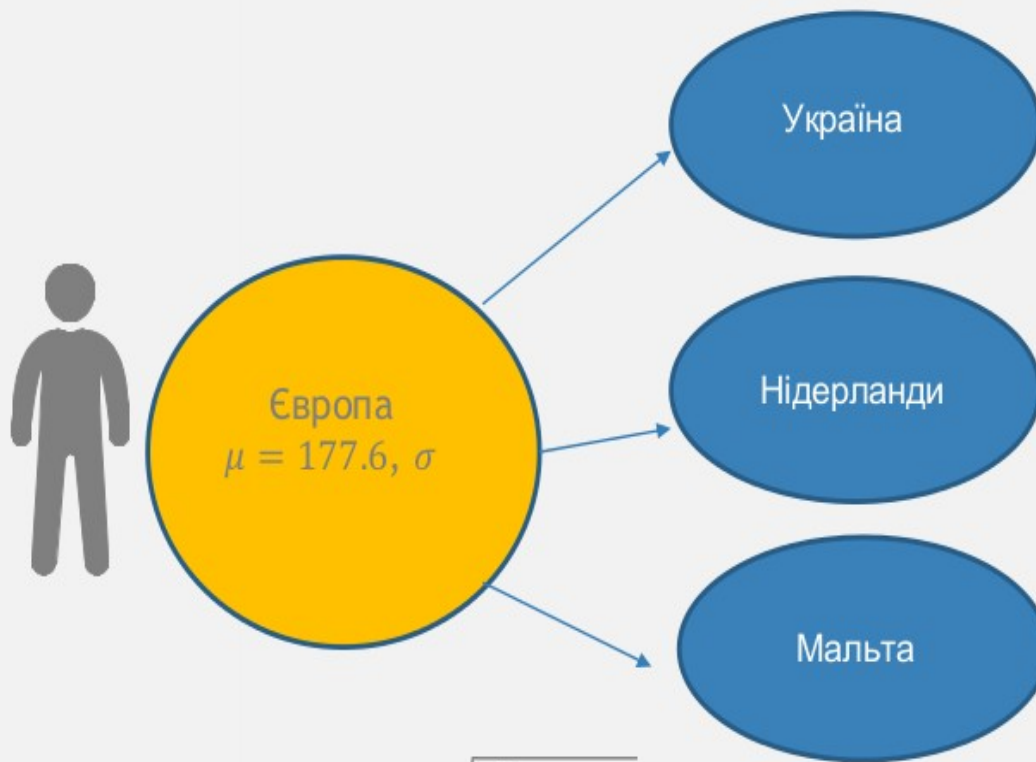
$$1) P\left(\left|\frac{\sum_{j=1}^n z_j}{n} - \frac{\sum_{j=1}^n Ez_j}{n}\right| < \varepsilon\right) \longrightarrow \Phi\left(\frac{\bar{z} - \mu}{\sigma}\right)$$

$$2) P\left(\left|\frac{m}{n} - p\right| < \Delta\right) = 2\Phi\left(\Delta\sqrt{\frac{n}{pq}}\right)$$





# Середнє сукупності вибірок



$$\mu = \frac{x_1 + x_2 + \dots + x_N}{N} \quad \sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N}}$$

вибірка 1  
 $\bar{x} = 176.5, s$

вибірка 2  
 $\bar{x} = 183.8, s$

вибірка 3  
 $\bar{x} = 169.9, s$

середнє ( $\bar{x}$ )  $\approx \mu$   
 $sd(\bar{x}) < \sigma$

# Довірча ймовірність (confidence level) та розмір вибірки для оцінювання середнього

$$P(|\bar{x} - \mu| < \Delta) = \gamma - \text{довірча ймовірність (рівень довіри)}$$

$\Delta$  - гранична помилка (межа похибки)

$$D(\bar{x}) = \frac{1}{n^2} (D(x_1) + \dots + D(x_n)) = \frac{s^2}{n} - \text{дисперсія середнього}$$

$$\sigma = \frac{s}{\sqrt{n}} - \text{стандартна похибка середнього}$$

$$t = \frac{\bar{x} - \mu}{\sigma} = \frac{\Delta \sqrt{n}}{s} - \text{квантиль нормального розподілу}$$

порядку  $P$

$$n = \left( \frac{t s}{\Delta} \right)^2 - \text{розмір вибірки}$$



# Довірча ймовірність та розмір вибірки для оцінювання пропорції (успіху)

$$P\left(\left|\frac{m}{n} - p\right| < \Delta\right) = \gamma - \text{довірча ймовірність (рівень довіри)}$$

$\Delta$  - гранична помилка (межа похибки), виражена в частках одиниці

$$t = \Delta \sqrt{\frac{n}{pq}} - \text{квантиль нормального розподілу}$$

порядку  $P$

$$n = \left(\frac{t \sqrt{pq}}{\Delta}\right)^2 - \text{розмір вибірки}$$



# Приклад 1

- Якого розміру має бути вибірка, щоб оцінити середній об'єм випитого за місяць пива для людей певного регіону. Випадкова похибка має не перевищувати 0,1 л з довірчою ймовірністю 95%. Відомо, що  $s=2$ л

$t = 1,96$  - квантиль нормального розподілу

порядку  $P = 0.95$

$$n = \left( \frac{t s}{\Delta} \right)^2 = \left( \frac{1.96 * 2}{0.1} \right)^2 \approx 1600 - \text{розмір вибірки}$$



# Приклад 2

- Якого розміру має бути вибірка, щоб оцінити частку громадян України, які мають намір прийти на вибори. Випадкова похибка має не перевищувати 2%. З попередніх досліджень для частки осіб, які мають намір прийти має не перевищувати 0,9 з довірчою ймовірністю 95%.

$t = 1,96$  - квантиль нормального розподілу

порядку  $P = 0.95$

$$n = \left( \frac{t \sqrt{pq}}{\Delta} \right)^2 = \left( \frac{1.96 \sqrt{0.9 * 0.1}}{0.02} \right)^2 \approx 900 - \text{розмір вибірки}$$

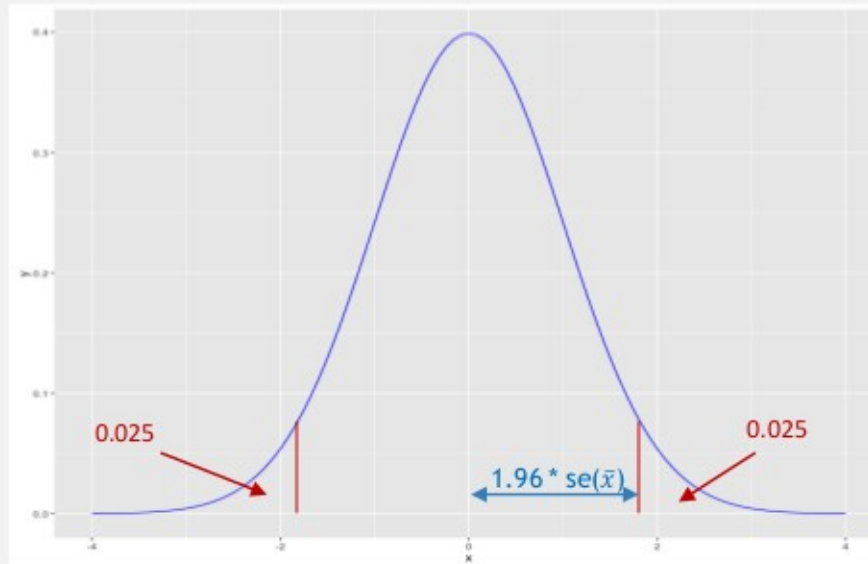


# довірчий інтервал для середнього значення

$$\bar{x} \pm Z_{95\%} * se(\bar{x}), \text{ де } se(\bar{x}) = \frac{\sigma}{\sqrt{n}}$$

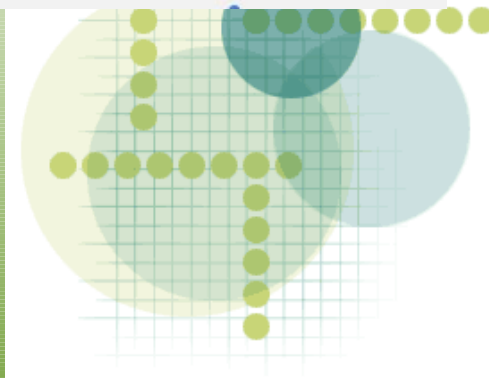


$$\bar{x} \pm 1.96 * se(\bar{x}), \text{ де } se(\bar{x}) = \frac{\sigma}{\sqrt{n}}$$



$$\mu_{\bar{x}} = \mu$$
$$se(\bar{x}) = \frac{\sigma}{\sqrt{n}}$$

$1.96 * se(\bar{x})$  - межа по  
для рівня довіри 95%

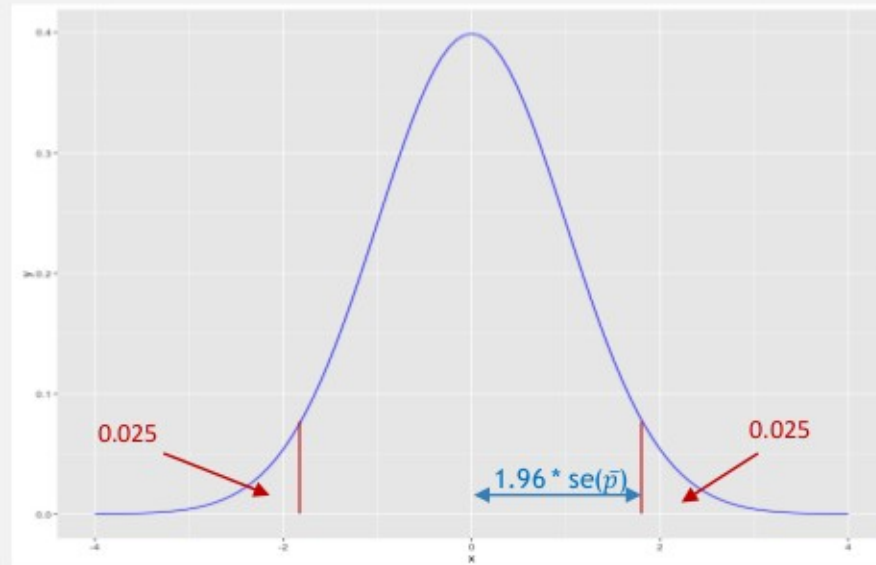


## довірчий інтервал для пропорції

$$p \pm Z_{95\%} * se(\bar{p}), \text{ де } se(\bar{p}) = \sqrt{\frac{p(1-p)}{n}}$$

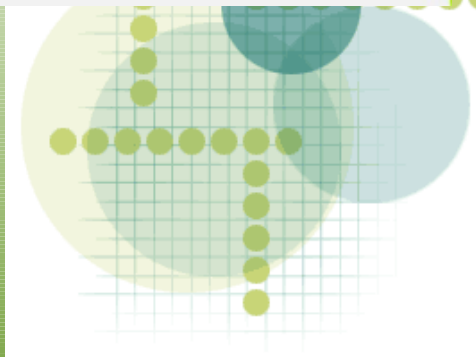


$$p \pm 1.96 * se(\bar{p}), \text{ де } se(\bar{p}) = \sqrt{\frac{p(1-p)}{n}}$$



$$\mu_{\bar{p}} = p, \\ se(\bar{p}) = \sqrt{\frac{p(1-p)}{n}}$$

$1.96 * se(\bar{p})$  - межа пох  
для рівня довіри 95%



# Приклад

## довірчий інтервал для пропорції

Серед 935 випадковим чином обраних респондентів на питання “чи вірите ви в існування розумного життя на інших планетах?” ствердно відповіли 60%

$$\bar{p} = 0.6, n = 935$$

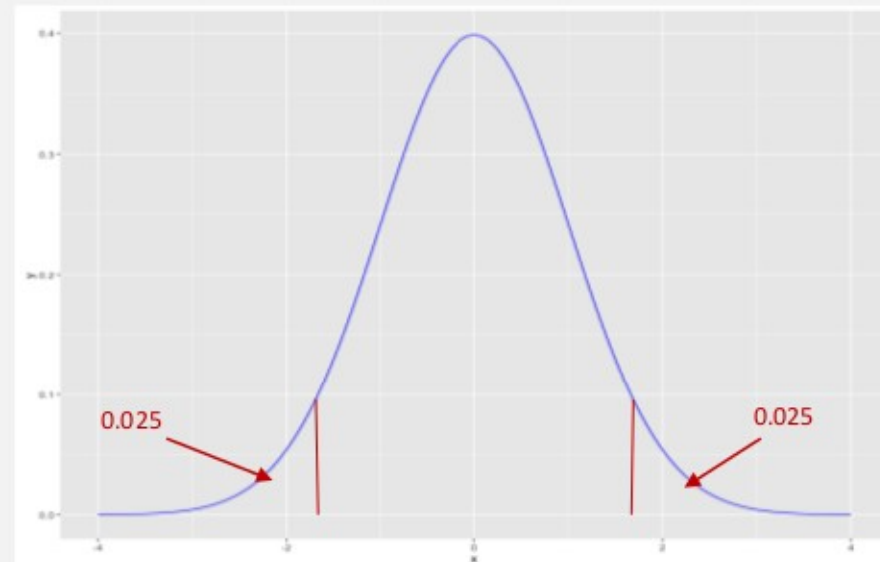
$$p \pm Z_{95\%} * se(\bar{p}), \text{ де } se(\bar{p}) = \sqrt{\frac{p(1-p)}{n}}$$

$$0.6 \pm 1.96 * se(\bar{p}), \text{ де } se(\bar{p}) = \sqrt{\frac{0.6(1-0.6)}{935}} = 0.016$$

$$0.6 \pm 1.96 * 0.016$$

$$0.6 \pm 0.03136$$

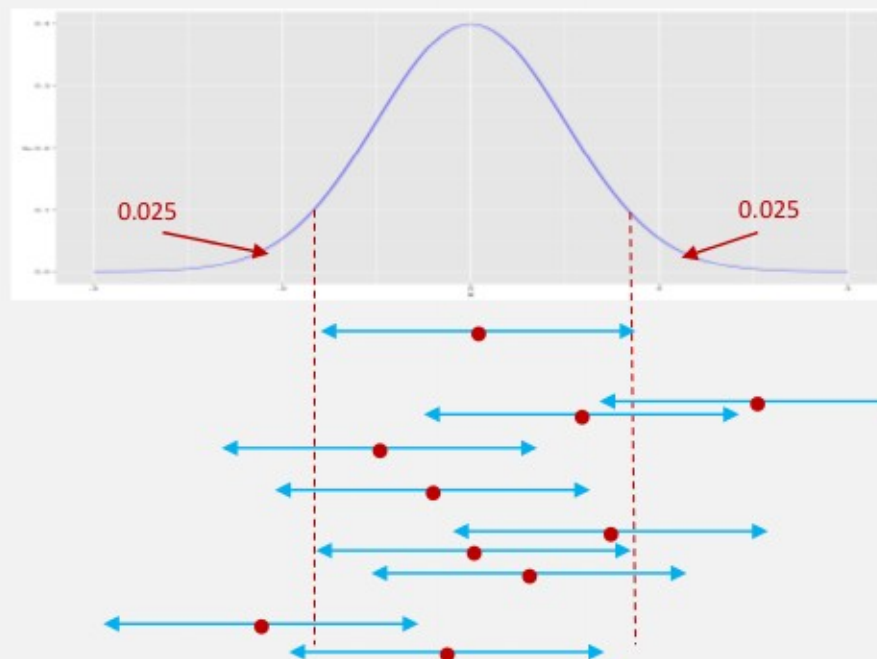
$$[0.56864, 0.63136]$$





# Довірчий інтервал

рівень довіри



Дякую за увагу!

