

# ЛІНІЙНА РЕГРЕСІЯ

*Лекція 3*

# ПЛАН

1. Машинне навчання
2. Регресія
3. Лінійна регресія



# МАШИННЕ НАВЧАННЯ

Машинне навчання (МН, Machine Learning, ML) - великий підрозділ штучного інтелекту, що вивчає методи побудови алгоритмів, здатних навчатися.

Машинне навчання - це «розділ ШІ, який досліджує методи, що дозволяють покращувати свої характеристики на основі отриманого досвіду».

**Машинне навчання** – це підрозділ штучного інтелекту, який розглядає побудову алгоритмів, які можуть навчатися на наявних даних.

Задача МН виглядає так: уявімо собі, що в нас є певний набір об'єктів прикладів і певний набір міток, тобто, реакцій, відповідей. Між прикладами-спостереженнями і відповідями є певна прихована залежність.

Задача МН – знайти цю приховану залежність для прогнозування відповідей на основі нових даних.

У найзагальнішому випадку розрізняють два типи машинного навчання: **навчання по прецедентах**, або **індуктивне навчання**, і **дедуктивне навчання**.

*Індуктивне* навчання знайоме кожному, адже воно полягає у спостереженні за світом та побудові певних моделей, які пояснюють причини тих чи інших явищ. Потім такі моделі неодноразово перевіряються, певні з них «виживають» і використовуються, покращуються. А деякі моделі згодом цілком відкидаються.

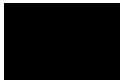
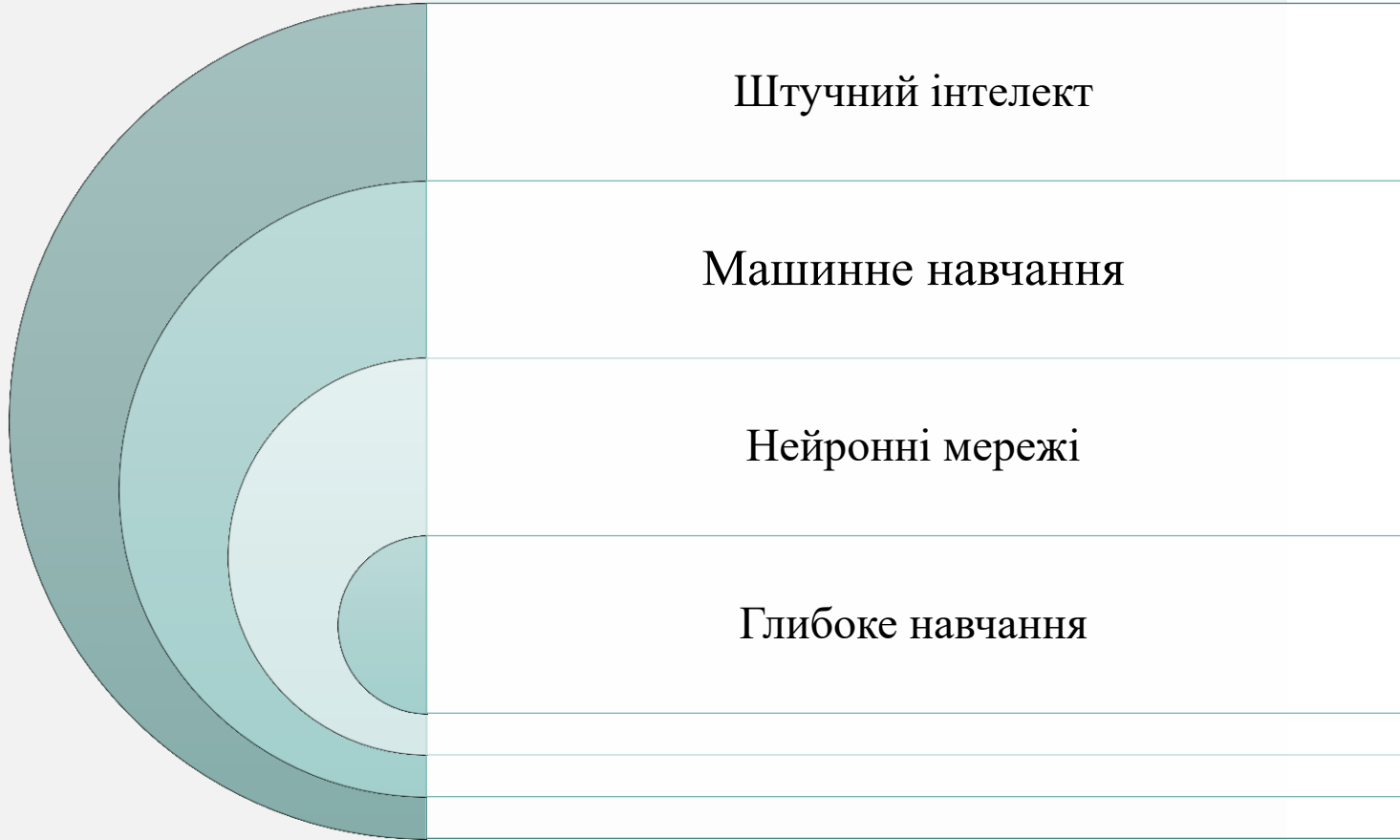
*Дедуктивне* навчання подібне до математики в школі, коли учню дають готові формули і розказують, як застосовувати їх на практиці

Навчання по прецедентах, в свою чергу, поділяють на три основних типи:

*контрольоване навчання*, або навчання з учителем (supervised learning),

*неконтрольоване навчання* (unsupervised learning), або навчання без учителя,

*навчання з підкріпленням* (reinforcement learning).





Вираховувати шахрайство з банківськими картками



Моделювати ризики для інвестицій або кредитування



Робити фінансові прогнози



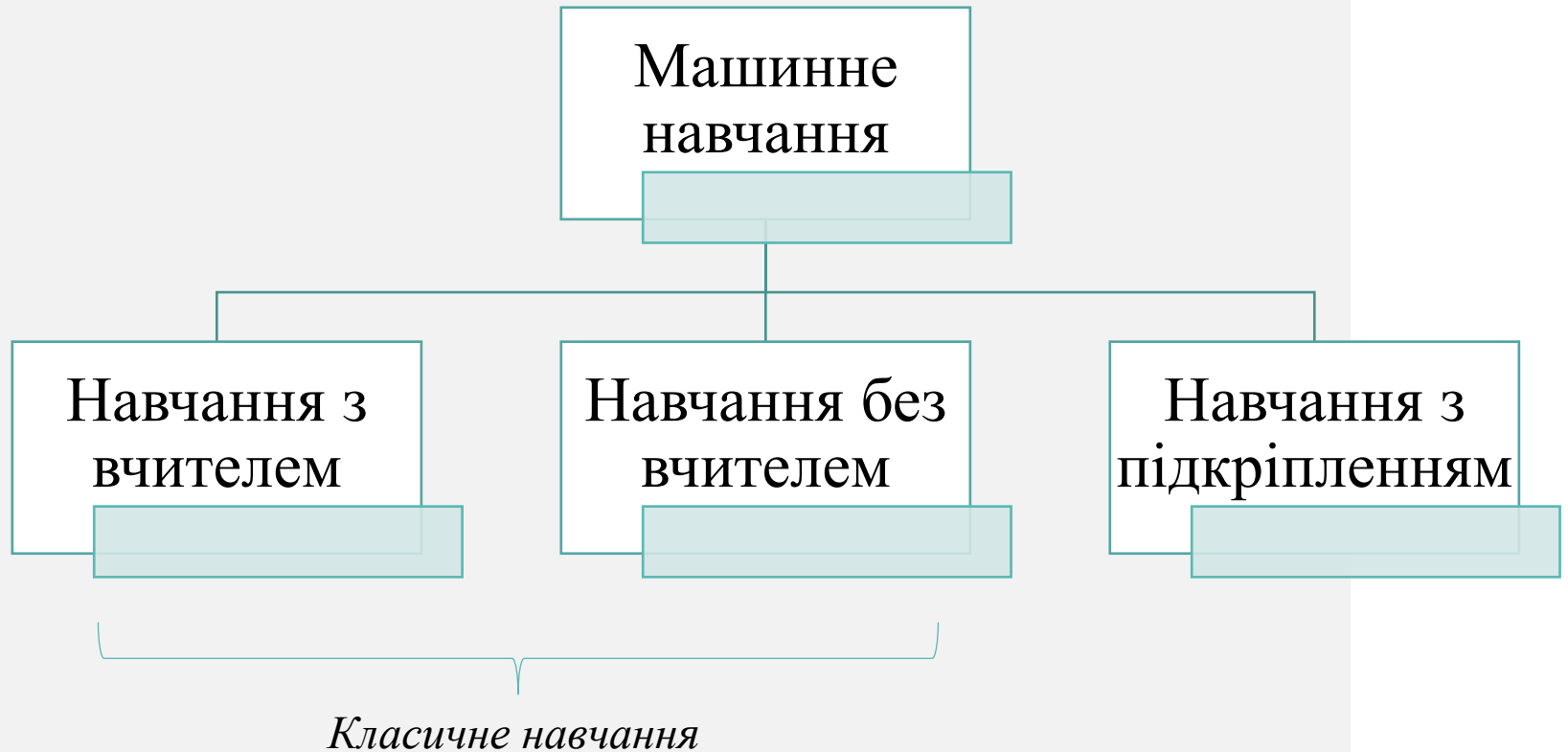
Сегментувати клієнтів



Створювати системи рекомендацій



# Основні типи машинного навчання



## Навчання з вчителем

Є набір прикладів, до кожного прикладу є правильна відповідь. Задача системи – навчитися по прикладах надавати правильну відповідь, задану вчителем. Вчителями є ми.

**SUPERVISED  
LEARNING**



## Навчання без вчителя

Є великий набір даних. В цих даних є приховані закономірності. Задача системи – знайти закономірності, наприклад, розбивши дані на певні групи чи кластери.

**UNSUPERVISED  
LEARNING**



## Навчання з підкріпленням

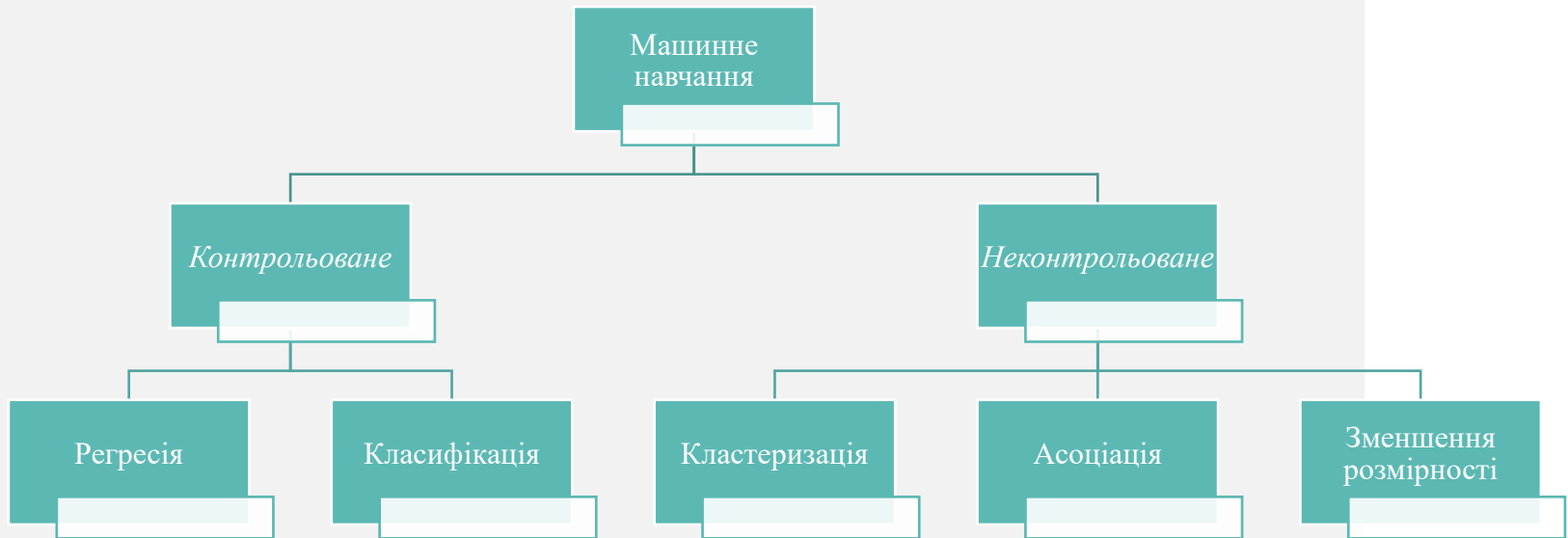
Є певне середовище, в якому є певний агент, що контролюється комп'ютером. Агент може вчиняти певні дії. Певні дії приводять до позитивних відкликів чи негативних відкликів. Задача – максимізувати позитивні і мінімізувати негативні відклики.

Приклад: гра, в якій треба максимізувати набрані бали, або виграти усю гру

**REINFORCEMENT  
LEARNING**



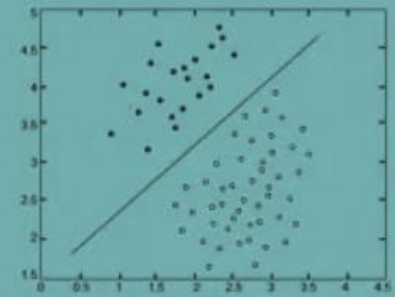
Задачі МН класифікуються ще на декілька типів по виду вирішуваної проблеми:



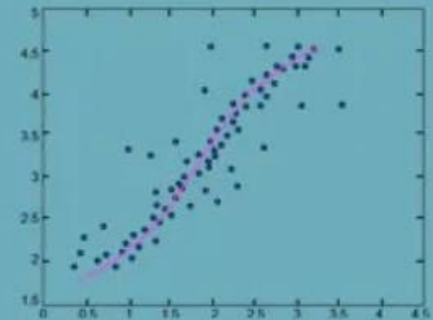
## Задачі

**Класифікація** - передбачення категорії об'єкта

**Регресія** - передбачення місця на числовій прямій.



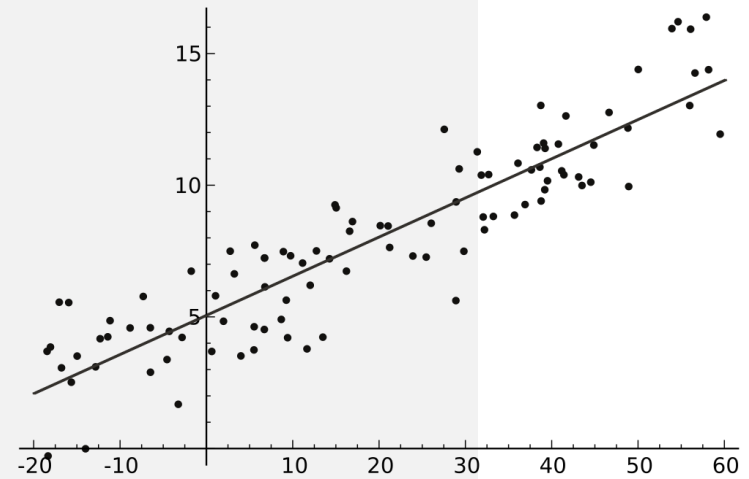
Classification



Regression

# РЕГРЕСІЯ

**Регресія** служить для визначення виду зв'язку між змінними і дає можливість для прогнозування значення однієї (залежної) змінної, відштовхуючись від значень інших (незалежних) змінних.



Регресію використовують для:

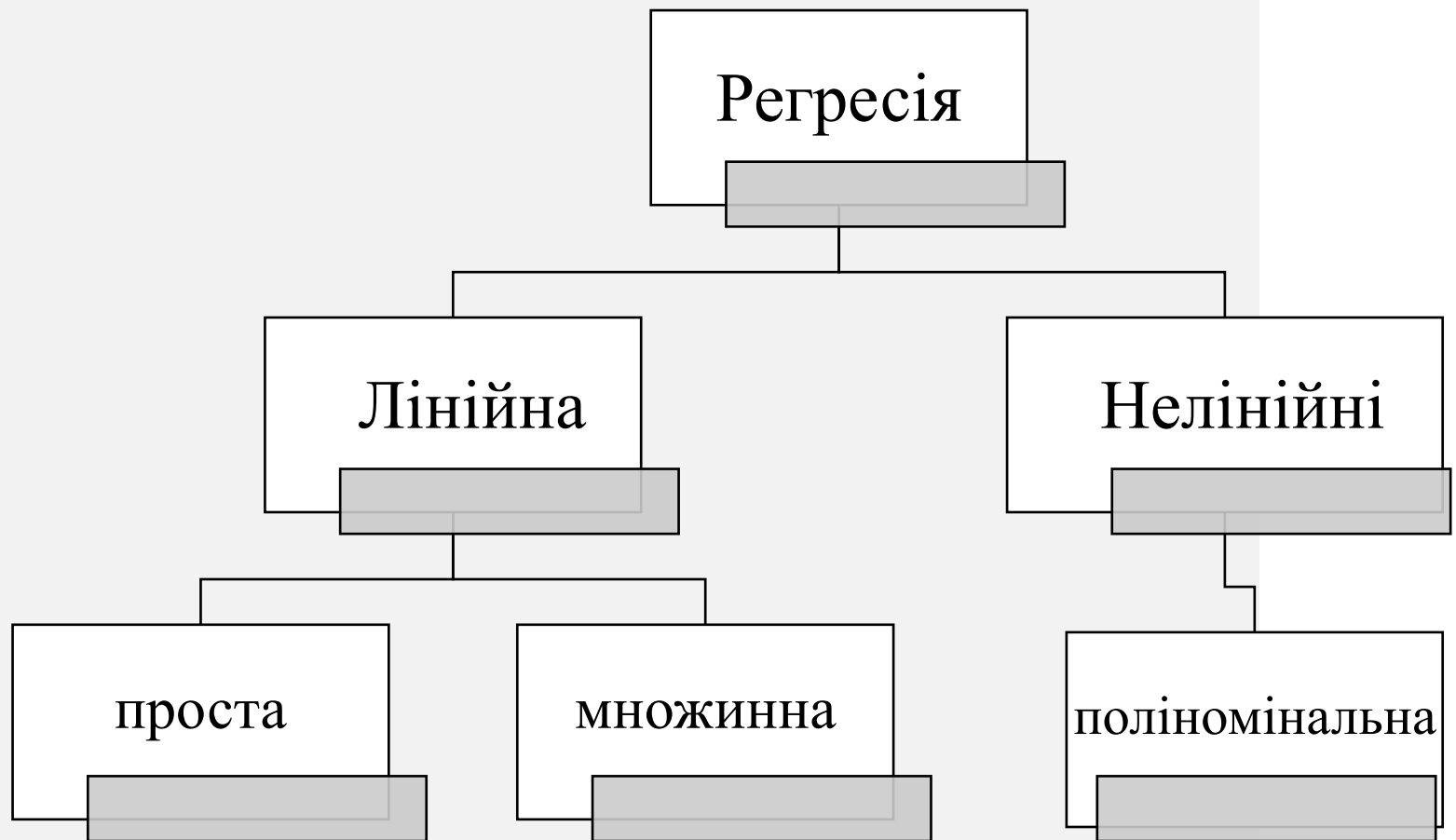
- Прогнозу вартості цінних паперів
- Аналізу попиту, обсягу продажів
- Встановлення медичних діагнозів
- Встановлення будь-якої залежності числа від часу

Регресію використовують у:

- Економіці
- Менеджменті
- Психології



# Моделі



Варіанти і узагальнень лінійної регресії

Метод найменших квадратів

LAD – Least Absolute Deviation,  
метод найменших модулів

Ridge - регресія

Lasso - регресія

...

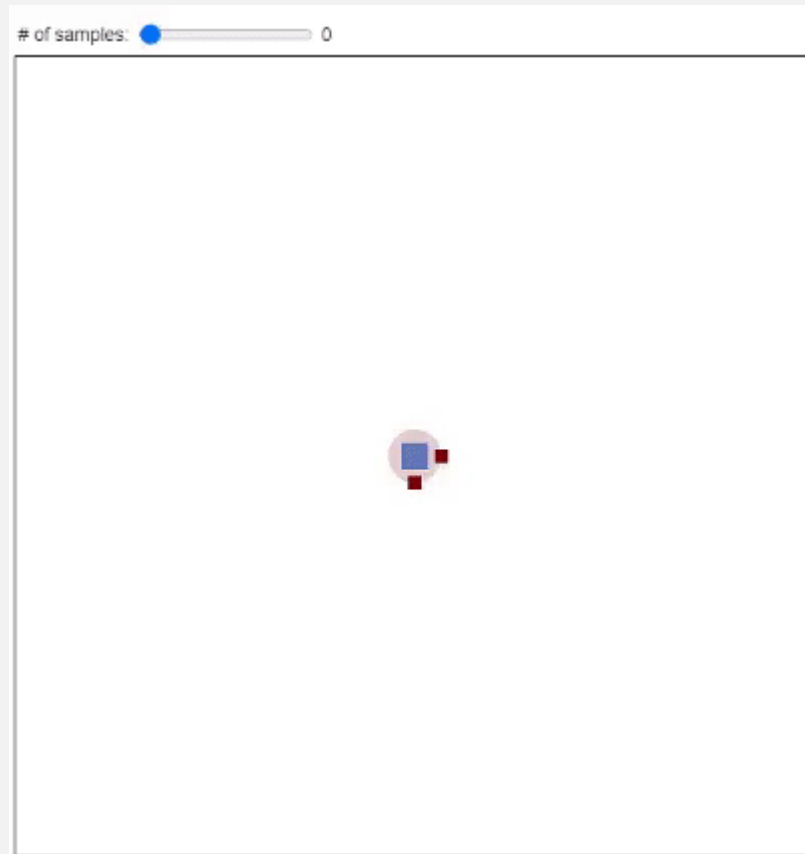
*Лінійна регресія* - «найстаріший» тип, що з'явився два з половиною століття тому.

Вперше метод найменших квадратів опублікував **Адрієн Марі Лежандр** в 1805 році, хоча Гаусс прийшов до нього раніше і успішно використовував для передбачення орбіти «комети» Церери.



**Проста (Парна) лінійна регресія** - це модель, що дозволяє моделювати взаємозв'язок у двовимірному просторі вибірки, утвореному однією незалежною змінною та однією залежною змінною (зазвичай  $x$  і  $y$  — координати у декартовій системі координат).

Модель призначена для знаходження лінійної функції залежності, яка якомога точніше прогнозує значення залежної змінної як функції незалежної змінної.



GoogleColab -

[https://colab.research.google.com/github/fbeilstein/machine\\_learning/blob/master/workbook\\_09\\_linear\\_regression.ipynb#scrollTo=26\\_N4hKfD-Fe](https://colab.research.google.com/github/fbeilstein/machine_learning/blob/master/workbook_09_linear_regression.ipynb#scrollTo=26_N4hKfD-Fe)

# ЛІНІЙНА РЕГРЕСІЯ

Якщо коефіцієнт кореляції дає розуміння чи є лінійна залежність між двома змінними, то лінійна регресія дає модель для оцінки як зміниться одна змінна при зміні іншої.

## *Проста лінійна регресія*

*Лінійною регресією назвемо регресію виду:*

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m$$

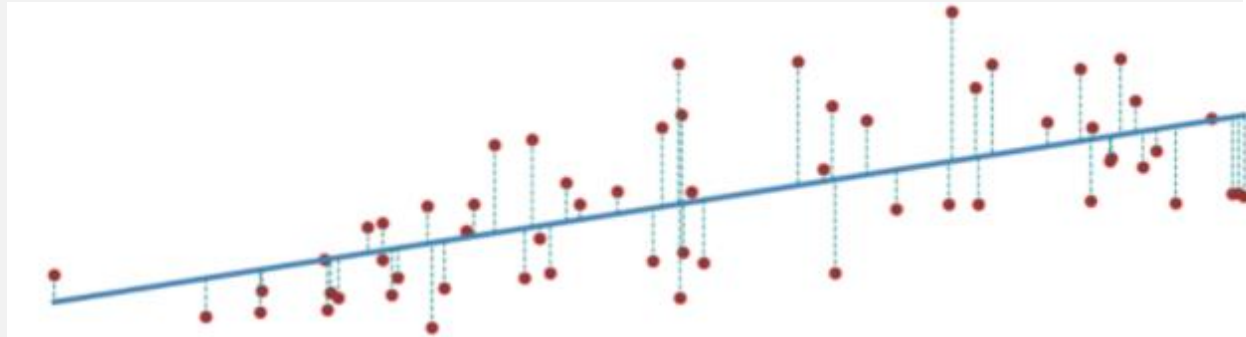
*У випадку парної (простой) регресії вираз для лінійної регресії набуває вигляду:*

$$y = \beta_0 + \beta_1 x + \varepsilon$$

*Де  $\beta_0$  зсув по осі ординат,  $\beta_1$  - регресійні коефіцієнти,  $\varepsilon$  - випадкова помилка змінної  $Y$  в  $i$ -м спостереженні.*

*$\beta_0$  і  $\beta_1$  називаються модельними коефіцієнтами.*

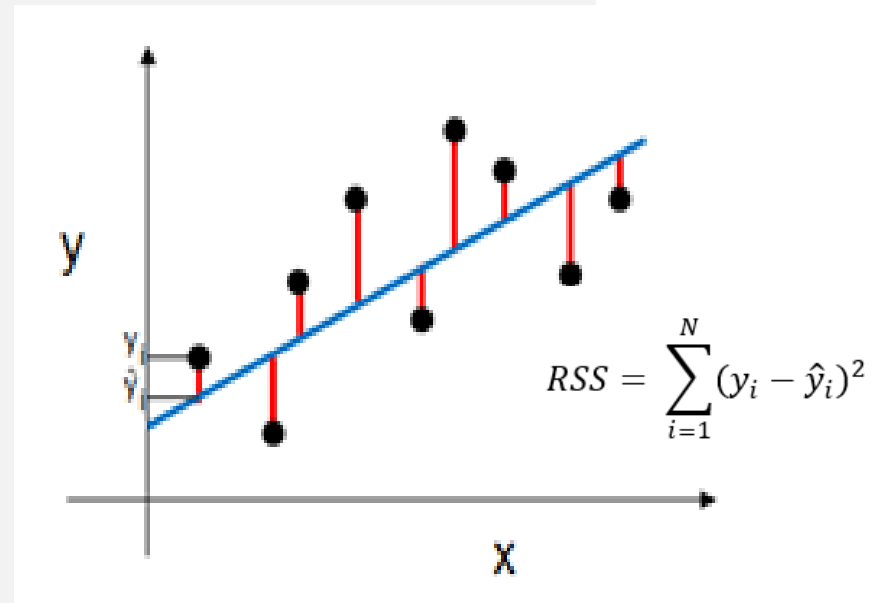
*Щоб створити модель, необхідно дізнатися значення цих коефіцієнтів. І як тільки ці коефіцієнти знайдені, можна використовувати модель для прогнозування.*





## Оцінка ("навчання") модельних коефіцієнтів

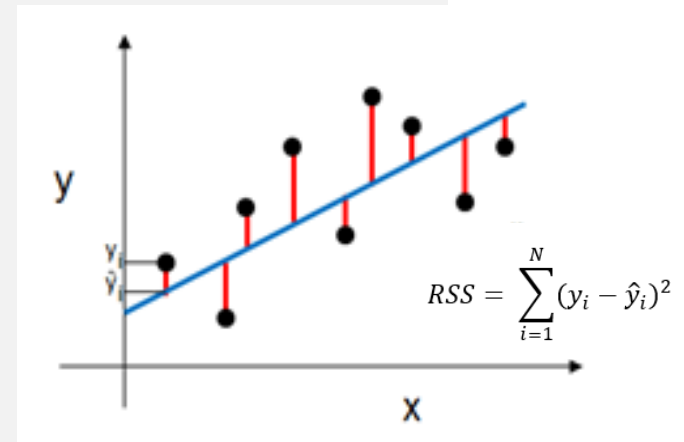
Коефіцієнти оцінюються з використанням методу найменших квадратів, що означає, що необхідно знайти лінію (математично), яка мінімізує суму квадратів помилок (англ. *Residual Sum of Squares (RSS)*)



RSS дозволяє визначити рівень варіативності даних на відстані від лінії.

Лінія найкращої відповідності має найменше значення RSS.

Математичний метод, заснований на мінімізації суми квадратів відхилень деяких функцій від шуканих змінних, називається методом найменших квадратів (англ. Ordinary Least Squares (OLS)).

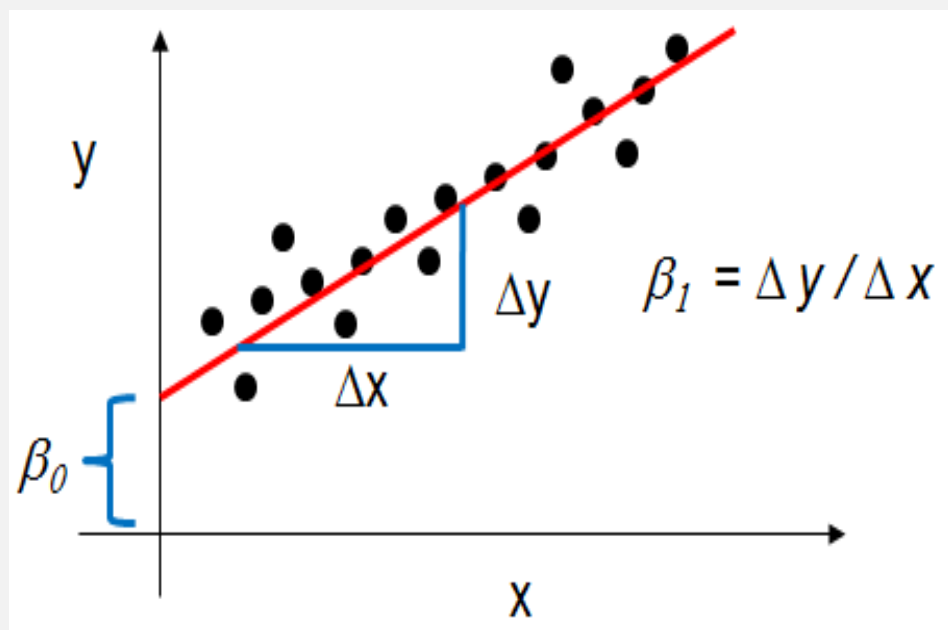


## *Як модельні коефіцієнти відносяться до лінії найменших квадратів?*

*$\beta_0$  є значення  $y$  при  $x=0$*

*$\beta_1$  - нахил (зміна  $y$  поділена на зміну  $x$ )*

*Графічне зображення цих розрахунків:*



# Приклад

В результаті дослідження, було отримано чотири точки  $(x,y)$  даних:  $(1,6)$ ,  $(2,5)$ ,  $(3,7)$  і  $(4,10)$ .

Необхідно знайти пряму  $y=\beta_0+\beta_1x$ , яка найкраще підходить для цих точок. Для цього необхідно знайти  $\beta_0$  і  $\beta_1$  і розв'язати систему рівнянь

$$\beta_0+1\beta_1=6$$

$$\beta_0+2\beta_1=5$$

$$\beta_0+3\beta_1=7$$

$$\beta_0+4\beta_1=10$$

## Метод найменших квадратів:

розв'язання полягає у спробі зробити якомога меншою суму квадратів похибок між правою і лівою сторонами цієї системи, тобто необхідно знайти мінімум функції

$$S(\beta_0, \beta_1) = [6 - (\beta_0 + 1\beta_1)]^2 + [5 - (\beta_0 + 2\beta_1)]^2 + [7 - (\beta_0 + 3\beta_1)]^2 + [10 - (\beta_0 + 4\beta_1)]^2.$$

$\beta_0'$

$$[6 - (\beta_0 + \beta_1)]^2 + [5 - (\beta_0 + 2\beta_1)]^2 + [7 - (\beta_0 + 3\beta_1)]^2 + [10 - (\beta_0 + 4\beta_1)]^2 =$$

$$2[6 - (\beta_0 + \beta_1)] + 2[5 - (\beta_0 + 2\beta_1)] + 2[7 - (\beta_0 + 3\beta_1)] + 2[10 - (\beta_0 + 4\beta_1)] =$$

$$12 - 2(\beta_0 + \beta_1) + 10 - 2(\beta_0 + 2\beta_1) + 14 - 2(\beta_0 + 3\beta_1) + 20 - 2(\beta_0 + 4\beta_1) =$$

$$56 - 2\beta_0 - 2\beta_1 - 2\beta_0 - 4\beta_1 - 2\beta_0 - 6\beta_1 - 2\beta_0 - 8\beta_1 =$$

$$56 - 8\beta_0 - 20\beta_1$$

$$\frac{\partial S}{\partial \beta_0} = 0 = 8\beta_0 + 20\beta_1 - 56$$

$\beta_1$

$$[6-(\beta_0+\beta_1)]^2+[5-(\beta_0+2\beta_1)]^2+[7-(\beta_0+3\beta_1)]^2+[10-(\beta_0+4\beta_1)]^2=$$

$$2[6-(\beta_0+\beta_1)]+2*2[5-(\beta_0+2\beta_1)]+2*3[7-(\beta_0+3\beta_1)]+2*4[10-(\beta_0+4\beta_1)]=$$

$$12-2\beta_0-2\beta_1+20-4\beta_0-8\beta_1+42-6\beta_0-18\beta_1+80-8\beta_0+32\beta_1=$$

$$154-20\beta_0-60\beta_1$$

$$\frac{\partial S}{\partial \beta_1} = 0 = 20\beta_0 + 60\beta_1 - 154$$

Мінімум визначають через обчислення часткової похідної від  $S(\beta_0, \beta_1)$  щодо  $\beta_0$  і  $\beta_1$  і прирівнюванням її до нуля

$$\frac{\partial S}{\partial \beta_0} = 0 = 8\beta_0 + 20\beta_1 - 56$$

$$\frac{\partial S}{\partial \beta_1} = 0 = 20\beta_0 + 60\beta_1 - 154$$

Це приводить до системи з двох рівнянь і двох невідомих, які називаються нормальними рівняннями.



В результаті рішення системи з двох рівнянь отримуємо

$$\beta_0=3.5$$

$$\beta_1=1.4$$

Рівняння лінії :  $y=3.5+1.4x$

Мінімальна сума квадратів похибок:

$$S(3.5,1.4)=1.1^2+(-1.3)^2+(-0.7)^2+0.9^2=4.2.$$



# Самостійна робота

Розібрати приклад:

`scipy.stats.linregress` - <https://docs.scipy.org/doc/scipy-0.19.1/reference/generated/scipy.stats.linregress.html>

Що потрібно знати о регресії –

1. Для чого застосовується регресія?
2. Що таке лінійна регресія?
3. У чому суть методу найменших квадратів?
4. Що таке нахил у рівнянні лінійної регресії?
5. Як розраховуються коефіцієнти рівняння лінійної регресії?
6. Які переваги і недоліки методу найменших квадратів?