



Машинне навчання

Лекція 5: Метод найближчих сусідів

Кочура Юрій Петрович
iuriy.kochura@gmail.com
[@y_kochura](#)

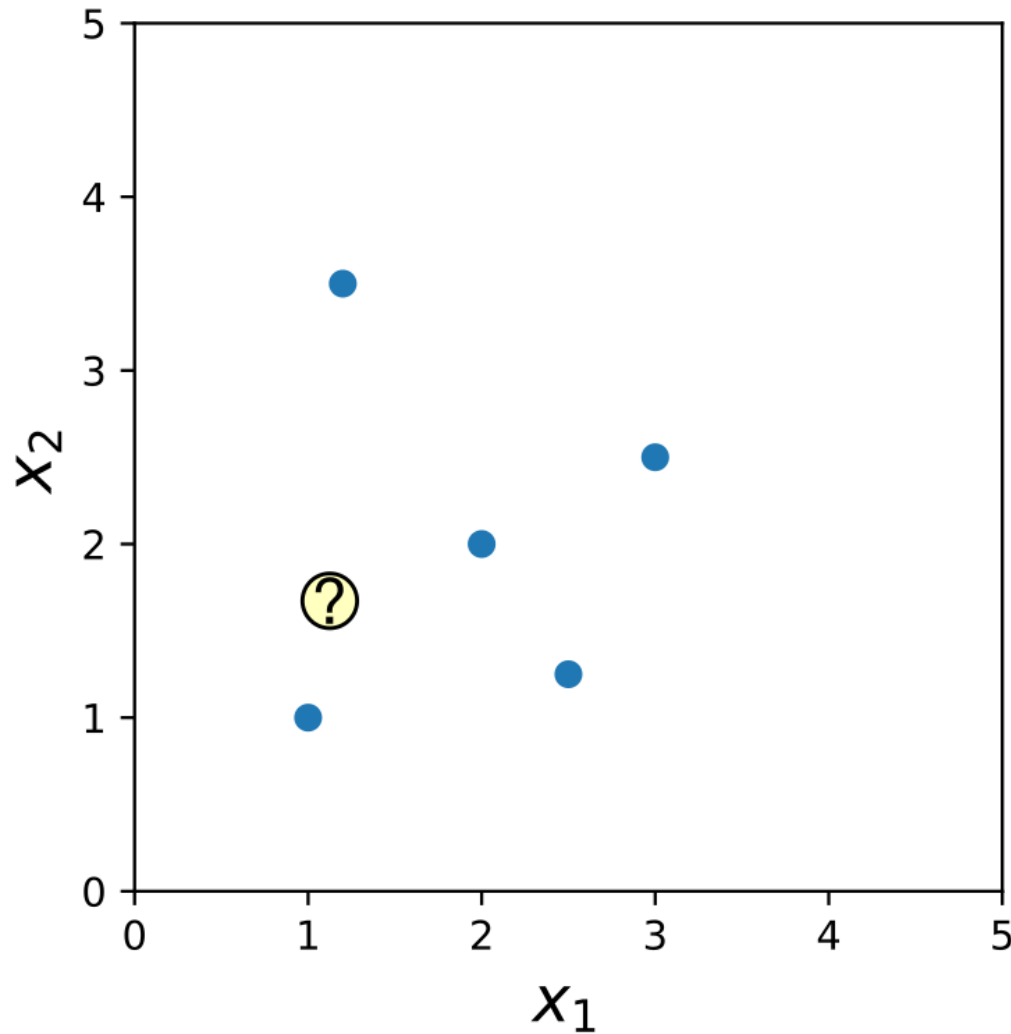
СЬОГОДНІ

- 🎙️ k-найближчих сусіди
- 🎙️ Неперервні міри відстані
- 🎙️ Дискретні міри відстані
- 🎙️ Редагування та покращення ефективності kNN

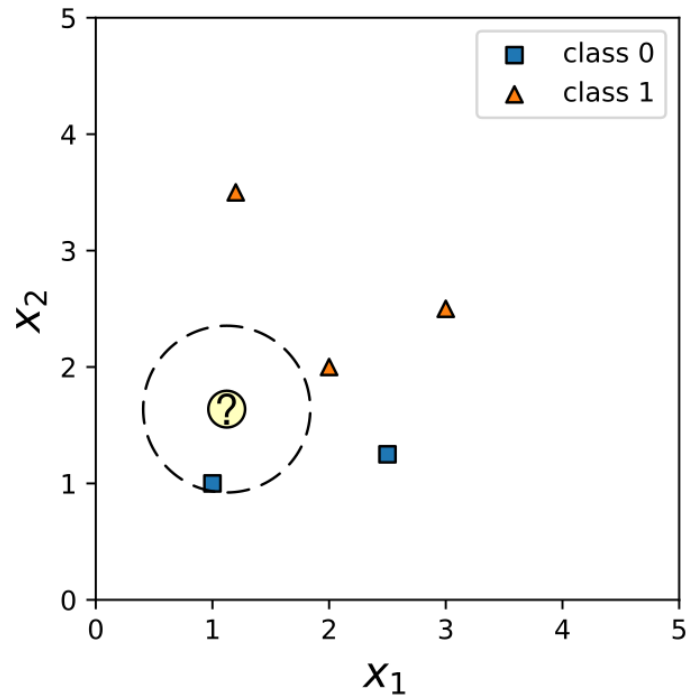
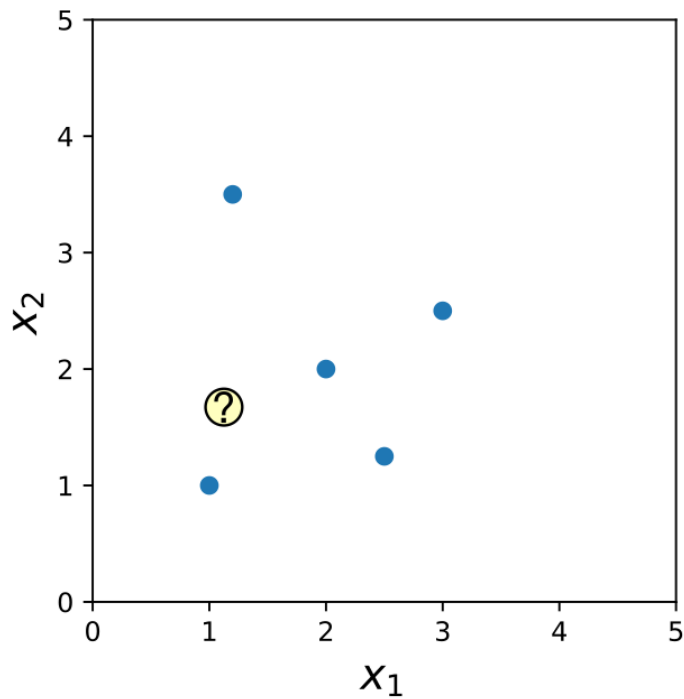
к-найближчих сусіди

k-nearest neighbors (kNN)

1-найближчий сусід



1-найближчий сусід



Навчальна вибірка

$$\left(\mathbf{X}^{(i)}, y^{(i)} \right) \in \mathcal{D},$$

$$|\mathcal{D}| = n$$

Псевдокод: 1-найближчий сусід

closest_point := None

closest_distance := ∞

- for $i = 1, \dots, n$:

current_distance := $d(\mathbf{X}^{(i)}, \mathbf{X}^{(q)})$

if current_distance < closest_distance:

- closest_distance := current_distance
- closest_point := $\mathbf{X}^{(i)}$
- return $f(\mathbf{X}^{(i)})$

Мітка найближчого сусіда (closest_point): $\left(\mathbf{X}^{(i)}, f(\mathbf{X}^{(i)}) \right)$

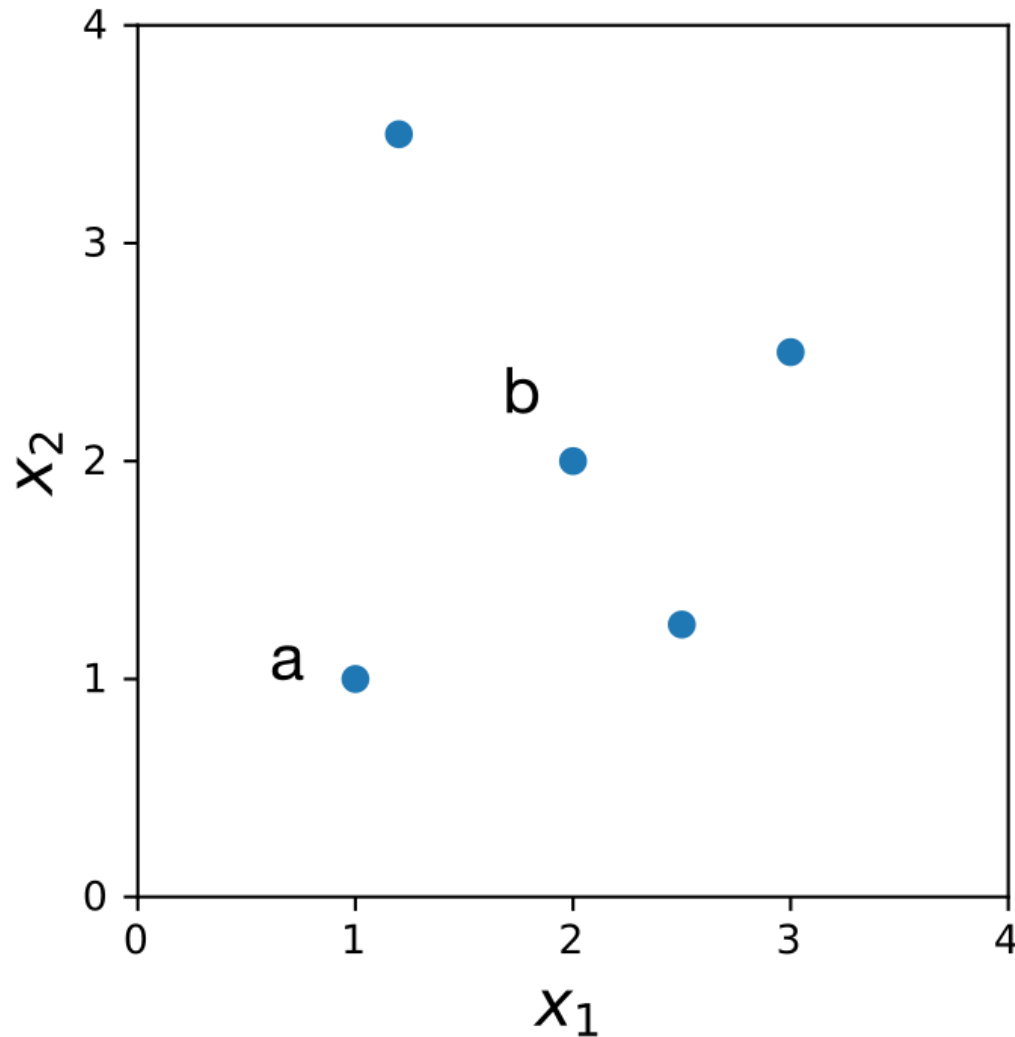
Евклідова відстань (L2)

$$\begin{aligned} d(\mathbf{X}^{(a)}, \mathbf{X}^{(b)}) &= \sqrt{\sum_{j=1}^m \left(\mathbf{x}_j^{(a)} - \mathbf{x}_j^{(b)} \right)^2} = \\ &= \left\| \mathbf{X}^{(a)} - \mathbf{X}^{(b)} \right\|_2 \end{aligned}$$

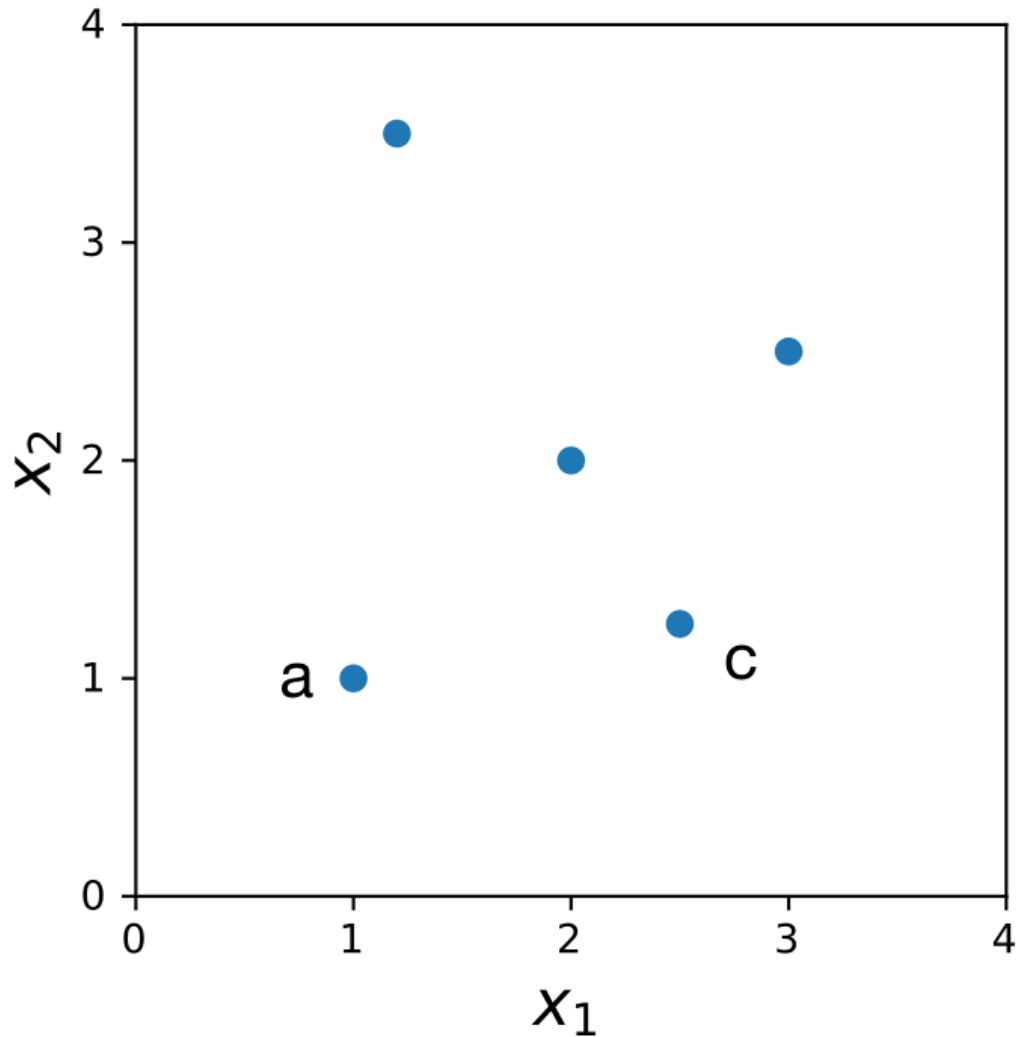
де $\mathbf{X}^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_m^{(i)}) \in \mathbb{R}^m$ – m -вимірний вектор ознак.

Граница рішень методу найближчих сусідів

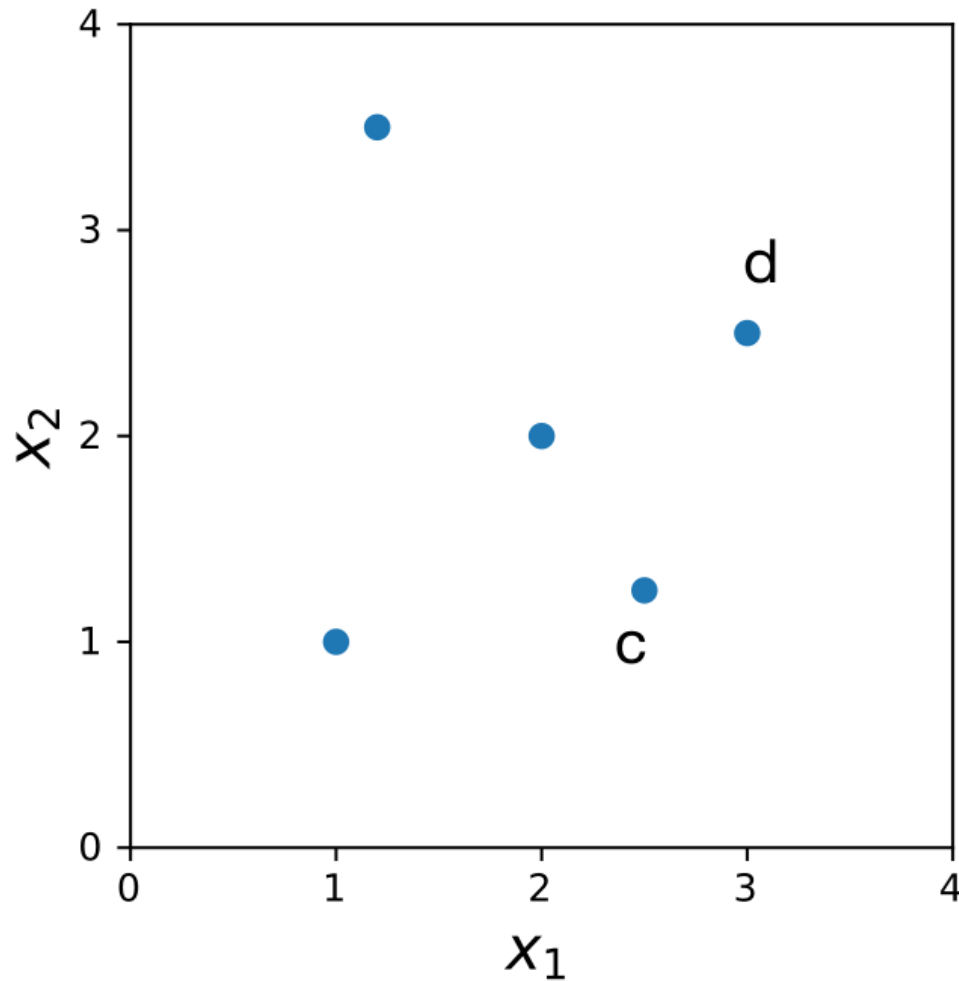
Границя рішень між (a) та (b)



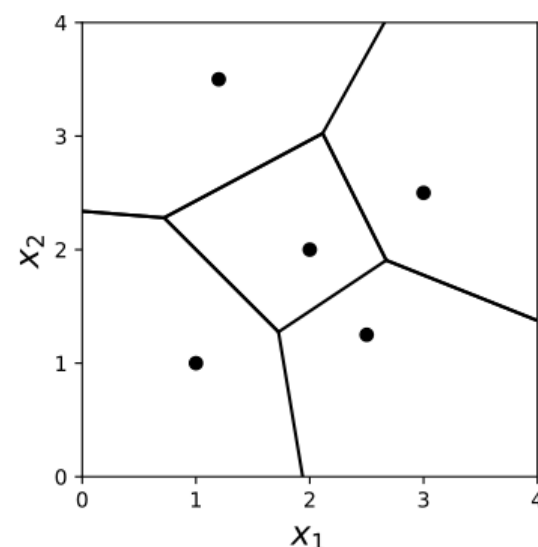
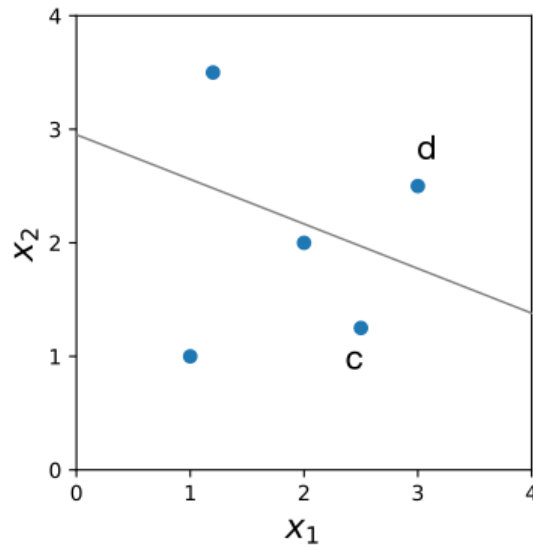
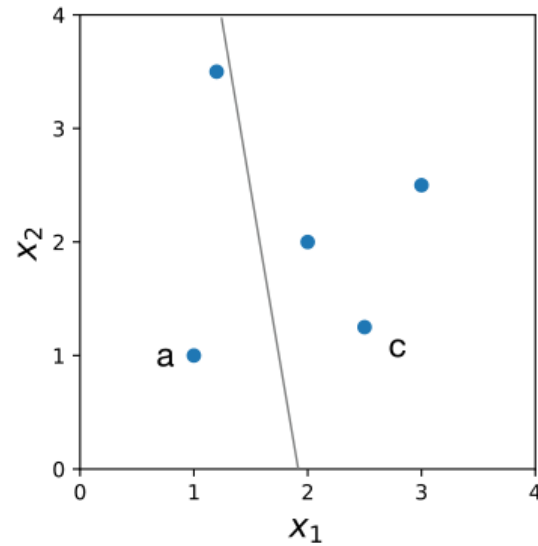
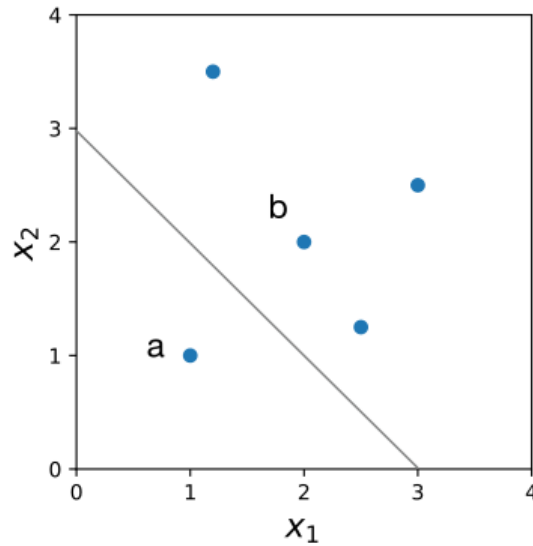
Границя рішень між (a) та (c)



Границя рішень між (c) та (d)



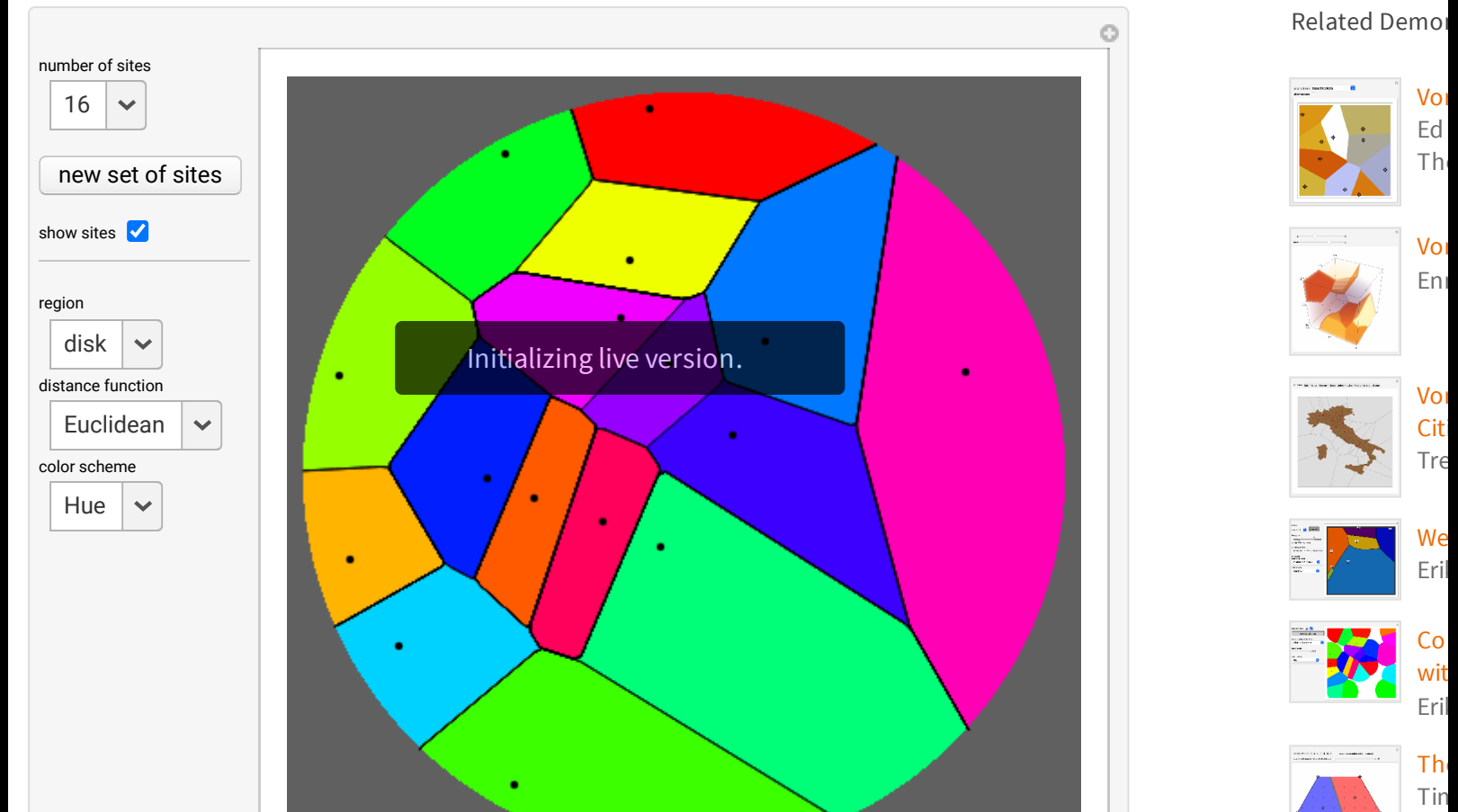
Границя рішень для 1-NN





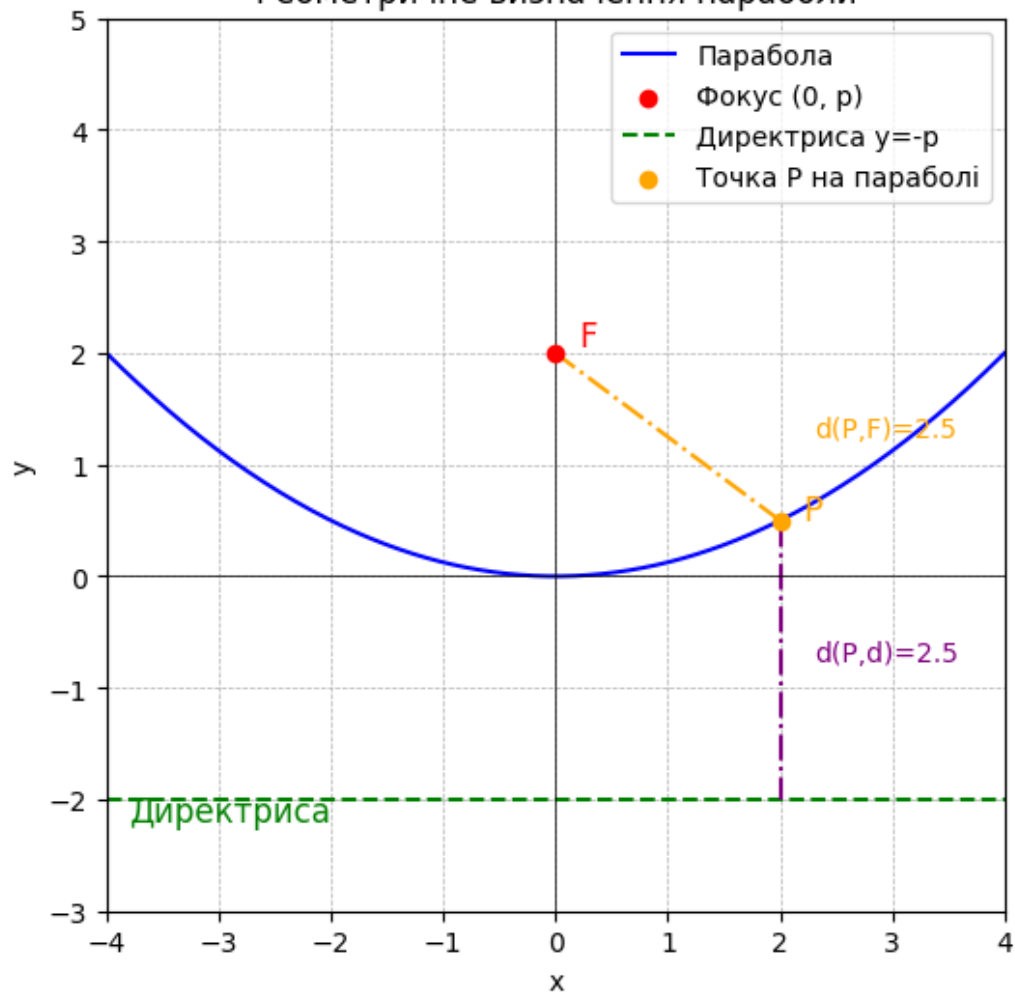
Георгій Феодосійович Вороний
(1868-1908)

Voronoi Diagrams in Two-Dimensional Regions



Діаграми Вороного в двовимірних областях

Геометричне визначення параболи



$$y = \frac{x^2}{4p}$$

Fortune's Algorithm for Voronoi Diagrams

[Open Notebook in Cloud](#)[Copy Manipulate to Clipboard](#)[Initializing live version](#)[Source Code](#)

Given a finite set of points (called sites) in a plane, a Voronoi diagram divides the plane into regions around each site such that each region contains all points closer to that site than to any of the others. This Demonstration shows Fortune's algorithm for drawing Voronoi diagrams[1].

Two auxiliary curves are involved in the procedure:

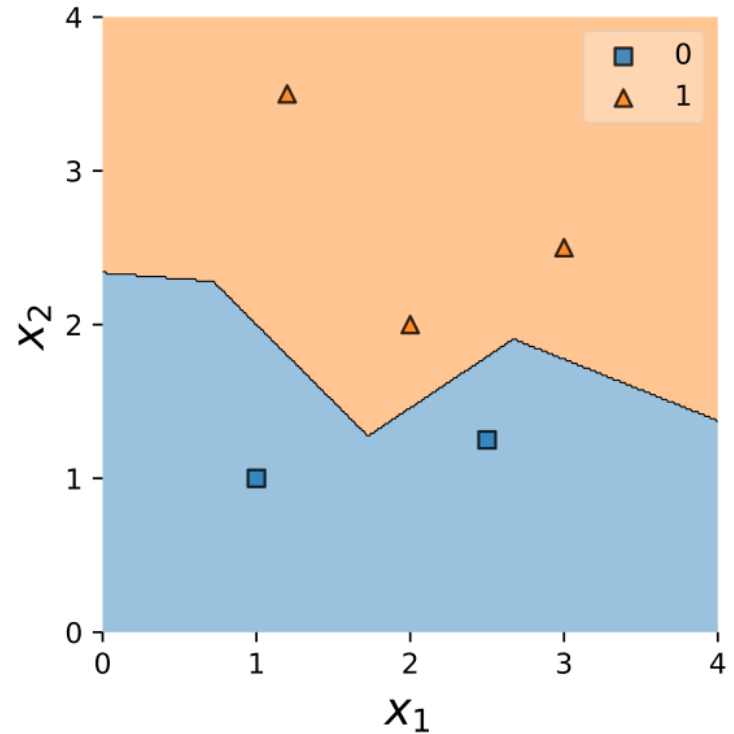
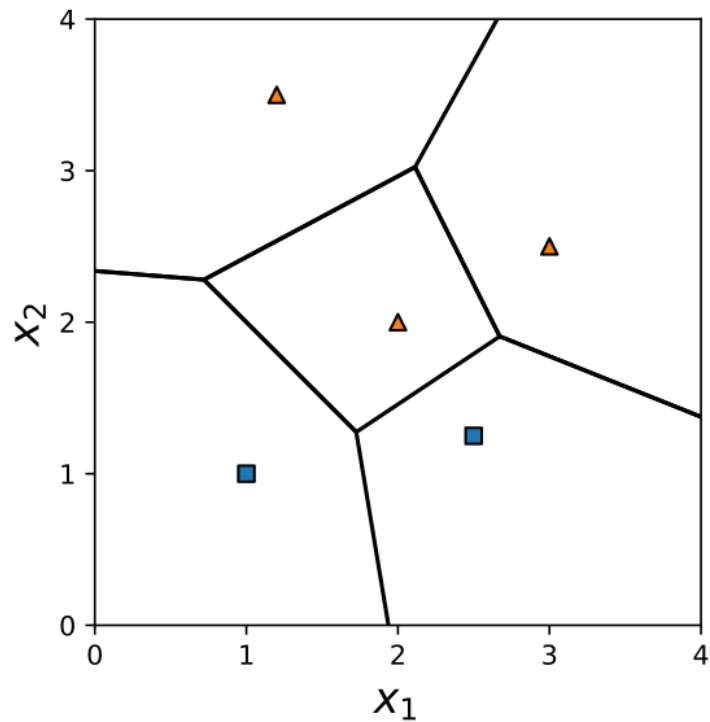
1. The sweep line (green) can move up or down; only sites above it are active.
2. The beach line (red) is a sequence of parabolic arcs. Each parabola has an active site for its locus and the sweep line as its directrix.

The intersections of two parabolic arcs (the breakpoints on the beach line) have equal distance to two sites (the two parabolas) and to the sweep line. Thus, as the sweep line moves down, these intersection points determine the rest of the diagram.

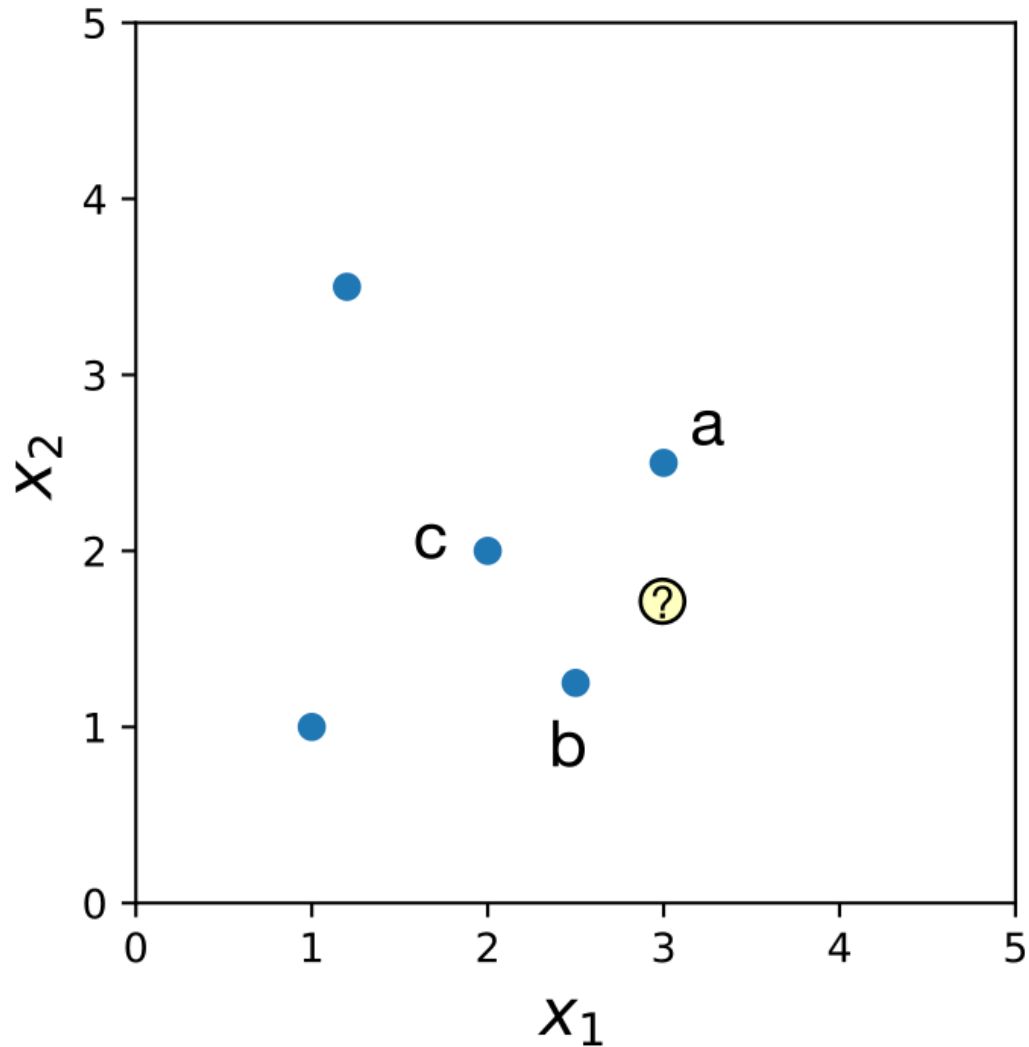
Contributed by: [Erik Mahieu](#) (2016)

Open content licensed under [CC BY-NC-SA](#)

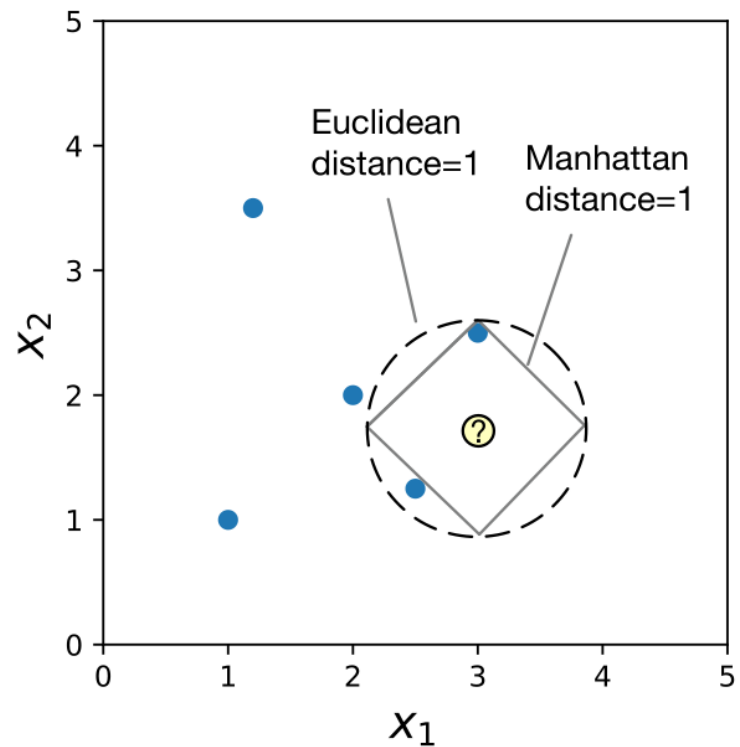
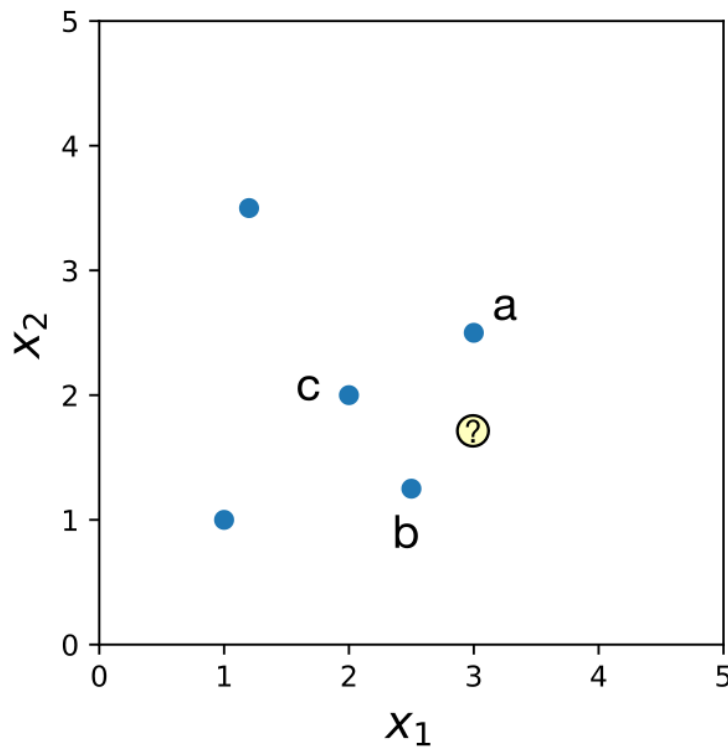
Границя рішень для 1-NN



Який екземпляр найближчий?



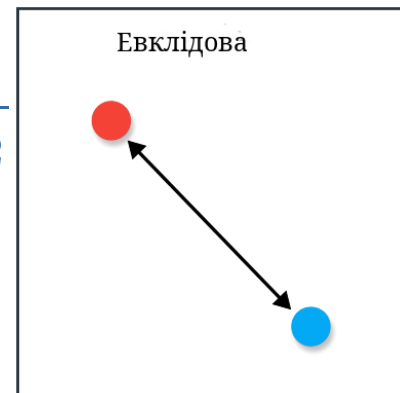
Залежить від міри відстані!



Неперервні міри відстані

Евклідова відстань:

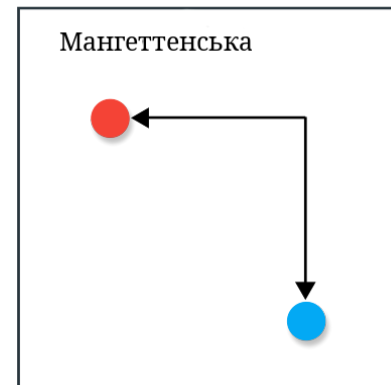
$$d(\mathbf{X}^{(a)}, \mathbf{X}^{(b)}) = \sqrt{\sum_{j=1}^m \left(\mathbf{x}_j^{(a)} - \mathbf{x}_j^{(b)} \right)^2}$$



Неперервні міри відстані

Мангеттенська відстань:

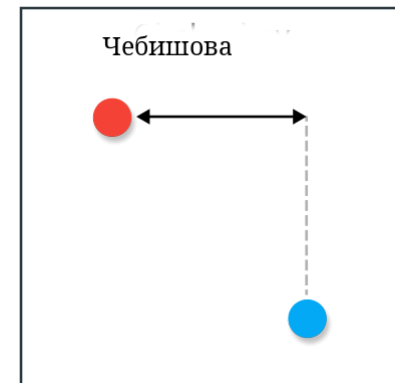
$$d(\mathbf{X}^{(a)}, \mathbf{X}^{(b)}) = \sum_{j=1}^m |\mathbf{X}_j^{(a)} - \mathbf{X}_j^{(b)}|$$



Неперервні міри відстані

Чебишова відстань:

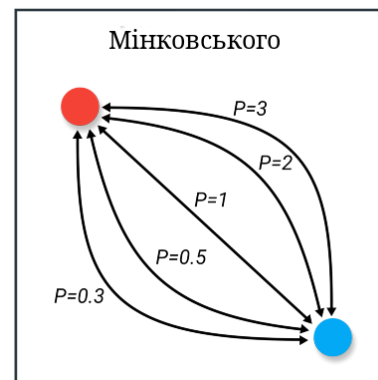
$$d(\mathbf{X}^{(a)}, \mathbf{X}^{(b)}) = \max_{j=1}^m |\mathbf{X}_j^{(a)} - \mathbf{X}_j^{(b)}|$$



Неперервні міри відстані

Мінковського відстань:

$$d(\mathbf{X}^{(a)}, \mathbf{X}^{(b)}) = \left[\sum_{j=1}^m (|\mathbf{X}_j^{(a)} - \mathbf{X}_j^{(b)}|)^p \right]^{\frac{1}{p}}$$



Неперервні міри відстані

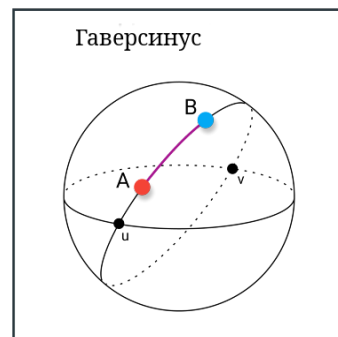
Гаверсинус кута θ визначається як:

$$\text{hav}(\theta) = \sin^2 \left(\frac{\theta}{2} \right)$$

Відстань гаверсинуса:

$$d = 2r \cdot \arcsin \sqrt{\text{hav}(\phi_b - \phi_a) + \cos(\phi_a) \cos(\phi_b) \text{hav}(\lambda_b - \lambda_a)}$$

де r — радіус сфери, ϕ — широта, λ — довгота.



Неперервні міри відстані

Відстань Махаланобіса:

$$d(\mathbf{X}^{(a)}, \mathbf{X}^{(b)}) = \sqrt{(\mathbf{X}^{(a)} - \mathbf{X}^{(b)})^T S^{-1} (\mathbf{X}^{(a)} - \mathbf{X}^{(b)})},$$

де S^{-1} — коваріаційна матриця.

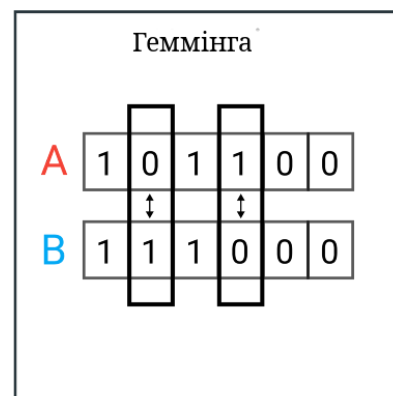
Дискретні міри відстані

Відстань Геммінга:

$$d(\mathbf{X}^{(a)}, \mathbf{X}^{(b)}) = \sum_{j=1}^m \delta(\mathbf{X}_j^{(a)}, \mathbf{X}_j^{(b)}),$$

де

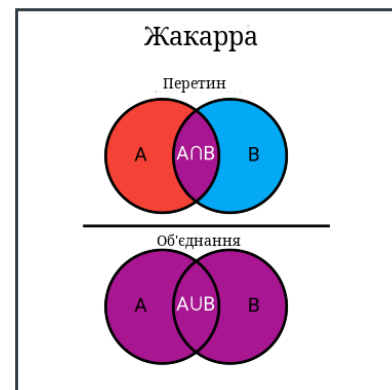
$$\delta(\mathbf{X}_j^{(a)}, \mathbf{X}_j^{(b)}) = \begin{cases} 1, & \text{якщо } \mathbf{X}_j^{(a)} \neq \mathbf{X}_j^{(b)} \\ 0, & \text{якщо } \mathbf{X}_j^{(a)} = \mathbf{X}_j^{(b)} \end{cases}$$



Дискретні міри відстані

Відстань Жаккара/Танімото:

$$d(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

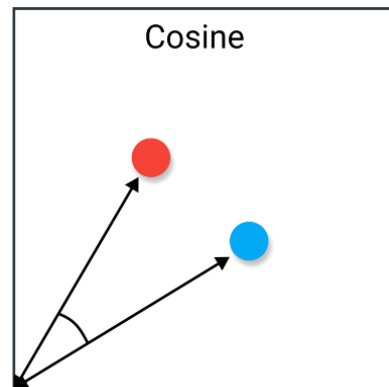


Дискретні міри відстані

Косинусна відстань:

$$d(A, B) = 1 - \cos(\theta),$$

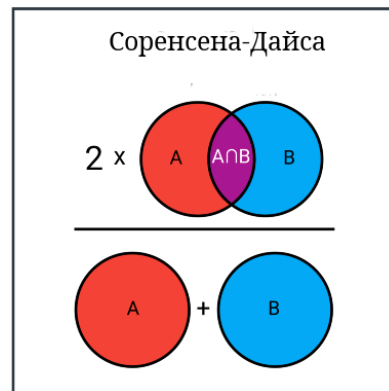
$$\text{де } \cos(\theta) = \frac{A \cdot B}{\|A\| \cdot \|B\|}$$



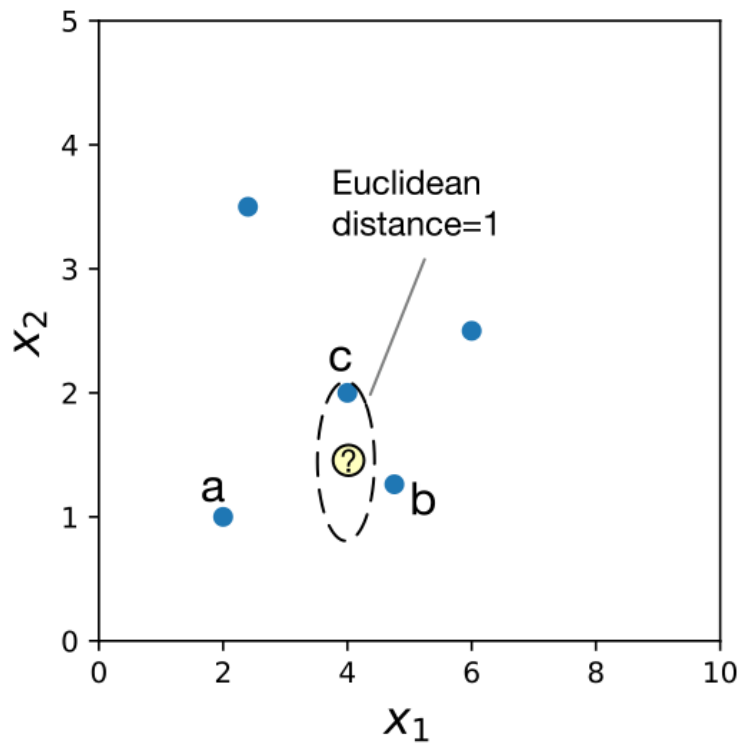
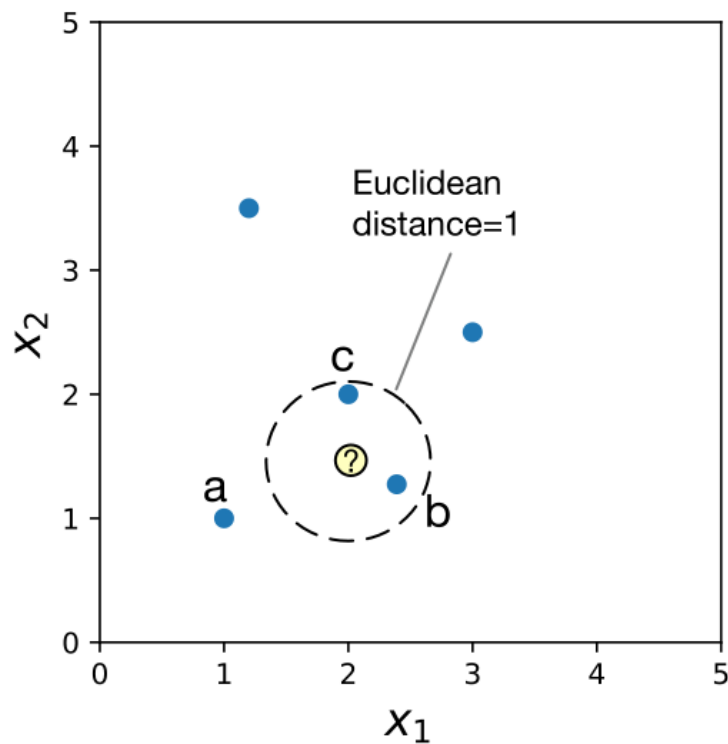
Дискретні міри відстані

Відстань Соренсена-Дайса:

$$d(A, B) = \frac{2|A \cap B|}{|A| + |B|}$$

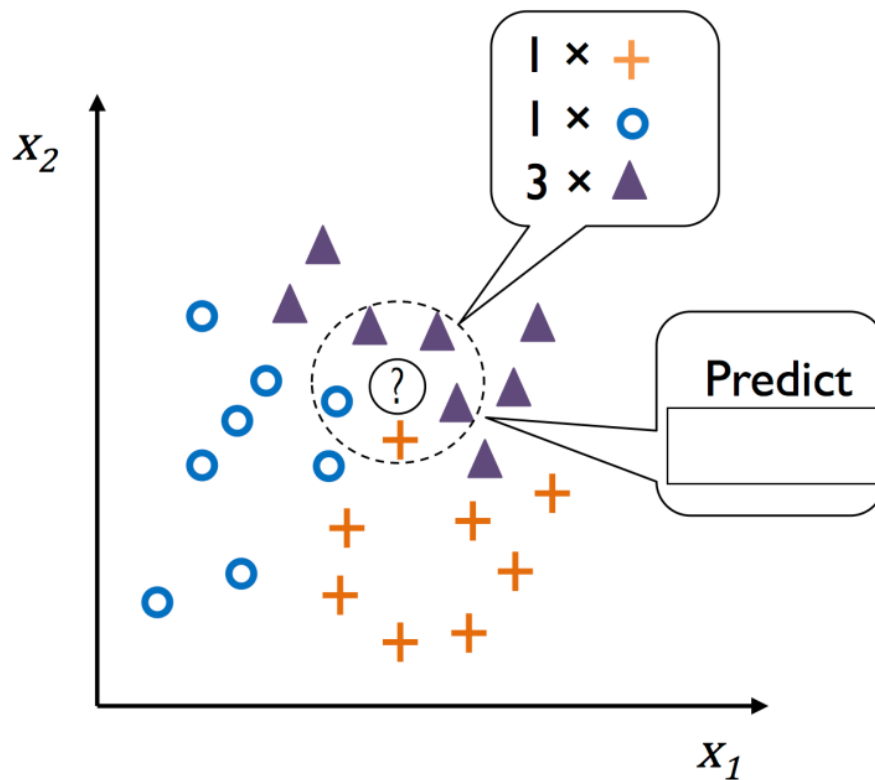


Масштабування ознак




k-найближчі сусіди


$k = 5$



A



Голосування більшістю: 

Кількісна більшість: 

B

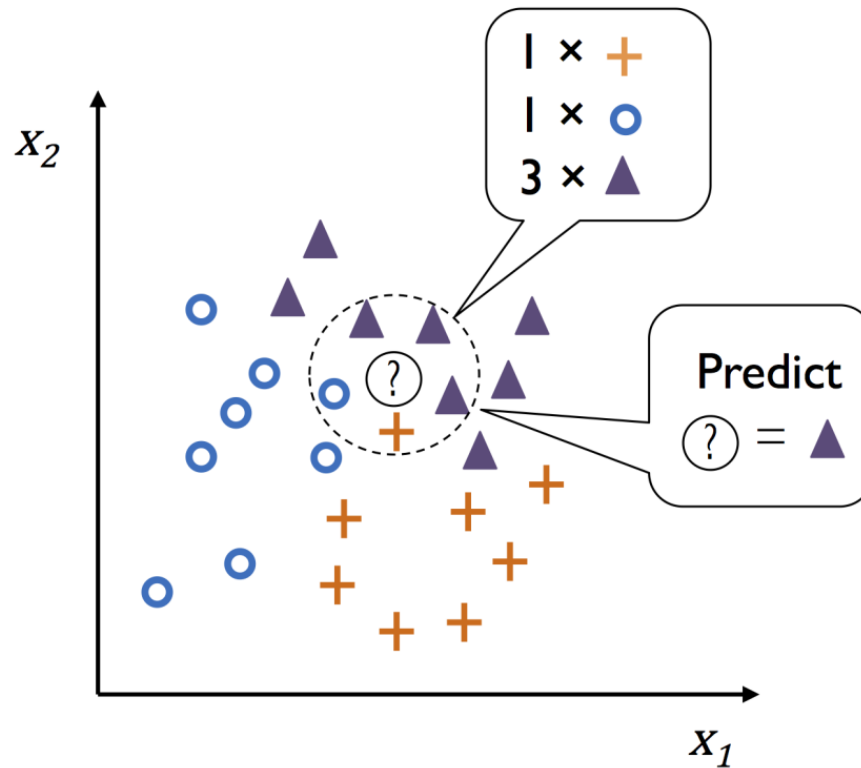


Голосування більшістю: None

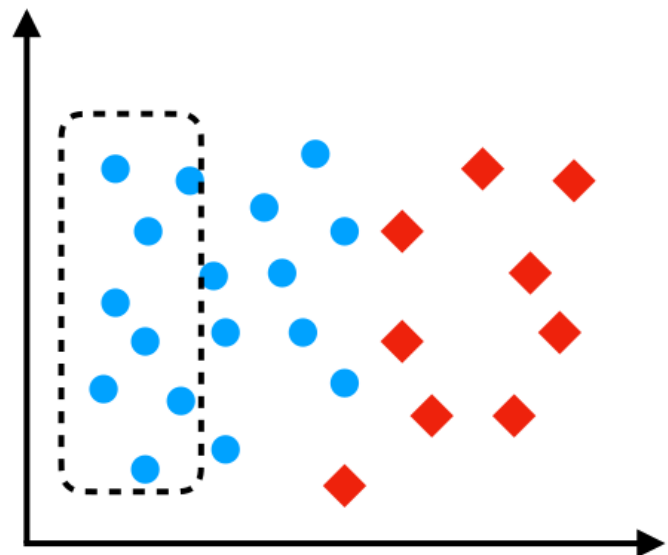
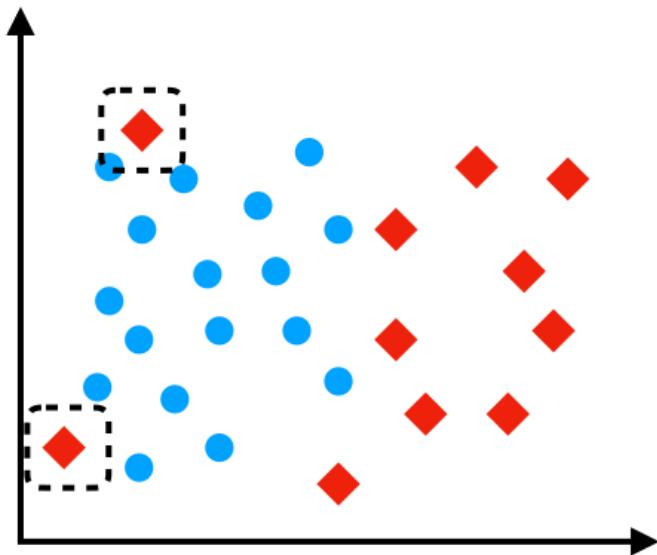
Кількісна більшість: 

k-найближчі сусіди

$k = 5$



Редагування kNN



Покращення ефективності прогнозування

- Підібрати оптимальне значення k
- Масштабування осей ознак
- Вибір метрики для визначення відстані
- Зважування міри відстані

Кінець