

Київський національний університет  
імені Тараса Шевченка

Звіт

до лабораторної роботи 6,  
дисципліни: “Технології аналізу та візуалізації даних”,  
студента: Євчика Олексія,  
групи Інформатика, Магістри

**Тема:** Візуалізація розподілених даних з використанням гістограм, ядерної оцінки щільності (KDE) та діаграм розмахів (Boxplot) (R).

**Мета роботи:**

1. Засвоїти принципи відображення розподілу даних за допомогою гістограм, KDE та boxplot.
2. Навчитися налаштовувати параметри бінінгу, ядра та «вусів» для коректної інтерпретації.
3. Оволодіти базовими прийомами попередньої обробки, виявлення викидів і побудови порівняльних візуалізацій у R.

Ознайомитися з методами візуалізації розподілених даних у R, навчитися застосовувати гістограми, ядерну оцінку щільності та діаграми розмахів, дослідити вплив параметрів (*binwidth*, *bandwidth*, *adjust*, правила визначення вусів) на форму графіків та інтерпретацію даних.

**Теоретична частина**

**1. Гістограма**

**Ідея.** Гістограма є методом візуалізації розподілу числової змінної шляхом дискретизації осі значень на інтервали (біни) та підрахунку кількості спостережень або оцінки щільності у кожному біні.

**Ключові параметри.**

- *binwidth* — ширина біну;
- *breaks / number of bins* — кількість та розташування інтервалів;
- *limits* — межі відображення;
- *count vs density* — шкала частот або нормована щільність.

**Зауваги.** Занадто широкий бін може приховати багатомодальність, тоді як надто вузький — створює шум та “зубчастість” графіка.

## 2. Ядерна оцінка щільності (KDE)

**Ідея.** KDE — непараметричний метод оцінки щільності ймовірності, що базується на згладженні спостережень за допомогою ядра (частіше Gaussian).

**Параметри згладження.**

- *bandwidth* (ширина вікна) визначає ступінь згладження;
- у *ggplot2* масштабування *bandwidth* задається через аргумент *adjust*.

**Компроміс *bandwidth*.**

- більший *bw* → гладка крива, але втрата локальних деталей;
- менший *bw* → нерівний контур і чутливість до шуму.

## 3. Діаграма розмахів (Boxplot)

**Складові.**

- медіана (лінія усередині коробки),
- межі кватилів Q1–Q3 (коробка),
- “вуса” (зазвичай до 1.5 IQR або іншого правила),
- точки-викиди (спостереження за межами вусів).

**Призначення.** Порівняння розподілів між групами, робастне відображення центральної тенденції та розкиду.

**Варіації.**

- boxplot із накладеними точками (*stripplot*, *jitter*);

- violin plot для візуалізації форми щільності.

#### 4. Порівняння методів

- **Гістограма:** дає загальну форму розподілу, але залежить від вибору бінів.
  - **KDE:** гладка безбінова оцінка, чутлива до bandwidth.
  - **Boxplot:** відображає позиційні характеристики та викиди, але не показує детальну форму розподілу.
- 

#### 5. Filled Contour (залиті контури)

##### 5.1. Ідея методу

Filled contour — це спосіб візуалізації *двовимірної поверхні* ( $z = f(x, y)$ ) у вигляді рівневих ліній із заливкою між рівнями. Метод дозволяє відобразити структуру просторових даних, теплових карт, топографії та поверхонь щільності.

Графік формує набір рівнів (levels), між якими простір зафарбовується окремими кольорами, що полегшує сприйняття градієнтів і переходів значень.

##### 5.2. Вибір палітри

Колірна палітра є критичною для коректної інтерпретації таких карт.

Основні вимоги до палітри для filled contour:

- **послідовність (perceptual uniformity)** — однаковий крок кольору відповідає однаковій зміні значення;

- **дружність до дальтонізму** — кольори повинні бути розрізняваними навіть у випадках порушення колірного зору;
- **монотонність яскравості** — щоб забезпечити правильне сприйняття “високих” і “низьких” значень;
- **читабельність у друці та в градаціях сірого.**

Палітри сімейства **viridis** (`viridisLite::viridis`) повністю задовольняють ці вимоги. Вони були спеціально розроблені для наукової візуалізації та є стандартом де-факто у багатьох бібліотеках.

### 3. Дані та попередня обробка

У даній роботі були використані випадково згенеровані дані.

Попередня обробка включала:

- **фільтрацію** (видалення некоректних або явно помилкових записів),
- **масштабування** (за потреби — стандартне масштабування або нормалізація),
- **обробку викидів**: попередній аналіз за допомогою `boxplot`; рішення — видалення / заміна / залишення без змін, залежно від мети аналізу.
- **приведення типів даних** (числові, фактори, дати).

## 4. Виявлення викидів за IQR-правилом

### 4.1. Мета

Мета індивідуального завдання — ідентифікувати викиди у вибраній числовій змінній з використанням IQR-правила (Interquartile Range Rule), визначити порогові значення та проаналізувати природу знайдених викидів.

---

## 4.2. Теоретичні відомості

IQR-правило ґрунтується на інтерквартильному розмаху, що визначається як:

$$IQR = Q3 - Q1,$$

де:

- $Q1$  — перший квартиль (25-й перцентиль),
- $Q3$  — третій квартиль (75-й перцентиль).

Порогові значення для виявлення викидів визначаються як:

$$\text{Lower Bound} = Q1 - 1.5 \cdot IQR,$$

$$\text{Upper Bound} = Q3 + 1.5 \cdot IQR.$$

Будь-які спостереження, що лежать нижче нижнього порогу або вище верхнього, вважаються потенційними викидами.

Переваги методу: простота, робастність до шуму, нечутливість до великих амплітудних викидів.

Недоліки: IQR-правило добре працює для одномодальних симетричних розподілів, але може давати багато "штучних" викидів у сильно асиметричних розподілах.

## 4.3. Кроки виконання

### Обчислення квартилів:

Отримано значення  $Q1$  та  $Q3$  для обраної змінної.

## Обчислення IQR:

$$IQR = Q3 - Q1$$

### 1. Визначення порогів:

- нижній поріг =  $Q1 - 1.5 \cdot IQR$
- верхній поріг =  $Q3 + 1.5 \cdot IQR$

### 2. Підрахунок кількості викидів:

Визначено кількість спостережень нижче та вище порогів.

### 3. Формування таблиці результатів та коротка інтерпретація.

## 6.4. Очікуваний результат

Нижче наведено загальний формат таблиці (числа заповнюються після обчислень у R).

Код:

```
## -----  
## 1. Підготовка середовища  
## -----  
packages <- c("ggplot2","dplyr","tidyr","readr","scales","ggpubr")  
to_install <- setdiff(packages, rownames(installed.packages()))  
if(length(to_install)) install.packages(to_install)  
lapply(packages, library, character.only = TRUE)  
set.seed(123)
```

```

11 ▾ ## -----
12 ## 2. Дані-примірки
13 ▾ ## -----
14
15 # 1) Вбудований набір даних
16 data("faithful") # колонки: eruptions, waiting
17
18 # 2) Синтетика
19 x1 <- rnorm(400, mean = 0, sd = 1)
20 x2 <- rnorm(300, mean = 3, sd = 0.7)
21 df_mix <- data.frame(x = c(x1, x2))
22
23
24
25 ## 3. Гістограми: ширина біну, count vs density
26 ▾ ## -----
27
28 # Count
29 p_h1 <- ggplot(faithful, aes(x = eruptions)) +
30   geom_histogram(binwidth = 0.2, fill = "grey70", color = "grey30") +
31   labs(title = "Гістограма (лічильник)",
32        x = "Тривалість виверження (хв)", y = "Кількість") +
33   theme_minimal(base_size = 12)
34
35 # Density
36 p_h2 <- ggplot(faithful, aes(x = eruptions, y = after_stat(density))) +
37   geom_histogram(binwidth = 0.2, fill = "grey70", color = "grey30") +
38   labs(title = "Гістограма (щільність)",
39        x = "Тривалість (хв)", y = "Щільність") +
40   theme_minimal(12)
41
42 # Вплив binwidth - малий
43 p_bw_small <- ggplot(df_mix, aes(x)) +
44   geom_histogram(binwidth = 0.15, fill = "steelblue", color = "white") +
45   labs(title = "Малий binwidth (деталі + шум)", x = "x", y = "Count") +
46   theme_minimal(12)
47
48
49
50 # Вплив binwidth - великий
51 p_bw_large <- ggplot(df_mix, aes(x)) +
52   geom_histogram(binwidth = 0.6, fill = "steelblue", color = "white") +
53   labs(title = "Великий binwidth (згладження + ризик втрати мод)", x = "x", y = "Count") +
54   theme_minimal(12)

```



```

56 ## 4. KDE та накладання на гістограму
57 ## -----
58
59 p_kde <- ggplot(df_mix, aes(x)) +
60   geom_density(linewidth = 1) +
61   labs(title = "KDE (оцінка щільності)", x = "x", y = "Щільність") +
62   theme_minimal(12)
63
64 p_hist_kde <- ggplot(df_mix, aes(x, y = after_stat(density))) +
65   geom_histogram(binwidth = 0.3, fill = "grey80", color = "grey40") +
66   geom_density(linewidth = 1) +
67   labs(title = "Гістограма + KDE", x = "x", y = "Щільність") +
68   theme_minimal(12)
69
70 # Вплив adjust
71 p_kde_adj06 <- ggplot(df_mix, aes(x)) +
72   geom_density(adjust = 0.6, linewidth = 1) +
73   labs(title = "KDE: adjust = 0.6 (детальніше)", x = "x", y = "Щільність") +
74   theme_minimal(12)
75
76 p_kde_adj10 <- ggplot(df_mix, aes(x)) +
77   geom_density(adjust = 1.0, linewidth = 1) +
78   labs(title = "KDE: adjust = 1.0 (баланс)", x = "x", y = "Щільність") +
79   theme_minimal(12)
80

```

```

87 ## -----
88 ## 5. Boxplot + IQR-викиди
89 ## -----
90
91 p_box <- ggplot(iris, aes(x = Species, y = Sepal.Length)) +
92   geom_boxplot(outlier.colour = "red", width = 0.6) +
93   geom_jitter(width = 0.1, alpha = 0.4) +
94   labs(title = "Boxplot: Sepal.Length за видами",
95         x = "Вид", y = "Sepal.Length") +
96   theme_minimal(12)
97

```

```

97
98 ## --- Визначення викидів за IQR ---
99 x <- faithful$eruptions
100
101 Q1 <- quantile(x, 0.25)
102 Q3 <- quantile(x, 0.75)
103 IQRv <- IQR(x)
104
105 lower <- Q1 - 1.5 * IQRv
106 upper <- Q3 + 1.5 * IQRv
107
108 out_idx <- which(x < lower | x > upper)
109 out_values <- x[out_idx]
110
111 ## Таблиця порогів IQR
112 iqr_table <- data.frame(
113   Metric = c("Q1", "Q3", "IQR", "Нижній поріг", "Верхній поріг", "Кількість викидів"),
114   Value = c(Q1, Q3, IQRv, lower, upper, length(out_idx))
115 )
116

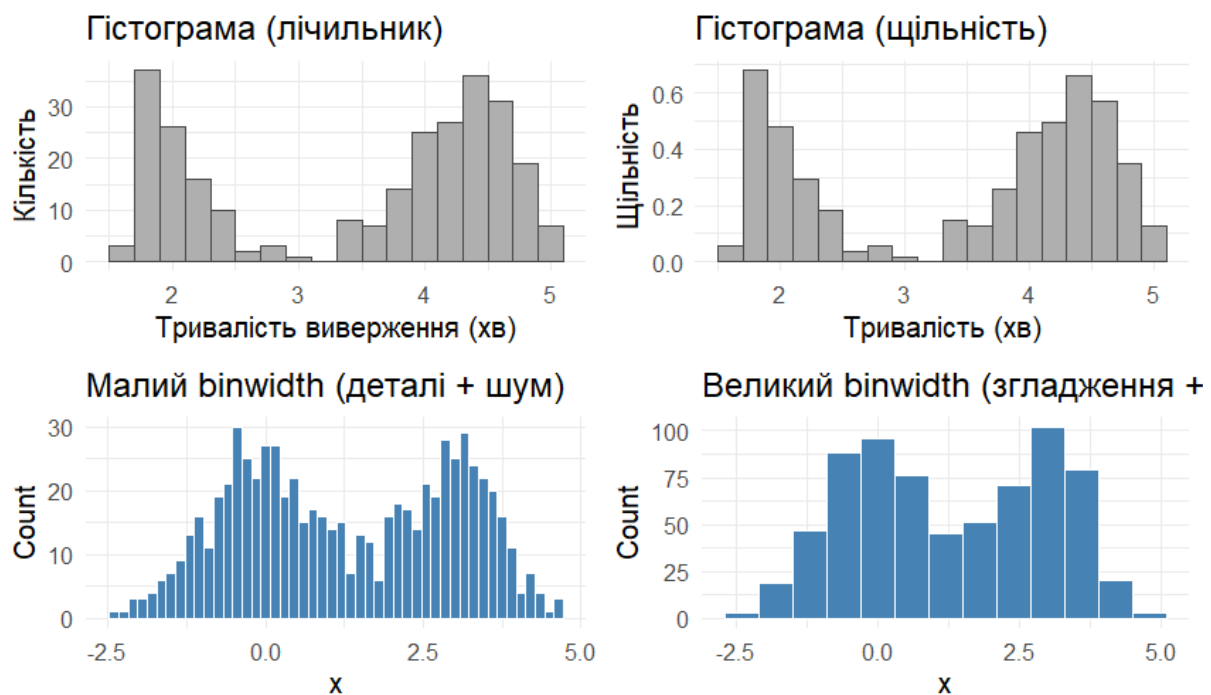
```

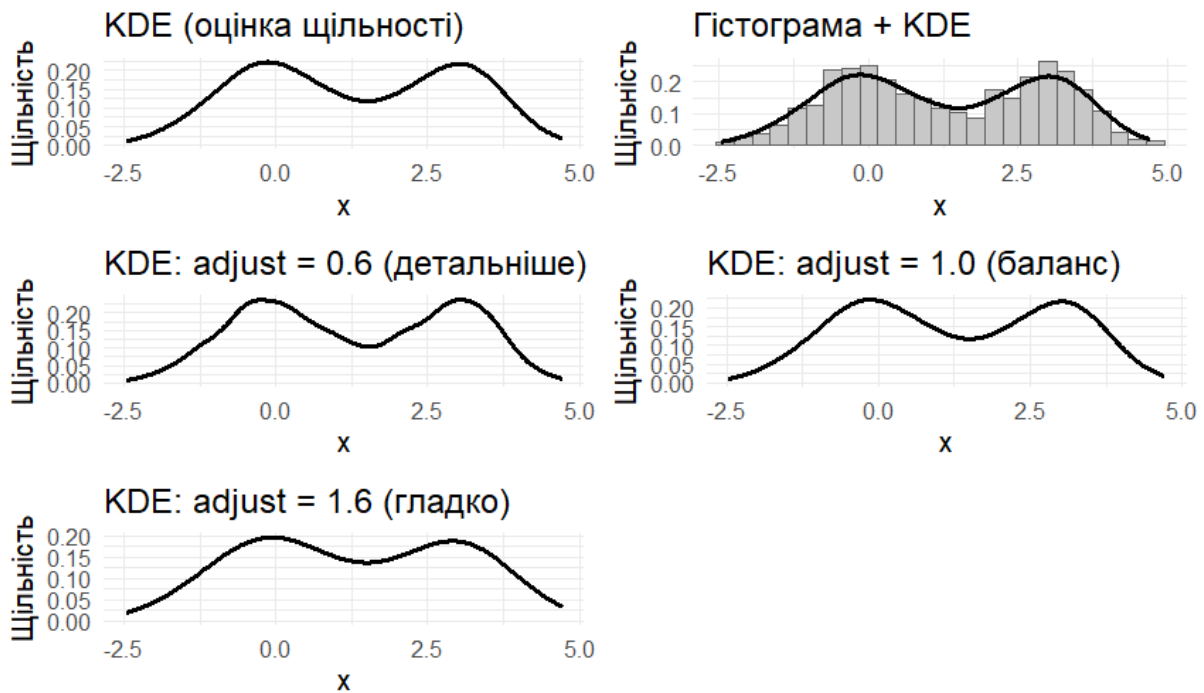
```

111 ## Таблиця порогів IQR
112 iqr_table <- data.frame(
113   Metric = c("Q1", "Q3", "IQR", "Нижній поріг", "Верхній поріг", "Кількість викидів"),
114   Value   = c(Q1, Q3, IQRv, lower, upper, length(out_idx))
115 )
116
117 ## -----
118 ## 6. Компоновка фігур
119 ## -----
120
121 ggpubr::ggarrange(p_h1, p_h2, p_bw_small, p_bw_large, ncol = 2, nrow = 2)
122 ggpubr::ggarrange(p_kde, p_hist_kde, p_kde_adj06, p_kde_adj10, p_kde_adj16,
123   ncol = 2, nrow = 3)
124
125 ## -----
126 ## 7. Збереження результатів
127 ## -----
128 ggsave("hist_density_overlay.png", p_hist_kde, width = 8, height = 5, dpi = 300)
129 ggsave("box_iris.png", p_box, width = 7, height = 5, dpi = 300)
130

```

Результат:





Контрольні запитання та відповіді:

1. Як вибір binwidth впливає на форму та інтерпретацію гістограми?

- Малий binwidth → гістограма стає «зашумленою», показує дрібні коливання, може створити хибні моди.
- Великий binwidth → занадто згладжує розподіл, приховує важливі структури (наприклад, мультимодальність).
- Тому binwidth визначає баланс між деталізацією та згладженістю.

2. У чому різниця між відображенням частоти (count) і щільності (density) на гістограмі?

- Count — показує абсолютну кількість спостережень у кожному біні.
- Density — масштабує площу гістограми до 1 (щоб можна було порівнювати з KDE та між вибірками різного розміру).

Density обов'язкова, якщо гістограму хочуть накладати на KDE.

3. Яку роль відіграє параметр `adjust` у KDE і як його обирати?

- `adjust` множить базову ширину ядра (`bandwidth`).
- Менше 1 → менше згладження, більше деталей, але більше шуму.
- Більше 1 → сильніше згладження, менше шуму, але можливе зникнення мод.
- Обирають за метою аналізу:
  - Для виявлення структур → 0.6–1.0
  - Для гладкої презентації → 1.2–1.6

4. Чому KDE чутлива до ширини ядра? Що таке компроміс «зміщення–дисперсія»?

- Малий `bandwidth` → низьке зміщення, висока дисперсія (занадто змінна крива).
- Великий `bandwidth` → високе зміщення, низька дисперсія (занадто згладжена крива).
- Компроміс зміщення–дисперсія означає, що не можна одночасно мати і низьке зміщення, і низьку дисперсію.

5. За яким правилом формуються «вуса» у `boxplot` за замовчуванням?

У класичному Tukey `boxplot`:

- Нижній вус =  $Q1 - 1.5 \times IQR$
- Верхній вус =  $Q3 + 1.5 \times IQR$
- Точки поза вусами вважаються потенційними викидами.

6. Чому `boxplot` вважається робастним до викидів?

Бо він використовує медіану та квартилі, а не середнє і стандартне відхилення.

Ці статистики стійкі: значення викидів не відсувають їх суттєво.

7. Коли доцільно комбінувати histogram + KDE на одному графіку?

- Коли потрібно одночасно показати:
  - емпіричний розподіл (гістограма)
  - гладке наближення розподілу (KDE)
- Особливо корисно при аналізі:
  - мультимодальності
  - порівняння вибірок
  - демонстрації різних bandwidth

8. Які підходи до виявлення та обробки викидів ви знаєте (IQR, z-score тощо)?

Методи виявлення:

- IQR-правило (Tukey)
- z-score:  $|z| > 3$
- robust z-score (MAD)
- діаграми: boxplot, scatter + LOF
- кластеризація (DBSCAN)
- моделі зважування (robust regression)

Методи обробки:

- вилучення
- winsorizing (обрізання)

- лог-трансформації
- заміна на медіану/квартиль
- окреме моделювання груп

#### 9. Які ризики неправильного масштабування осей/форматування підписів?

- Неправильне масштабування може створити візуальні викривлення, зокрема:
  - перебільшення різниць
  - приховання варіації
  - некоректне інтерпретування трендів
- Погані підписи призводять до:
  - неправильного розуміння одиниць
  - плутанини між density та count
  - хибних висновків у порівняннях

#### 10. Як забезпечити відтворюваність графіків у звіті?

- фіксувати `set.seed()`
- зберігати версії пакетів (`sessionInfo()`)
- зберігати всі дані/генерацію у файлі
- зберігати всі параметри графіків (`binwidth`, `adjust`, `limits`)
- використовувати RMarkdown або Quarto, де виконання коду повністю контрольоване
- уникати випадкових аспектів (джиттер з `width`, `alpha` без `seed`)