

TRƯỜNG ĐẠI HỌC THỦY LỢI
KHOA CÔNG NGHỆ THÔNG TIN



ĐỀ TÀI:
KHAI PHÁ DỮ LIỆU NHẬN BIẾT TỈ LỆ MẮC UNG
THƯ PHÔI

Nhóm sinh viên thực hiện:

1. Phan Thanh Toàn – 185P1063393 – 60TH
2. Nguyễn Hồng Sơn – 185P1063376 – 60TH
3. Nguyễn Ngọc Anh – 185P1063479 – 60TH

Giảng viên hướng dẫn: Giảng viên Trần Mạnh Tuấn

Hà Nội, 03/2022

LỜI MỞ ĐẦU

Trong những năm gần đây cùng với phát triển nhanh chóng của khoa học kỹ thuật là sự bùng nổ về tri thức. Kho dữ liệu, nguồn tri thức của nhân loại cũng trở nên đồ sộ, vô tận làm cho vấn đề khai thác các nguồn tri thức đó ngày càng trở nên nóng bỏng và đặt ra thách thức lớn cho nền công nghệ thông tin thế giới.

Nhu cầu về tìm kiếm và xử lý thông tin, cùng với yêu cầu về khả năng kịp thời khai thác chúng để mang lại những năng suất và chất lượng cho công tác quản lý, hoạt động kinh doanh... đã trở nên cấp thiết trong xã hội hiện đại. Để đáp ứng phần nào yêu cầu này, người ta đã xây dựng các công cụ tìm kiếm và xử lý thông tin nhằm giúp cho người dùng tìm kiếm được các thông tin cần thiết cho mình.

Vì vậy nhóm em chọn đề tài: “***Khai phá dữ liệu nhận biết tỉ lệ mắc ung thư phổi***”, để làm báo cáo môn học của mình

Báo cáo gồm 5 chương:

Chương 1: Tổng quan về khai phá dữ liệu.

Chương 2: Tiền xử lý dữ liệu.

Chương 3: Khai phá dữ liệu bằng thuật toán K-means clustering

Chương 4: Khai phá dữ liệu bằng thuật toán phân lớp

Chương 5: Kết luận và hướng phát triển

CHƯƠNG 1: TỔNG QUAN VỀ KHAI PHÁ DỮ LIỆU

1.1 Phát hiện tri thức và khai phá dữ liệu.

Phát hiện tri thức (*Knowledge Discovery*) trong các cơ sở dữ liệu là một quy trình nhận biết các mẫu hoặc các mô hình trong dữ liệu với các tính năng: hợp thức, mới, khả ích, và có thể hiểu được.

Khai phá dữ liệu (*Data mining*) được định nghĩa như sau: “*Data mining là một quá trình tìm kiếm, phát hiện các tri thức mới, tiềm ẩn, hữu dụng trong CSDL lớn*”.

Khai phá dữ liệu có thể được sử dụng cho các lĩnh vực y tế, phân tích thị trường, xây dựng ... có thể được xem như là kết quả của sự tiến triển tự nhiên của công nghệ thông tin.

1.2 Quy trình khám phá tri thức trong CSDL

Quá trình phát hiện tri thức bao gồm các bước:

- **Làm sạch dữ liệu:** Các nhiễu và dữ liệu không nhất quán sẽ được loại bỏ.
- **Tích hợp dữ liệu:** Dữ liệu từ nhiều nguồn khác có thể được tổ hợp lại.
- **Lựa chọn dữ liệu:** Những dữ liệu thích hợp với nhiệm vụ phân tích sẽ được trích rút ra từ CSDL.
- **Chuyển đổi dữ liệu:** Dữ liệu sau khi được chọn lọc sẽ được chuyển đổi hay hợp nhất về dạng thích hợp cho việc khai phá.
- **Khai phá dữ liệu:** Quá trình cốt lõi, tất yếu trong đó các phương pháp thông minh sẽ được áp dụng nhằm trích rút ra các mẫu dữ liệu.
- **Đánh giá mẫu:** Các nhà phân tích dữ liệu sẽ dựa trên một số độ đo nào đó để xác định lợi ích thực sự, độ quan trọng của các mẫu biểu diễn tri thức.
- **Biểu diễn tri thức:** Giai đoạn này các kỹ thuật biểu diễn và hiển thị tri thức sẽ được sử dụng để đưa tri thức đã lấy ra đến người dùng.

1.3 Mô tả bài toán chuẩn đoán tỉ lệ.

1.3.1 Tổng quan bài toán.

Dataset gồm các mô tả về các thuộc tính tương ứng với chẩn đoán ung thư phổi. Áp dụng các thuật toán để xác định xem đối tượng mắc ung thư phổi: YES, NO

1.3.2 Phân tích dữ liệu thô.

Nguồn dữ liệu thô:

https://www.kaggle.com/datasets/mysarahmadbhat/lung-cancer?fbclid=IwAR0Ap29_cgYPW-aaoEuL4-QXDpYSHBbhP7fzpzPIJlCqwsBffg1wdHDiTRk

+ *Hiểu dữ liệu*: Dữ liệu dự đoán mắc ung thư phổi dựa trên giá trị các dữ liệu thuộc tính.

+ *Dữ liệu gồm*: Dữ liệu bao gồm 309 bản ghi cùng 16 thuộc tính chẩn đoán ung thư phổi.

GENDER	AGE	SMOKING	YELLOW_FINGERS	ANXIETY	PEER_PRESSURE	CHRONIC_DISEASE	FATIGUE	ALLERGY	WHEEZING	ALCOHOL_CONSUMING	COUGHING	SHORTNESS_OF_BREATH	SWALLOWING_DIFFICULTY	CHEST_PAIN	LUNG_CANCER
0	69	1	2	2	1	1	2	1	2	2	2	2	2	2	YES
0		2	1	1	1	2	2	2	1	1	1	2	2	2	YES
1	59	1	1	1	2	1	2	1	2	1	2	2	1	2	NO
0	63	2	2	2	1	1	1	1	1	2	1	1	2	2	NO
1	63	1	2	1	1	1	1	1	2	1	2	2	1	1	NO
1	75	1	2	1	1	2	2	2	2	1	2	2	1	1	YES
0	52	2	1	1	1	1	2	1	2	2	2	2	1	2	YES
1	51	2	2	2	2	1	2	2	1	1	1	2	2	1	YES
1		2	1	2	1	1	2	1	1	1	1	1	1	1	NO
0	53	2	2	2	2	2	1	2	1	2	1	1	2	2	YES
1	61	2	2	2	2	2	2	1	2	1	2	2	2	1	YES
0	72	1	1	1	1	2	2	2	2	2	2	2	1	2	YES
1	60	2	1	1	1	1	2	1	1	1	1	2	1	1	NO

Hình 1.2 Dữ liệu ban đầu.

Hiểu các thuộc tính:

STT	Thuộc tính	Ý nghĩa thuộc tính
1	GENDER	Giới tính của đối tượng.
2	AGE	Tuổi của đối tượng.
3	SMOKING	Thói quen hút thuốc của đối tượng.
4	YELLOW_FINGERS	Màu sắc ngón tay của đối tượng.
5	ANXIETY	Trạng thái lo âu của đối tượng.
6	PEER PRESSURE	Áp lực cuộc sống của đối tượng.
7	CHRONIC DISEASE	Tình trạng mãn tính của đối tượng.
8	FATIGUE	Sự mệt mỏi của đối tượng.
9	ALLERGY	Dị ứng của đối tượng.
10	WHEEZING	Hơi thở khò khè của của đối tượng.
11	ALCOHOL CONSUMING	Mức độ tiêu thụ rượu của đối tượng.
12	COUGHING	Tần suất ho của đối tượng.
13	SHORTNESS OF BREATH	Tần suất thở gấp của đối tượng.
14	SWALLOWING DIFFICULTY	Sự khó khăn khi nuốt của đối tượng.
15	CHEST PAIN	Sự tức ngực của đối tượng.
16	LUNG _ CANCER	Phân loại xem đối tượng mắc ung thư hoặc không. (YES hoặc NO)

CHƯƠNG 2: TIỀN XỬ LÝ DỮ LIỆU

2.1 Làm sạch dữ liệu.

Là quá trình nhận dạng dữ liệu đã có để tiến hành xử lý các dữ liệu bị thiếu (missing data) xử lý dữ liệu bị nhiễu (noisy data) và không nhất quán.

(1) Xử lý dữ liệu bị thiếu (missing data)

(2) Xử lý dữ liệu, không nhất quán (inconsistent data).

Thực hiện:

- Xử lý trên excel:
 - + Xử lý nhất quán dữ liệu.
 - + Tiền xử lý dữ liệu trên Weka.
- Đưa dữ liệu vào weka:

Viewer													
Relation: survey lung cancer - survey lung cancer													
No.	1: GENDER Numeric	2: AGE Numeric	3: SMOKING Numeric	4: YELLOW_FINGERS Numeric	5: ANXIETY Numeric	6: PEER_PRESSURE Numeric	7: CHRONIC_DISEASE Numeric	8: FATIGUE Numeric	9: ALLERGY Numeric	10: WHEEZING Numeric	11: ALCOHOL_CONSUMING Numeric	12: COUGHING Numeric	13: SHORT Numeric
1	0.0	69.0	1.0	2.0	2.0	1.0	1.0	2.0	1.0	2.0		2.0	2.0
2	0.0		2.0	1.0	1.0	1.0		2.0	2.0	2.0	1.0	1.0	1.0
3	1.0	59.0	1.0	1.0	1.0	2.0	1.0	2.0	1.0	2.0		1.0	2.0
4	0.0	63.0	2.0	2.0	2.0	1.0	1.0	1.0	1.0	1.0		2.0	1.0
5	1.0	63.0	1.0	2.0	1.0	1.0	1.0	1.0	1.0	2.0		1.0	2.0
6	1.0	75.0	1.0	2.0	1.0	1.0	2.0	2.0	2.0	2.0		1.0	2.0
7	0.0	52.0	2.0	1.0	1.0	1.0	1.0	2.0	1.0	2.0		2.0	2.0
8	1.0	51.0	2.0	2.0	2.0	2.0	1.0	2.0	2.0	1.0		1.0	1.0
9	1.0		2.0	1.0	2.0	1.0	1.0	2.0	1.0	1.0		1.0	1.0
10	0.0	53.0	2.0	2.0	2.0	2.0	2.0	1.0	2.0	1.0		2.0	1.0
11	1.0	61.0	2.0	2.0	2.0	2.0	2.0	2.0	1.0	2.0		1.0	2.0
12	0.0	72.0	1.0	1.0	1.0	1.0	2.0	2.0	2.0	2.0		2.0	2.0
13	1.0	60.0	2.0	1.0	1.0	1.0	1.0	2.0	1.0	1.0		1.0	1.0
14	0.0	58.0	2.0	1.0	1.0	1.0	1.0	2.0	2.0	2.0		2.0	2.0
15	0.0	69.0	2.0	1.0	1.0	1.0	1.0	1.0	2.0	2.0		2.0	2.0
16	1.0	48.0	1.0	2.0	2.0	2.0	2.0	2.0	2.0	2.0		1.0	2.0
17	0.0	75.0	2.0	1.0	1.0	1.0	2.0	1.0	2.0	2.0		2.0	2.0
18	0.0	57.0	2.0	2.0	2.0	2.0	2.0	1.0	1.0	1.0		2.0	1.0
19	1.0	68.0			2.0	2.0	2.0	2.0	1.0	1.0		1.0	2.0
20	1.0		1.0	1.0	1.0	1.0	2.0	2.0	1.0	1.0		1.0	1.0
21	1.0	44.0	2.0	2.0	2.0	2.0	2.0	2.0	1.0	1.0		1.0	1.0
22	1.0	64.0	1.0	2.0	2.0	2.0	1.0	1.0	2.0	2.0		1.0	2.0
23	1.0	21.0	2.0	1.0	1.0	1.0	2.0	2.0	2.0	1.0		1.0	1.0

+ Các dữ liệu thiếu thay thế bằng giá trị trung bình của thuộc tính dùng bộ lọc ReplaceMissingValue.

Viewer

Relation: survey_lung_cancer - survey_lung_cancer-weka.filters.unsupervised.attribute.ReplaceMissingValues

No.	1: GENDER	2: AGE	3: SMOKING	4: YELLOW_FINGERS	5: ANXIETY	6: PEER_PRESSURE	7: CHRONIC_DISEASE	8: FATIGUE	9: ALLERGY	10: WHEEZING	11: ALCOHOL	Nur
	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Nur
1	0.0	69.0	1.0	2.0	2.0	1.0	1.0	2.0	1.0	2.0	2.0	
2	0.0	62.6...	2.0	1.0	1.0	1.0	2.0	2.0	2.0	1.0	1.0	
3	1.0	59.0	1.0	1.0	1.0	2.0	1.0	2.0	1.0	2.0	2.0	
4	0.0	63.0	2.0	2.0	2.0	1.0	1.0	1.0	1.0	1.0	1.0	
5	1.0	63.0	1.0	2.0	1.0	1.0	1.0	1.0	1.0	2.0	2.0	
6	1.0	75.0	1.0	2.0	1.0	1.0	2.0	2.0	2.0	2.0	2.0	
7	0.0	52.0	2.0	1.0	1.0	1.0	1.0	2.0	1.0	2.0	2.0	
8	1.0	51.0	2.0	2.0	2.0	2.0	1.0	2.0	2.0	1.0	1.0	
9	1.0	62.6...	2.0	1.0	2.0	1.0	1.0	2.0	1.0	1.0	1.0	
10	0.0	53.0	2.0	2.0	2.0	2.0	2.0	1.0	2.0	1.0	1.0	
11	1.0	61.0	2.0	2.0	2.0	2.0	2.0	2.0	1.0	2.0	2.0	
12	0.0	72.0	1.0	1.0	1.0	1.0	2.0	2.0	2.0	2.0	2.0	
13	1.0	60.0	2.0	1.0	1.0	1.0	1.0	2.0	1.0	1.0	1.0	
14	0.0	58.0	2.0	1.0	1.0	1.0	1.0	2.0	2.0	2.0	2.0	
15	0.0	69.0	2.0	1.0	1.0	1.0	1.0	1.0	2.0	2.0	2.0	
16	1.0	48.0	1.0	2.0	2.0	2.0	2.0	2.0	2.0	2.0	2.0	
17	0.0	75.0	2.0	1.0	1.0	1.0	2.0	1.0	2.0	2.0	2.0	
18	0.0	57.0	2.0	2.0	2.0	2.0	2.0	1.0	1.0	1.0	1.0	
19	1.0	68.0	2.0	2.0	2.0	2.0	2.0	2.0	1.0	1.0	1.0	
20	1.0	62.6...	1.0	1.0	1.0	1.0	2.0	2.0	1.0	1.0	1.0	
21	1.0	44.0	2.0	2.0	2.0	2.0	2.0	2.0	1.0	1.0	1.0	
22	1.0	64.0	1.0	2.0	2.0	2.0	1.0	1.0	2.0	2.0	2.0	
23	1.0	21.0	2.0	1.0	1.0	1.0	2.0	2.0	2.0	1.0	1.0	
24	0.0	60.0	2.0	1.0	1.0	1.0	1.0	2.0	2.0	2.0	2.0	
25	0.0	72.0	2.0	2.0	2.0	2.0	2.0	1.0	2.0	2.0	2.0	
26	0.0	65.0	1.0	2.0	2.0	1.0	1.0	2.0	1.0	2.0	2.0	
27	1.0	61.0	2.0	2.0	2.0	1.0	1.0	2.0	2.0	1.0	1.0	
28	0.0	69.0	1.0	1.0	1.0	2.0	1.0	2.0	1.0	2.0	2.0	
29	1.0	53.0	2.0	2.0	2.0	1.0	2.0	1.0	1.0	2.0	2.0	

2.2 Tích hợp dữ liệu.

Là quá trình trộn dữ liệu từ các nguồn khác nhau vào một kho dữ liệu sẵn sàng cho quá trình khai phá dữ liệu.

- (1) Tích hợp lược đồ và so trùng đối tượng.
- (2) Vấn đề dư thừa.
- (3) Phát hiện và xử lý mâu thuẫn giá trị dữ liệu.

=> Dữ liệu lấy từ một nguồn nên không cần thực hiện quá trình này.

2.3 Biến đổi dữ liệu.

Là quá trình biến đổi hay kết hợp dữ liệu vào những dạng thích hợp cho quá trình khai phá dữ liệu.

Các chiến lược thu giảm:

- + Làm trơn dữ liệu.
- + Kết hợp dữ liệu.
- + Tổng quát hóa dữ liệu.
- + Chuẩn hóa dữ liệu.
- + Xây dựng thuộc tính đặc tính.

⇒ Thực hiện chuẩn hóa dữ liệu:

$$v' = \frac{v - \min}{\max - \min} (\text{new_max} - \text{new_min}) + \text{new_min}$$

Trong đó: $v = [\min A, \max A]$ là giá trị cũ.

$v' = [0, 1]$ là giá trị mới.

Ví dụ: Thuộc tính AGE (tuổi): Có $v = [21, 87]$.

Chuẩn hóa về giá trị: $v' = [0, 1]$.

Với $v = 69$.

$$\Rightarrow v' = (69 - 21) / (87 - 21) * (1 - 0) + 0 = 0.727273$$

Tiến hành chuẩn hóa các thuộc tính số về đoạn $[0, 1]$ bằng phương pháp chuẩn hóa min-max bằng bộ lọc Normalize, lưu file lại dưới định dạng csv.

Viewer

Relation: survey_lung_cancer - survey_lung_cancer-weka.filters.unsupervised.attribute.ReplaceMissingValues-weka.filters.unsupervised.at

No.	1: GENDER Numeric	2: AGE Numeric	3: SMOKING Numeric	4: YELLOW_FINGERS Numeric	5: ANXIETY Numeric	6: PEER_PRESSURE Numeric	7: CHRONIC_DISEASE Numeric
1	0.0	0.7272727272727273	0.0	1.0	1.0	0.0	0.0
2	0.0	0.6305274971941639	1.0	0.0	0.0	0.0	1.0
3	1.0	0.5757575757575758	0.0	0.0	0.0	1.0	0.0
4	0.0	0.6363636363636364	1.0	1.0	1.0	0.0	0.0
5	1.0	0.6363636363636364	0.0	1.0	0.0	0.0	0.0
6	1.0	0.8181818181818182	0.0	1.0	0.0	0.0	1.0
7	0.0	0.4696969696969697	1.0	0.0	0.0	0.0	0.0
8	1.0	0.4545454545454545...	1.0	1.0	1.0	1.0	0.0
9	1.0	0.6305274971941639	1.0	0.0	1.0	0.0	0.0
10	0.0	0.484848484848484...	1.0	1.0	1.0	1.0	1.0
11	1.0	0.6060606060606061	1.0	1.0	1.0	1.0	1.0
12	0.0	0.7727272727272727	0.0	0.0	0.0	0.0	1.0
13	1.0	0.5909090909090909	1.0	0.0	0.0	0.0	0.0
14	0.0	0.5606060606060606	1.0	0.0	0.0	0.0	0.0
15	0.0	0.7272727272727273	1.0	0.0	0.0	0.0	0.0
16	1.0	0.4090909090909091	0.0	1.0	1.0	1.0	1.0
17	0.0	0.8181818181818182	1.0	0.0	0.0	0.0	1.0
18	0.0	0.5454545454545454	1.0	1.0	1.0	1.0	1.0
19	1.0	0.7121212121212122	1.0	1.0	1.0	1.0	1.0
20	1.0	0.6305274971941639	0.0	0.0	0.0	0.0	1.0
21	1.0	0.3484848484848485	1.0	1.0	1.0	1.0	1.0
22	1.0	0.6515151515151515	0.0	1.0	1.0	1.0	0.0

Right click (or left+alt) for context menu

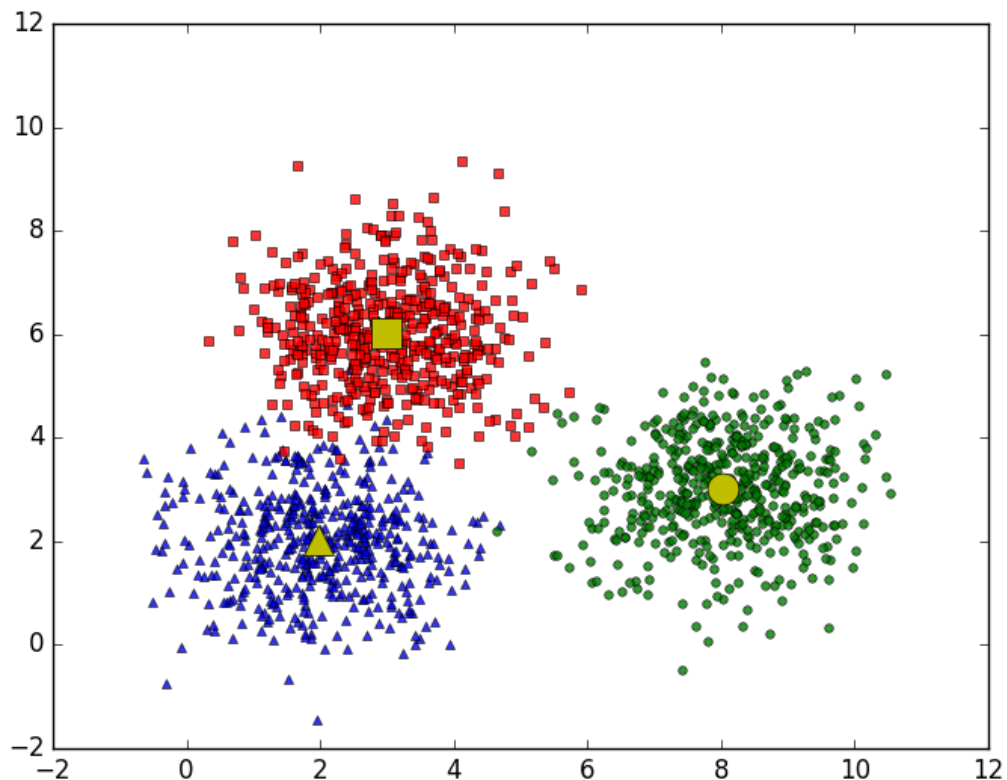
CHƯƠNG 3: KHAI PHÁ DỮ LIỆU BẰNG THUẬT TOÁN K-MEANS CLUSTERING

3.1 Giới thiệu về kỹ thuật K-means clustering.

- **Khái niệm:**

+ Trong thuật toán K-means clustering, chúng ta không biết nhãn (label) của từng điểm dữ liệu. Mục đích là làm thế nào để phân dữ liệu thành các cụm (cluster) khác nhau sao cho dữ liệu trong cùng một cụm có tính chất giống nhau.

+ Giả sử mỗi cluster có một điểm đại diện (center) màu vàng. Và những điểm xung quanh mỗi center thuộc vào cùng nhóm với center đó. Một cách đơn giản nhất, xét một điểm bất kỳ, ta xét xem điểm đó gần với center nào nhất thì nó thuộc về cùng nhóm với center đó.



Bài toán với 3 clusters.

- **Mô hình và phương trình:**

Thuật toán tối ưu hàm mất mát

Đây là một bài toán khó tìm điểm tối ưu vì nó có thêm các điều kiện ràng buộc. Bài toán này thuộc loại mix - integer programming (điều kiện biến là số nguyên) - là loại rất khó tìm nghiệm tối ưu toàn cục (global optimal point, tức nghiệm làm cho hàm mất mát đạt giá trị nhỏ nhất có thể). Tuy nhiên, trong một số trường hợp chúng ta vẫn có thể tìm được phương pháp để tìm được nghiệm gần đúng hoặc điểm cực tiểu. (Nếu chúng ta vẫn nhớ chương trình toán ôn thi đại học thì điểm cực tiểu chưa chắc đã phải là điểm làm cho hàm số đạt giá trị nhỏ nhất).

Giả sử đã tìm được các centers, hãy tìm các label vector để hàm mất mát đạt giá trị nhỏ nhất. Điều này tương đương với việc tìm cluster cho mỗi điểm dữ liệu.

Khi các centers là cố định, bài toán tìm label vector cho toàn bộ dữ liệu có thể được chia nhỏ thành bài toán tìm label vector cho từng điểm dữ liệu xi như sau:

$$\mathbf{y}_i = \arg \min_{\mathbf{y}_i} \sum_{j=1}^K y_{ij} \|\mathbf{x}_i - \mathbf{m}_j\|_2^2 \quad (3)$$

$$\text{subject to: } y_{ij} \in \{0, 1\} \quad \forall j; \quad \sum_{j=1}^K y_{ij} = 1$$

$$j = \arg \min_j \|\mathbf{x}_i - \mathbf{m}_j\|_2^2$$

Vì $\|\mathbf{x}_i - \mathbf{m}_j\|_2^2$ chính là bình phương khoảng cách tính từ điểm xi tới center \mathbf{m}_j ta có thể kết luận rằng mỗi điểm xi thuộc vào cluster có center gần nó nhất! Từ đó ta có thể dễ dàng suy ra label vector của từng điểm dữ liệu.

Giả sử đã tìm được cluster cho từng điểm, hãy tìm center mới cho mỗi cluster để hàm mất mát đạt giá trị nhỏ nhất.

Một khi chúng ta đã xác định được label vector cho từng điểm dữ liệu, bài toán tìm center cho mỗi cluster được rút gọn thành:

$$\mathbf{m}_j = \arg \min_{\mathbf{m}_j} \sum_{i=1}^N y_{ij} \|\mathbf{x}_i - \mathbf{m}_j\|_2^2.$$

Tới đây, ta có thể tìm nghiệm bằng phương pháp giải đạo hàm bằng 0, vì hàm cần tối ưu là một hàm liên tục và có đạo hàm xác định tại mọi điểm. Và quan trọng hơn, hàm này là hàm convex (lồi) theo \mathbf{m}_j nên chúng ta sẽ tìm được giá trị nhỏ nhất và điểm tối ưu tương ứng.

Đặt $l(\mathbf{m}_j)$ là hàm bên trong dấu argmin ta có đạo hàm:

$$\frac{\partial l(\mathbf{m}_j)}{\partial \mathbf{m}_j} = 2 \sum_{i=1}^N y_{ij} (\mathbf{m}_j - \mathbf{x}_i)$$

Giải phương trình đạo hàm bằng 0 ta có:

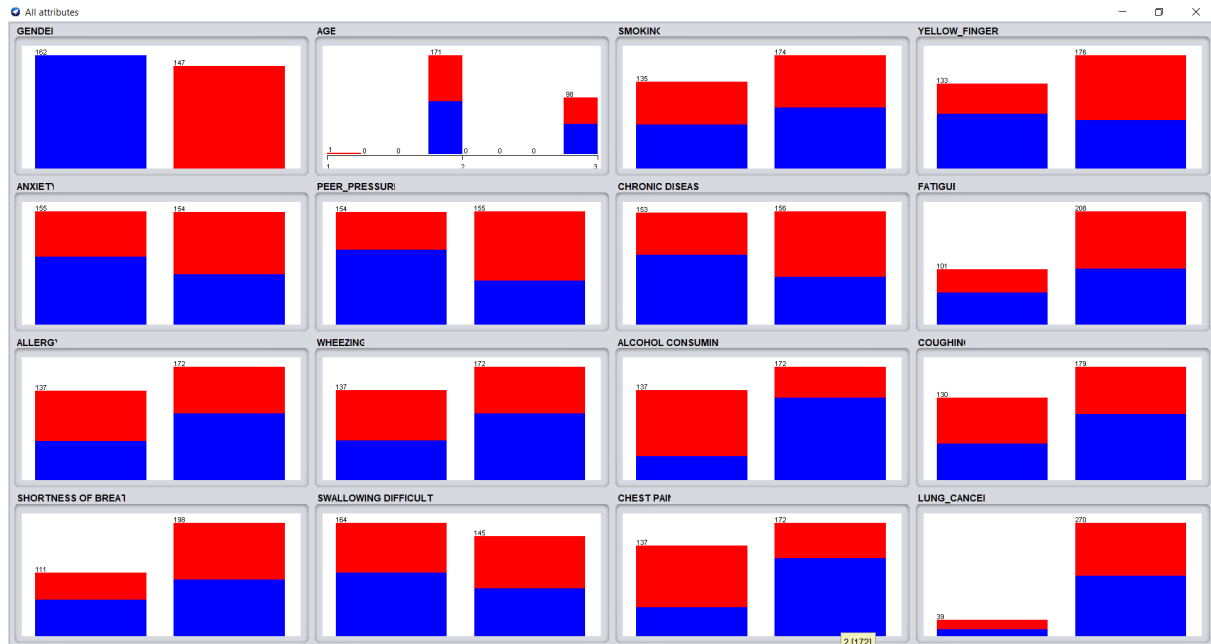
$$\begin{aligned} \mathbf{m}_j \sum_{i=1}^N y_{ij} &= \sum_{i=1}^N y_{ij} \mathbf{x}_i \\ \Rightarrow \mathbf{m}_j &= \frac{\sum_{i=1}^N y_{ij} \mathbf{x}_i}{\sum_{i=1}^N y_{ij}} \end{aligned}$$

\mathbf{m}_j là trung bình cộng của các điểm trong cluster j .

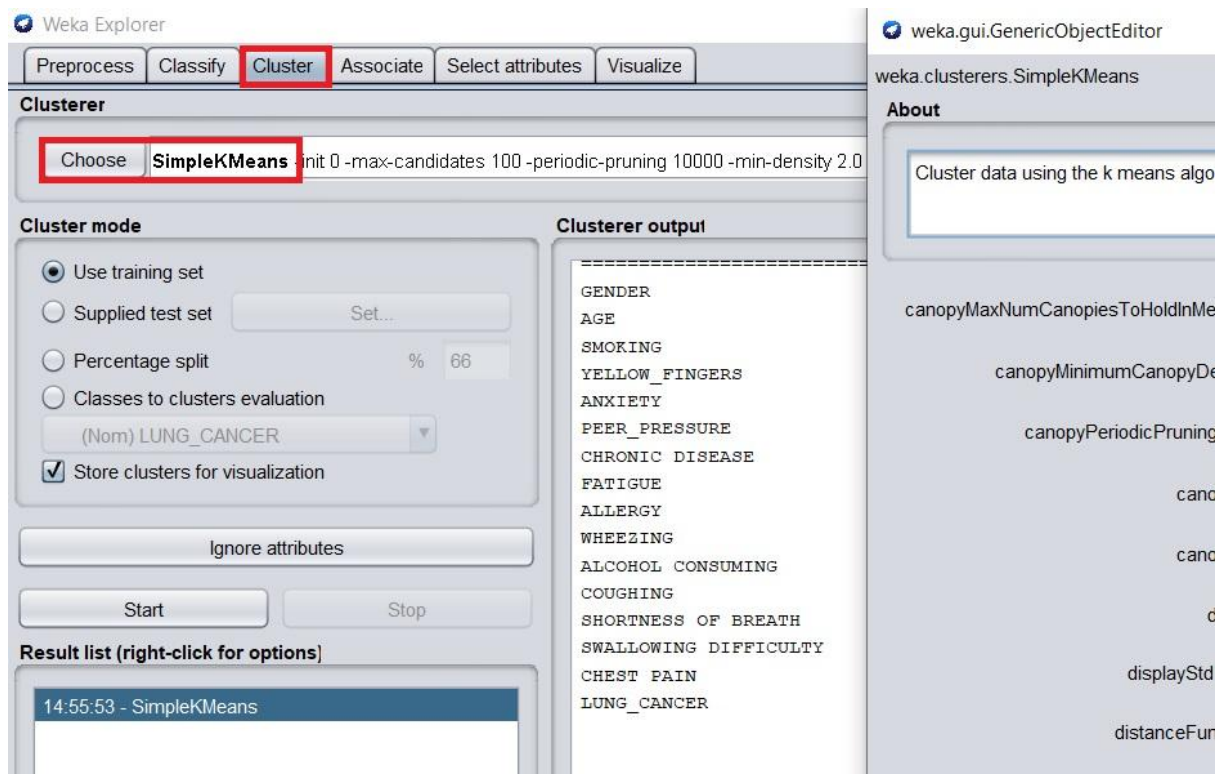
3.2 Thuật toán phân cụm K-means clustering

- Các bước thực hiện:

Bước 1: Đọc dữ liệu vào weka. Dữ liệu trong dataset như sau:



Bước 2: Trên giao diện weka chọn Cluster và chọn SimpleKMeans.



Bước 3: Tiến hành đánh giá hiệu quả của thuật toán đối với tập dữ liệu được dùng theo phương pháp:

Phương pháp: Theo Percentage split.

- Cho biết tỷ lệ phân chia là bao nhiêu phần trăm thì đạt hiệu quả phân lớp cao nhất.

(2) Percentage split:

Train/Test: 70/30%.

The screenshot shows the WEKA Clusterer interface. In the 'Cluster mode' section, 'Percentage split' is selected with a value of 70%. The 'Start' button is highlighted. The 'Clusterer output' section displays a table of attributes and their cluster assignments, along with a summary of clustered instances.

Attribute	Full Data (216.0)	Cluster# (60.0)
GENDER	0	0
AGE	2.3211	2.3309
SMOKING	2	2
YELLOW_FINGERS	2	1
ANXIETY	2	1
PEER_PRESSURE	1	1
CHRONIC_DISEASE	1	1
FATIGUE	2	2
ALLERGY	2	2
WHEEZING	2	2
ALCOHOL_CONSUMING	2	2
COUGHING	2	2
SHORTNESS_OF_BREATH	2	2
SWALLOWING_DIFFICULTY	1	1
CHEST_PAIN	2	2
LUNG_CANCER	1	1

Time taken to build model (percentage split) : {

Clustered Instances

Cluster	Count	Percentage
0	29	(31%)
1	24	(26%)
2	23	(25%)
3	8	(9%)
4	9	(10%)

Kết quả chạy:

=== Run information ===

Scheme: weka.clusterers.SimpleKMeans -init 0 -max-candidates 100

-periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 5 -A

"weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -S 10

Relation: survey lung cancer - survey lung

cancer-weka.filters.unsupervised.attribute.MathExpression-Eifelse(A<30,1,ifels

e(A<65,2,3))-V-R2-weka.filters.unsupervised.attribute.NumericToNominal-R2-

V-weka.filters.unsupervised.attribute.ReplaceMissingValues-weka.filters.unsupe

rvised.attribute.Normalize-S1.0-T0.0-weka.filters.unsupervised.attribute.Replac

eMissingValues

Instances: 309

Attributes: 16

GENDER

AGE

SMOKING

YELLOW_FINGERS

ANXIETY

PEER_PRESSURE

CHRONIC DISEASE

FATIGUE

ALLERGY

WHEEZING

ALCOHOL CONSUMING

COUGHING

SHORTNESS OF BREATH

SWALLOWING DIFFICULTY

CHEST PAIN

LUNG _ CANCER

Test mode: evaluate on training data

=== Clustering model (full training set) ===

kMeans

=====

Number of iterations: 6

Within cluster sum of squared errors: 1033.5341863729307

Initial starting points (random):

Cluster 0: 1,1,2,2,2,2,2,1,1,1,1,2,2,1,1

Cluster 1: 0,0.5,1,2,1,1,2,1,1,2,1,2,2,1,2,1

Cluster 2: 1,1,2,2,2,2,1,2,1,2,1,2,2,2,1,1

Cluster 3: 1,1,2,1,1,2,1,2,2,2,2,2,1,2,2,1

Cluster 4: 1,0.67963,2,1,2,1,1,2,1,1,1,1,1,1,0

Missing values globally replaced with mean/mode

Final cluster centroids:

	Cluster#						
Attribute	Full Data	0	1	2	3	4	
	(309.0)	(56.0)	(68.0)	(73.0)	(71.0)	(41.0)	
=====							
=====							
GENDER	0	1	0	1	0	1	
AGE	0.6796	0.6442	0.737	0.6972	0.6732	0.6126	
SMOKING	2	2	1	2	2	2	
YELLOW_FINGERS		2	2	2	2	1	1
ANXIETY	1	2	1	2	1	1	
PEER_PRESSURE		2	2	1	2	1	1
CHRONIC DISEASE		2	2	2	1	1	1
FATIGUE	2	1	2	2	2	2	
ALLERGY	2	1	2	1	2	1	
WHEEZING	2	1	2	2	2	1	
ALCOHOL CONSUMING		2	2	2	1	2	1

COUGHING	2	1	2	2	2	1	
SHORTNESS OF BREATH		2	1	2	2	2	2
SWALLOWING DIFFICULTY			1	2	1	2	1 1
CHEST PAIN	2	2	2	1	2	1	
LUNG _ CANCER		1	1	1	1	1	0

Time taken to build model (full training data) : 0.02 seconds

=== Model and evaluation on training set ===

Clustered Instances

0 56 (18%)
1 68 (22%)
2 73 (24%)
3 71 (23%)
4 41 (13%)