

KHAI PHÁ DỮ LIỆU

VO DUC QUANG – VINH UNIVERSITY

QUANGVD@VINHUNI.EDU.VN



Nội dung

- Chương 1: Tổng quan về Data Mining
- Chương 2: Dữ liệu và tiền xử lý dữ liệu
- Chương 3: Bài toán phân lớp dữ liệu
- Chương 4: Bài toán phân cụm dữ liệu
- Chương 5: Khai phá luật kết hợp

Chương 3 – Phân lớp (Classification)

- Bài toán Phân lớp
- Phân lớp sử dụng Cây quyết định
- Một số giải thuật khác
 - Naïve Bayes
 - KNN
- Đánh giá mô hình phân lớp

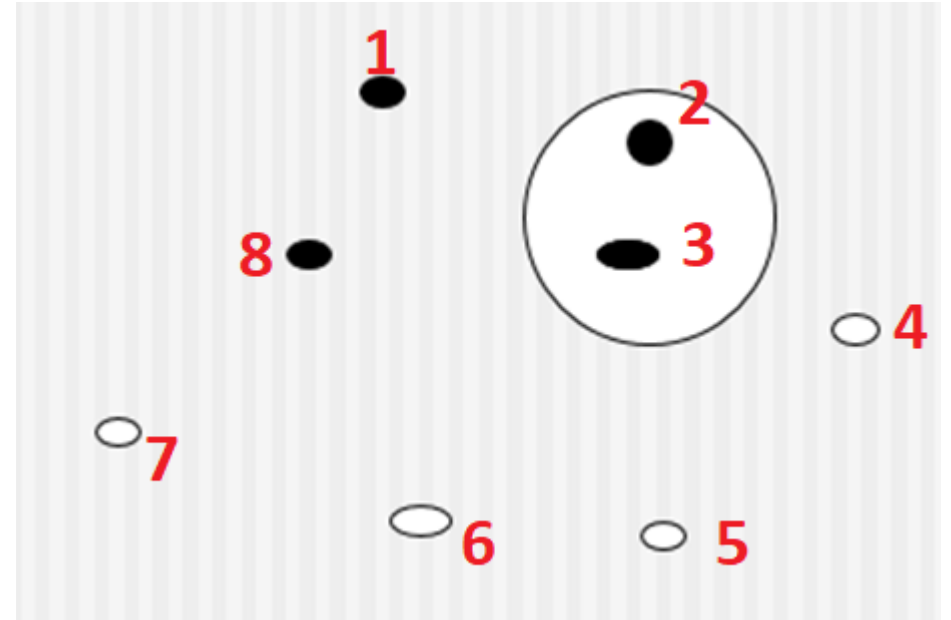
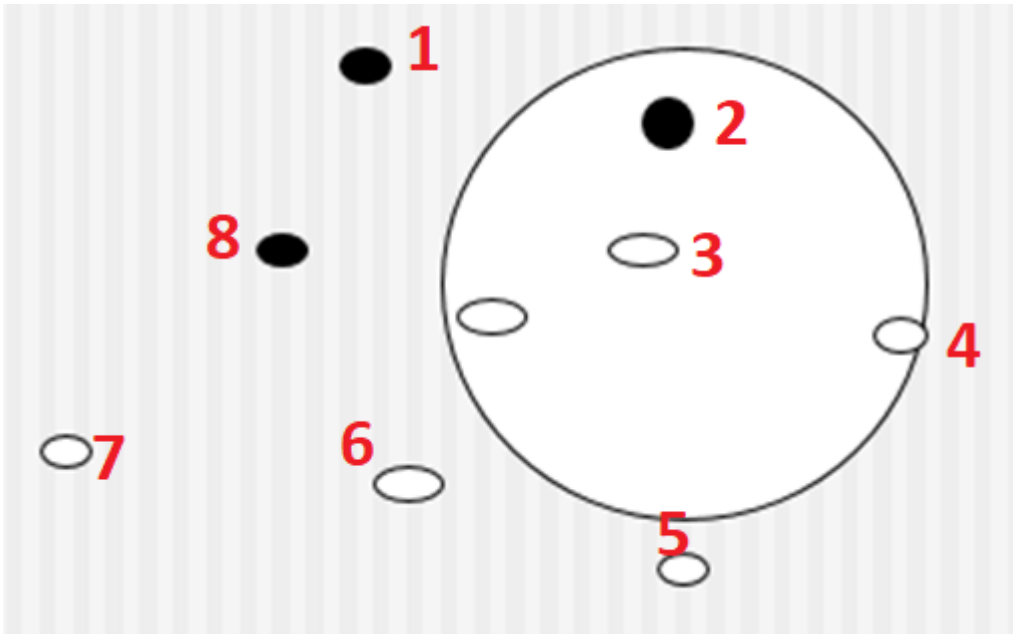
Ôn tập khái niệm

- Xem ví dụ (Machinelearningcoban.com)

Ngôn ngữ người	Ngôn ngữ Máy Học	in Machine Learning
Câu hỏi	Điểm dữ liệu	Data point
Đáp án	Đầu ra, nhãn	Output, Label
Ôn thi	Huấn luyện	Training
Tập tài liệu mang vào phòng thi	Tập dữ liệu tập huấn	Training set
Đề thi	Tập dữ liệu kiểm thử	Test set
Câu hỏi trong đề thi	Dữ liệu kiểm thử	Test data point
Câu hỏi có đáp án sai	Nhiều	Noise, Outlier
Câu hỏi gần giống	Điểm dữ liệu gần nhất	Nearest Neighbor

K-nearest neighbor (KNN)

- [lazy learning](#)
- Neighbor: Láng giềng
- K-nearest neighbor ?



$V_i^2 \rightarrow \text{Min}$

Min 20

$p_{1,2} = 0$

Điện 3 là hay xôn của 2

$$\textcircled{F} \text{ KC}(2,3) = \sqrt{(x_3 - x_2)^2 + (y_3 - y_2)^2}$$

~~* $K_c(2, 4) = \dots$~~

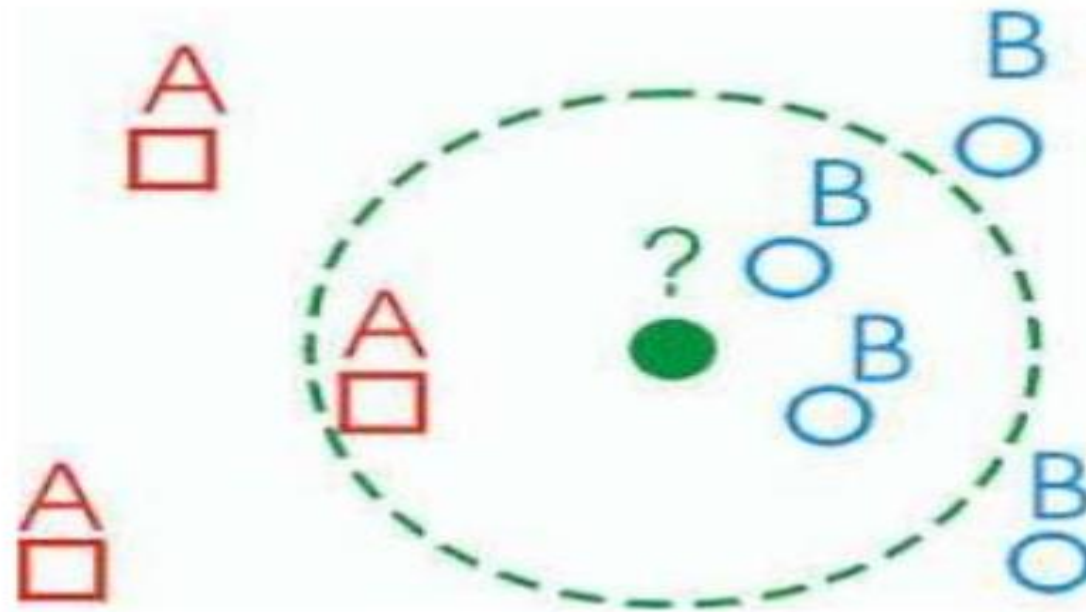
$$K_C(2, 1) = - \dots$$

En lid

$$d(2,3) = \sqrt{(x_3 - x_2)^2 + (y_3 - y_2)^2 + (z_3 - z_2)^2}$$

KNN

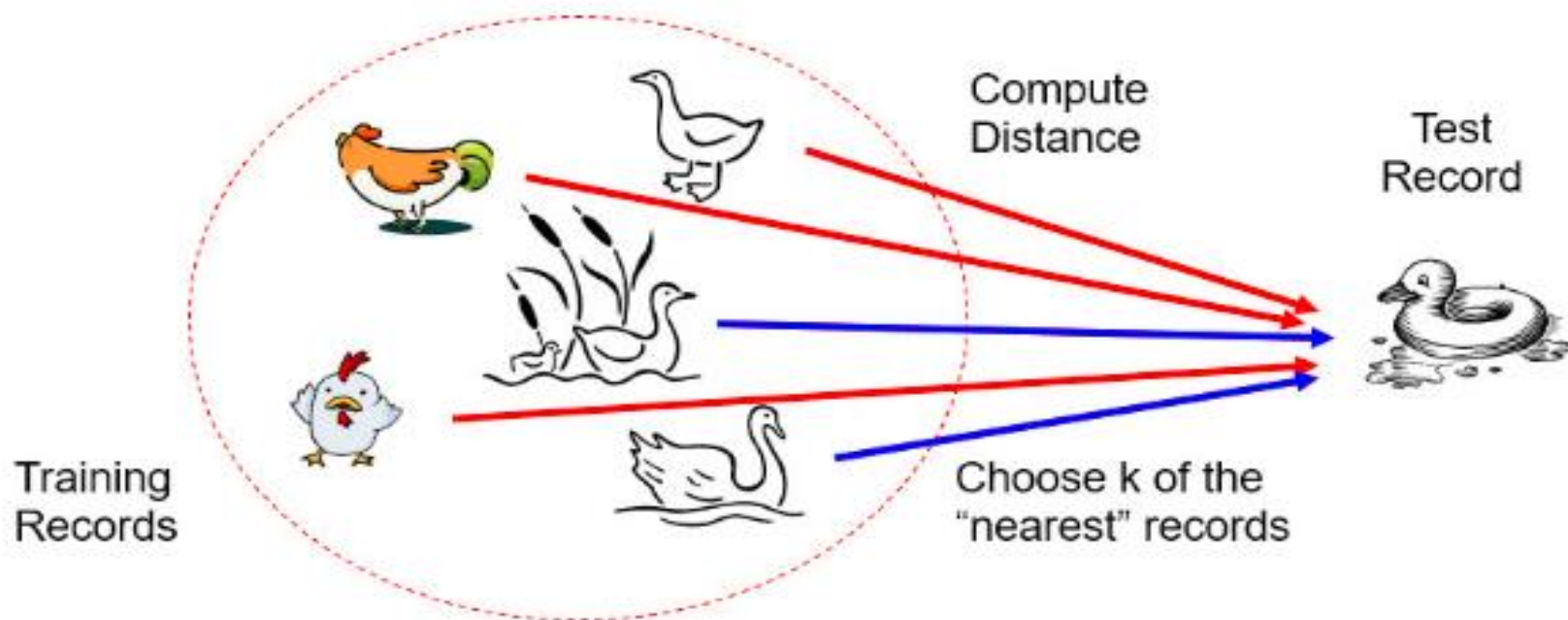
- Xác định nhãn cho một điểm?



KNN

K-Nearest Neighbors: *"Ask not what this is called, ask what this is like?"*

- Nếu nó đi như một con vịt, lang thang như một con vịt, thì nó có thể là một con vịt.



KNN

■ Demo code

Table 1. Euclidean distance matrix D listing all possible pairwise Euclidean distances between 19 samples.

	x ₁	x ₂	x ₃	x ₄	x ₅	x ₆	x ₇	x ₈	x ₉	x ₁₀	x ₁₁	x ₁₂	x ₁₃	x ₁₄	x ₁₅	x ₁₆	x ₁₇	x ₁₈
x ₂	1.5																	
x ₃	1.4	1.6																
x ₄	1.6	1.4	1.3															
x ₅	1.7	1.4	1.5	1.5														
x ₆	1.3	1.4	1.4	1.5	1.4													
x ₇	1.6	1.3	1.4	1.4	1.5	1.8												
x ₈	1.5	1.4	1.6	1.3	1.7	1.6	1.4											
x ₉	1.4	1.3	1.4	1.5	1.2	1.4	1.3	1.5										
x ₁₀	2.3	2.4	2.5	2.3	2.6	2.7	2.8	2.7	3.1									
x ₁₁	2.9	2.8	2.9	3.0	2.9	3.1	2.9	3.1	3.0	1.5								
x ₁₂	3.2	3.3	3.2	3.1	3.3	3.4	3.3	3.4	3.5	3.5	1.6							
x ₁₃	3.3	3.4	3.2	3.2	3.3	3.4	3.2	3.3	3.5	3.6	1.4	1.7						
x ₁₄	3.4	3.2	3.5	3.4	3.7	3.5	3.6	3.3	3.5	3.8	1.5	1.8	0.5					
x ₁₅	4.2	4.1	4.1	4.1	4.1	4.1	4.1	4.1	4.1	4.1	1.7	1.6	0.3	0.5				
x ₁₆	4.1	4.1	4.1	4.1	4.1	4.1	4.1	4.1	4.1	4.1	1.6	1.5	0.4	0.5	0.4			
x ₁₇	5.9	6.2	6.2	5.8	6.1	6.0	6.1	5.9	5.8	6.0	2.3	2.3	2.5	2.3	2.4	2.5		
x ₁₈	6.1	6.3	6.2	5.8	6.1	6.0	6.1	5.9	5.8	6.0	3.1	2.7	2.6	2.3	2.5	2.6	3.0	
x ₁₉	6.0	6.1	6.2	5.8	6.1	6.0	6.1	5.9	5.8	6.0	3.0	2.9	2.7	2.4	2.5	2.8	3.1	0.4

$k=3$

x_{11}

$\Rightarrow 15$

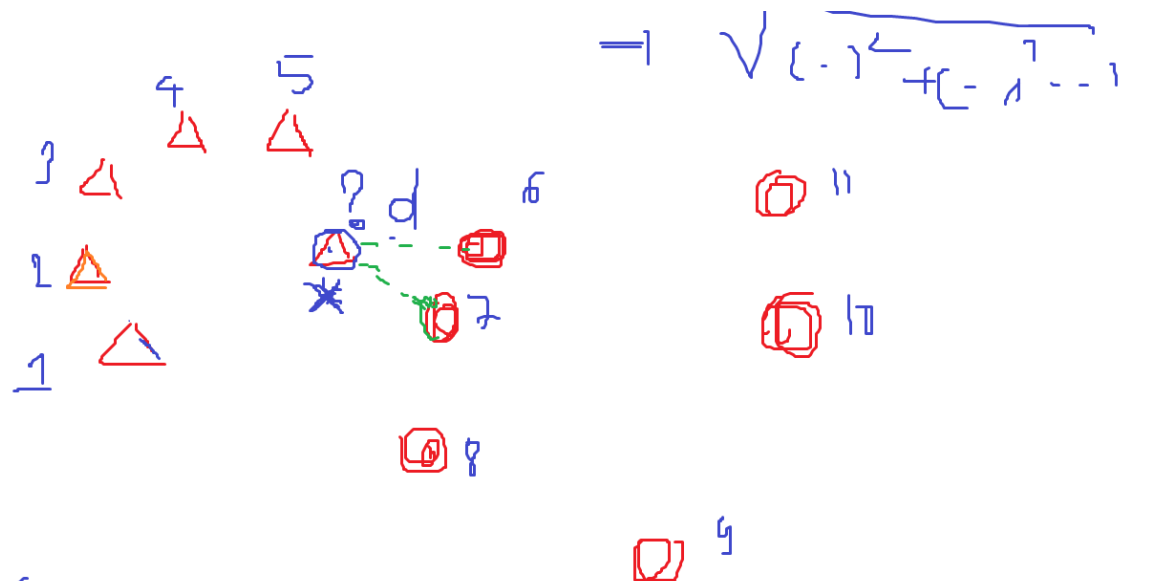
$$d(x_{\cancel{11}}, x_9)$$

$$= \sqrt{(x_{\cancel{11}} - x_9)^2 + (\dots)^2 + \dots}$$

KNN

KNN

(P, T, E)



$$= \sqrt{(-1)^2 + (-1)^2 + \dots}$$

$$K = 7 \quad ?$$

$$K = 6$$

$$\frac{A}{3} / \frac{B}{3} \quad ?$$

$$K = 1 \rightarrow d(x, 7) \text{ min} \Rightarrow \begin{cases} C_7 = \text{tròn} \\ \Rightarrow x \Rightarrow \text{tròn} \end{cases}$$

$$K = 3 \Rightarrow \left. \begin{array}{l} d(x, 7) \rightarrow T_n \\ d(x, 6) \rightarrow T_n \\ d(x, 5) \rightarrow T_n \end{array} \right\} \Rightarrow C_{765} \Rightarrow T_n$$

→ Tiêu chí phụ