

**ỦY BAN NHÂN DÂN THÀNH PHỐ HỒ CHÍ MINH**  
**TRƯỜNG ĐẠI HỌC SÀI GÒN**  
**KHOA CÔNG NGHỆ THÔNG TIN**

-----◆◆◆-----



**ĐỀ TÀI ĐỒ ÁN: PHÂN TÍCH DỮ LIỆU CÁC VIDEO**  
**THỊNH HÀNH TRÊN NỀN TẢNG YOUTUBE**

<b>Sinh viên thực hiện:</b>	<b>Trịnh Vĩnh Phúc</b>
<b>Mã số sinh viên:</b>	<b>3119410318</b>
<b>Học phần:</b>	<b>Khai phá dữ liệu</b>
<b>Giảng viên hướng dẫn:</b>	<b>TS. Vũ Ngọc Thanh Sang</b>

**Thành phố Hồ Chí Minh, tháng 04 năm 2023**

## MỤC LỤC

NHẬN XÉT CỦA GIẢNG VIÊN HƯỚNG DẪN .....	iii
Danh mục các bảng.....	iv
Danh mục sơ đồ hình ảnh .....	iv
LỜI MỞ ĐẦU .....	vi
CHƯƠNG 1: GIỚI THIỆU ĐỀ TÀI .....	1
1.1. Tìm hiểu đề tài .....	1
1.2. Mục đích đề tài .....	1
1.3. Phạm vi đề tài .....	1
CHƯƠNG 2: MÔ TẢ BỘ DỮ LIỆU .....	2
2.1. Nguồn gốc dữ liệu.....	2
2.2. Kích thước bộ dữ liệu.....	2
CHƯƠNG 3: TIỀN XỬ LÝ DỮ LIỆU .....	3
3.1. Sử dụng những thư viện và tệp tin dữ liệu .....	3
3.2. Loại bỏ các trường thông tin không cần thiết.....	4
3.3. Xử lý dữ liệu thiếu .....	5
3.4. Kiểm tra các dữ liệu trùng.....	6
3.5. Xử lý dữ liệu ngoại lai.....	7
3.6. Xử lý dữ liệu nhiễu .....	10
CHƯƠNG 4: PHÂN TÍCH DỮ LIỆU KHÁM PHÁ .....	15
4.1. Phân tích đơn biến.....	15
4.2. Phân tích đa biến.....	17
CHƯƠNG 5: KHAI PHÁ DỮ LIỆU.....	20
5.1. Đặt vấn đề.....	20
5.2. Thuật toán sử dụng.....	20
CHƯƠNG 6: ĐÁNH GIÁ VÀ CHỌN THUẬT TOÁN .....	22
6.1. Áp dụng thuật toán.....	22
6.2. Tiến hành đánh giá kết quả.....	28

CHƯƠNG 7: KẾT QUẢ VÀ THẢO LUẬN.....	30
7.1. Đánh giá kết quả .....	30
7.2. Điểm mạnh của nghiên cứu .....	30
7.3. Điểm yếu của nghiên cứu .....	31
CHƯƠNG 8: KẾT LUẬN .....	32
8.1. Tổng kết kết quả.....	32
8.2. Kết luận hiệu quả .....	32
TÀI LIỆU THAM KHẢO.....	33

[illegible]

## **Danh mục các bảng**

Bảng 3.1 Liên kết với Google Drive .....	3
Bảng 3.2 Tạo đường dẫn đến thư mục chứa dữ liệu .....	3
Bảng 3.3 Vẽ biểu đồ các thuộc tính của dữ liệu.....	7
Bảng 3.4 Vẽ biểu đồ sau khi xử lý giá trị ngoại lai.....	9
Bảng 3.5 Xử lý dữ liệu nhiễu của biến views .....	11
Bảng 3.6 Xóa cột không cần thiết .....	14
Bảng 4.1 Bảng phân vị của các cột thuộc tính.....	16
Bảng 4.2 Bảng ma trận tương quan.....	17
Bảng 4.3 Bảng phân bố các ngày hình thành xu hướng của video .....	19
Bảng 6.1 Tạo tập data chứa dữ liệu cần thiết.....	22
Bảng 6.2 Chuẩn hóa giá trị trong miền từ 1 đến 10 .....	22
Bảng 6.3 Hàm random_centroids .....	23
Bảng 6.4 Hàm get_labels .....	23
Bảng 6.5 Hàm new_centroids .....	23
Bảng 6.6 Hàm plot_clusters .....	24
Bảng 6.7 Kết hợp các hàm để phân nhóm tập dữ liệu.....	24
Bảng 6.8 Dữ liệu các centroid được tạo .....	25
Bảng 6.9 Khai báo các thư viện .....	26
Bảng 6.10 Tạo biến X và y .....	26
Bảng 6.11 Tách tập training và testing .....	26
Bảng 6.12 Mô hình Support Vector Machines .....	27

## **Danh mục sơ đồ hình ảnh**

Hình 3.1 Chèn các thư viện được sử dụng.....	3
Hình 3.2 Thông tin chi tiết của bộ dữ liệu .....	4
Hình 3.3 Xóa các cột không cần thiết trong bộ dữ liệu .....	5
Hình 3.4 Hiển thị số lượng dữ liệu trống.....	5
Hình 3.5 Hiển thị lại số lượng dữ liệu trống sau khi loại bỏ thành công .....	6

Hình 3.6	Hiện thị số lượng dữ liệu trùng nhau .....	6
Hình 3.7	Xóa dữ liệu trùng lặp.....	7
Hình 3.8	Biểu đồ số lượng view của bộ dữ liệu.....	8
Hình 3.9	Biểu đồ số lượng like của bộ dữ liệu .....	8
Hình 3.10	Biểu đồ số lượng dislike của bộ dữ liệu .....	8
Hình 3.11	Biểu đồ số lượng view sau khi xử lý ngoại lai .....	9
Hình 3.12	Biểu đồ số lượng like sau khi xử lý ngoại lai .....	9
Hình 3.13	Biểu đồ số lượng dislike sau khi xử lý ngoại lai .....	10
Hình 3.14	Biểu đồ tần số cột views.....	11
Hình 3.15	Biểu đồ tần số cột views sau khi xử lý dữ liệu nhiễu.....	12
Hình 3.16	Biểu đồ tần số cột likes sau khi xử lý dữ liệu nhiễu .....	13
Hình 3.17	Biểu đồ tần số cột dislike sau khi xử lý dữ liệu nhiễu .....	13
Hình 4.1	Giá trị trung bình của các cột thuộc tính .....	15
Hình 4.2	Giá trị trung vị của các cột thuộc tính.....	15
Hình 4.3	Độ lệch chuẩn của các cột thuộc tính.....	16
Hình 4.4	Đồ thị nhiệt giữa ba thuộc tính chính của tập dữ liệu.....	18
Hình 4.5	Biểu đồ phân bố video vào các ngày xu hướng.....	19
Hình 6.1	Dữ liệu được trả về.....	22
Hình 6.2	Giá trị trả về .....	23
Hình 6.3	Đồ thị phân nhóm bằng K-means Clustering .....	25
Hình 6.4	Ví dụ minh họa việc phân nhóm thành công.....	26
Hình 6.5	Mô hình kết quả của thuật toán Support Vector Machines .....	27
Hình 6.6	Giá trị đánh giá của mô hình K-means Clustering .....	28
Hình 6.7	Giá trị đánh giá của mô hình Support Vector Machines.....	29

## LỜI MỞ ĐẦU

Công nghệ ngày càng phổ biến và không ai có thể phủ nhận được tầm quan trọng và những hiệu quả mà nó đem lại cho cuộc sống chúng ta. Sự phát triển nhanh chóng của mạng Internet đã sinh ra một khối lượng khổng lồ các dữ liệu dạng siêu văn bản. Cùng với sự thay đổi và phát triển hàng ngày, hàng giờ về nội dung cũng như số lượng các trang web trên Internet thì vấn đề tìm kiếm thông tin đối với người sử dụng lại vô cùng khó khăn. Có thể nói nhu cầu tìm kiếm thông tin trên một cơ sở dữ liệu phi cấu trúc đã được phát triển mạnh mẽ cùng với sự bành trướng của Internet. Thật vậy, với Internet, con người đã dần làm quen với các trang web cùng với vô vàn các thông tin. Trong những năm gần đây Internet đã trở thành một trong những kênh về khoa học, thông tin kinh tế, thương mại và quảng cáo chính, ảnh hưởng đến đời sống, kinh tế, giáo dục và ngay cả chính trị. Có thể nói, Internet như là cuốn từ điển Bách khoa toàn thư với thông tin đa dạng về mặt nội dung cũng như hình thức được trình bày dưới dạng văn bản, hình ảnh, âm thanh,...

Tuy nhiên, cùng với sự đa dạng và số lượng lớn thông tin như vậy, Internet đã nảy sinh vấn đề quá tải thông tin. Sự tích lũy dữ liệu này xảy ra với một tốc độ bùng nổ. Người ta ước tính rằng lượng thông tin trên toàn cầu tăng gấp đôi sau khoảng hai năm và theo đó số lượng cũng như kích cỡ của các cơ sở dữ liệu cũng tăng lên một cách nhanh chóng. Khai phá dữ liệu ra đời như một hướng giải quyết hữu hiệu cho vấn đề quá tải thông tin và xử lý thông tin nhiễu loạn. Khai phá dữ liệu là một tập hợp các kỹ thuật được sử dụng để tự động khai thác và tìm ra các mối quan hệ lẫn nhau của dữ liệu trong một tập hợp dữ liệu khổng lồ và phức tạp, đồng thời cũng tìm ra các mẫu tiềm ẩn trong tập dữ liệu đó.

Em xin cảm ơn TS. Vũ Ngọc Thanh Sang – giảng viên trường Đại học Sài Gòn khoa Công nghệ thông tin đã tạo điều kiện cho em có cơ hội thực hiện và tận tình giúp đỡ em hoàn thành dự án này, qua đó gặt hái được những kinh nghiệm và kỹ năng quý báu song hành với những kiến thức đã được học ở trường. Tuy nhiên, trong quá trình học tập, em sẽ không thể tránh khỏi sai sót nhưng sẽ cố gắng cải thiện và tiếp thu ý kiến từ mọi người để em có thêm động lực tiếp tục trên con đường chinh phục công nghệ.

# **CHƯƠNG 1: GIỚI THIỆU ĐỀ TÀI**

## **1.1. Tìm hiểu đề tài**

YouTube là một nền tảng chia sẻ video trực tuyến của Mỹ có trụ sở chính tại San Bruno, California. Nền tảng này được tạo ra vào tháng 2 năm 2005 và đã được Google mua lại vào tháng 11 năm 2006 với giá 1,65 tỷ đô la Mỹ và hiện hoạt động như một trong những công ty con của Google. YouTube là trang web được truy cập nhiều thứ hai sau Google Tìm kiếm.

Với sức ảnh hưởng và mức độ phủ sóng mạnh mẽ, YouTube đã dần chiếm lĩnh thị trường Hoa Kỳ. Tại Mỹ, việc tìm kiếm, xem - nghe video trên YouTube đã trở thành thói quen hằng ngày của rất nhiều người, đồng thời cũng không nằm ngoài xu thế chung trên thế giới. Sức ảnh hưởng này còn nằm ở bảng xếp hạng các video xu hướng nhất của YouTube. YouTube duy trì danh sách các video thịnh hành nhất xét theo từng mốc thời gian nhất định. Theo tạp chí Variety, “Để xác định các video thịnh hành nhất trong năm, YouTube sử dụng kết hợp nhiều yếu tố bao gồm đo lường tương tác của người dùng (số lượt xem, lượt chia sẻ, nhận xét và lượt thích)”.

Ở đề tài này, chúng ta sẽ tìm hiểu cách để một video trở nên thịnh hành dựa trên những thông số tương tác của người dùng trên video đó.

## **1.2. Mục đích đề tài**

Dựa vào các phương pháp khai phá và phân tích dữ liệu, chúng ta sẽ tìm ra mối tương quan giữa các yếu tố (views, likes, dislikes) trong việc xác định video nào có thể trở thành xu hướng. Chúng ta tiến hành phân tích sự liên kết và quan hệ giữa thuật toán thịnh hành của YouTube với những thông số tương tác của người dùng (lượt xem, số lượng thích, số lượng không thích,...).

## **1.3. Phạm vi đề tài**

Phạm vi của đề tài là dữ liệu của các video thịnh hành trên YouTube tại Hoa Kỳ từ ngày 14/11/2017 đến ngày 14/06/2018.

Dữ liệu bao gồm các thuộc tính cơ bản của một video trên YouTube: Chỉ mục video, ngày thịnh hành, tiêu đề, ngày xuất bản, số lượng lượt xem, số lượng yêu thích, số lượng không thích, số lượng bình luận, mô tả video.



## CHƯƠNG 2: MÔ TẢ BỘ DỮ LIỆU

### 2.1. Nguồn gốc dữ liệu

Tập dữ liệu này là bản ghi từ ngày 14/11/2017 đến ngày 14/06/2018 về các video thịnh hành nhất trên YouTube tại thị trường Hoa Kỳ. Nguồn từ website:

<https://www.kaggle.com/datasets/datasnaek/youtube-new?select=USvideos.csv>

Để có được dữ liệu trending video từ YouTube này, ta có thể sử dụng các thư viện hoặc công cụ hỗ trợ như Google API, YouTube API,... Dưới đây là các bước để crawl dữ liệu trending video từ YouTube:

- Đăng ký tài khoản Google API và YouTube API và lấy API key.
- Sử dụng API key để truy cập dữ liệu từ YouTube API, chẳng hạn như danh sách các video đang hot, thông tin về kênh, bình luận, thẻ,...
- Sử dụng Scrapy hoặc BeautifulSoup để crawl dữ liệu từ trang web của YouTube, chẳng hạn như tiêu đề video, tác giả, số lượt xem, số lượt like, số lượt dislike, thời gian đăng, thời gian cập nhật,...
- Lưu trữ dữ liệu được crawl vào một tệp hoặc cơ sở dữ liệu để phân tích và xử lý dữ liệu sau.

### 2.2. Kích thước bộ dữ liệu

Tập dữ liệu bao gồm 40949 hàng và 16 thuộc tính. Trong đó, các thuộc tính gồm có:

- *video\_id*: Chỉ mục video
- *trending\_date*: Ngày thịnh hành của video
- *title*: Tiêu đề video
- *channel\_title*: Tên kênh sở hữu
- *category\_id*: Chỉ mục danh mục
- *publish\_time*: Thời gian công khai
- *tags*: Các nhãn video
- *views*: Số lượng xem
- *likes*: Số lượng thích
- *dislikes*: Số lượng không thích
- *comment\_count*: Số lượng bình luận
- *thumbnail\_link*: Địa chỉ hình ảnh thu nhỏ
- *comments\_disabled*: Vô hiệu bình luận
- *ratings\_disabled*: Vô hiệu đánh giá
- *video\_error\_or\_removed*: Video lỗi hoặc bị gỡ
- *description*: Mô tả video

## CHƯƠNG 3: TIỀN XỬ LÝ DỮ LIỆU

### 3.1. Sử dụng những thư viện và tệp tin dữ liệu

Những thư viện được sử dụng trong dự án:

- *numpy*: bổ sung hỗ trợ cho các mảng lớn, đa chiều, cùng và một bộ các hàm toán học cấp cao.
- *pandas*: bổ sung các thao tác phân tích dữ liệu, cấu trúc dữ liệu và các phép toán để thao tác với các bảng số và chuỗi thời gian.
- *matplotlib.pyplot*: trực quan hóa dữ liệu và vẽ đồ thị.
- *seaborn*: xây dựng những hình ảnh trực quan đẹp mắt. Nó có thể được coi là một phần mở rộng của một thư viện khác có tên là Matplotlib.
- *pca*: giảm chiều và tăng khả năng trực quan hóa dữ liệu.

```
[ ] import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn
import pca
```

Hình 3.1 Chèn các thư viện được sử dụng

Trong dự án này, chúng ta sử dụng Google Colab, là một sản phẩm từ Google Research, nó cho phép chạy các dòng code python thông qua trình duyệt, đặc biệt phù hợp với Data analysis, machine learning và giáo dục. Để liên kết với Google Drive nhằm dễ dàng xử lý tập tin ngay trên Drive, ta sử dụng câu lệnh:

```
from google.colab import drive
drive.mount('/content/drive')
```

Bảng 3.1 Liên kết với Google Drive

Tạo đường dẫn đến tệp dữ liệu (“USvideos.csv”) và đọc file csv truyền vào:

```
path_data = '/content/drive/MyDrive'
df = pd.read_csv(f'{path_data}/USvideos.csv', header=0, index_col=0)
```

Bảng 3.2 Tạo đường dẫn đến thư mục chứa dữ liệu

Hiển thị thông tin các thuộc tính của tập dữ liệu:

```
[5] #chi tiết của dataframe
df.info()

<class 'pandas.core.frame.DataFrame'>
Index: 40949 entries, 2kyS6SvSYSE to ooyjaVdt-jA
Data columns (total 15 columns):
#   Column                Non-Null Count  Dtype
---  -
0   trending_date          40949 non-null  object
1   title                  40949 non-null  object
2   channel_title          40949 non-null  object
3   category_id            40949 non-null  int64
4   publish_time           40949 non-null  object
5   tags                   40949 non-null  object
6   views                  40949 non-null  int64
7   likes                  40949 non-null  int64
8   dislikes               40949 non-null  int64
9   comment_count          40949 non-null  int64
10  thumbnail_link         40949 non-null  object
11  comments_disabled      40949 non-null  bool
12  ratings_disabled       40949 non-null  bool
13  video_error_or_removed 40949 non-null  bool
14  description            40379 non-null  object
dtypes: bool(3), int64(5), object(7)
memory usage: 4.2+ MB
```

Hình 3.2 Thông tin chi tiết của bộ dữ liệu

### 3.2. Loại bỏ các trường thông tin không cần thiết

Sau khi phân tích, ta nhận thấy các trường thông tin không cần thiết:

- *category\_id*
- *thumbnail\_link*
- *comment\_count*
- *tags*
- *ratings\_disabled*
- *video\_error\_or\_removed*
- *comments\_disabled*
- *channel\_title*
- *publish\_time*

Chúng ta không sử dụng toàn bộ các biến đầu vào mà cần lọc ra những trường cần thiết cho việc phân tích và nghiên cứu. Những trường dữ liệu trên không có chức năng trong quá trình phân tích và không đóng vai trò quan trọng. Ta có thể loại bỏ chúng nhằm hoàn thiện dữ liệu, tiêu tốn ít bộ nhớ lưu trữ và giảm thời gian huấn luyện.

Việc giảm chiều dữ liệu từ không gian cao chiều xuống không gian thấp chiều như trên vẫn giữ được những đặc trưng chính của dữ liệu nhưng có thể tiết kiệm được chi phí huấn luyện và dự báo.

Thông qua đó, chúng ta nhận thấy được mức độ quan trọng của các trường dữ liệu:

- video\_id
- trending\_date
- title
- tags
- views
- likes
- dislikes

```
[ ] df = df.drop('category_id', axis='columns')
df = df.drop("thumbnail_link", axis='columns')
df = df.drop("comment_count", axis='columns')
df = df.drop("tags", axis='columns')
df = df.drop("ratings_disabled", axis='columns')
df = df.drop("video_error_or_removed", axis='columns')
df = df.drop("comments_disabled", axis='columns')
df = df.drop("channel_title", axis='columns')
df = df.drop("publish_time", axis='columns')
df.head()
```

Hình 3.3 Xóa các cột không cần thiết trong bộ dữ liệu

### 3.3. Xử lý dữ liệu thiếu

Nhằm tạo nên sự hoàn thiện của dữ liệu, ta cần tìm những dữ liệu bị thiếu. Trong tập dữ liệu trên, chúng ta có thể thấy có những dữ liệu thiếu xuất hiện dưới dạng dữ liệu trống (`null`). Cách tốt nhất để giải quyết dữ liệu dạng này là hàm `isnull()` và hàm `dropna()`.

Để xác định trong tập dữ liệu có xuất hiện dữ liệu `null` hay không, ta có thể sử dụng hàm `isnull()` để tìm kiếm:

```
[ ] df.isnull().sum()
```

```
trending_date    0
title            0
views           0
likes           0
dislikes        0
description     570
dtype: int64
```

Hình 3.4 Hiện thị số lượng dữ liệu trống

Sau khi đã tìm được những trường dữ liệu `null`, ta tiến hành loại bỏ chúng khỏi tập dữ liệu bằng hàm `dropna()`:

```
[ ] df = df.dropna()
    df.isnull().sum()
```

```
trending_date    0
title            0
views            0
likes            0
dislikes         0
description      0
dtype: int64
```

Hình 3.5 Hiển thị lại số lượng dữ liệu trống sau khi loại bỏ thành công

### 3.4. Kiểm tra các dữ liệu trùng

Trong quá trình làm sạch dữ liệu, việc loại bỏ dữ liệu trùng là vô cùng cần thiết. Việc này giúp tăng tính chính xác và hiệu quả cho việc phân tích dữ liệu. Loại bỏ dữ liệu trùng lặp cũng giảm bớt chi phí cho việc xử lý và huấn luyện sau này.

```
[ ] #Kiểm tra các dữ liệu trùng nhau
    #Có 47 dòng dữ liệu trùng nhau
    Dup_Rows = df[df.duplicated()]
    Dup_Rows.count()
```

```
trending_date    47
title            47
views            47
likes            47
dislikes         47
description      47
dtype: int64
```

Hình 3.6 Hiển thị số lượng dữ liệu trùng nhau

Sau khi sử dụng hàm `duplicated()` để tìm kiếm những hàng có dữ liệu trùng lặp với nhau, kết quả cho ta thấy trong tập dữ liệu có 47 dòng trùng với nhau. Để xử lý việc này vô cùng đơn giản, ta có thể sử dụng `drop_duplicates(keep='first')`. Trong đó, thông số `keep='first'` nhằm giữ lại dòng trùng lặp đầu tiên xóa tất cả hàng còn lại.

```
[ ] df_count = df.count()
    DF_RM_DUP = df.drop_duplicates(keep='first')
    df_count_remove_duplicate = DF_RM_DUP.count()

    print(f"Số dòng Dataframe trước loại bỏ Duplicate: {df_count}")
    print(f"Số dòng Dataframe sau loại bỏ Duplicate: {df_count_remove_duplicate}")

Số dòng Dataframe trước loại bỏ Duplicate: trending_date    40379
title    40379
views    40379
likes    40379
dislikes    40379
description    40379
dtype: int64
Số dòng Dataframe sau loại bỏ Duplicate: trending_date    40332
title    40332
views    40332
likes    40332
dislikes    40332
description    40332
dtype: int64
```

*Hình 3.7 Xóa dữ liệu trùng lặp*

Sau khi xóa bỏ dữ liệu trùng, tập dữ liệu giảm từ 40379 dòng xuống 40332 dòng. Dữ liệu đã được làm sạch tương đối, tiếp theo ta cần xử lý dữ liệu ngoại lai và dữ liệu nhiễu.

### 3.5. Xử lý dữ liệu ngoại lai

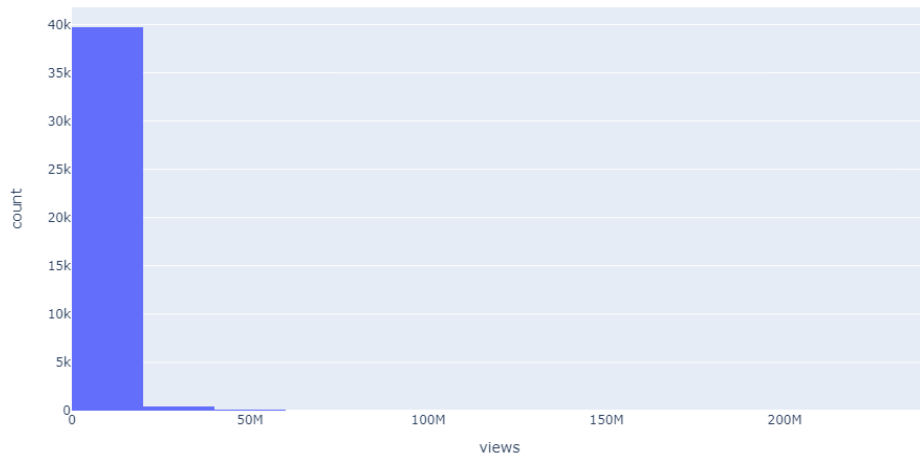
Ở quá trình này, ta sử dụng thư viện đồ họa mới plotly.express. Thư viện này giúp tạo ra các biểu đồ tương tác, đồ thị chất lượng cao trong xử lý dữ liệu. Việc xử lý giá trị ngoại lai là xác định và loại bỏ các giá trị khác xa với phần còn lại của các giá trị trong trường đó. Các giá trị này có tần xuất xảy ra vô cùng thấp trong cột dữ liệu. Đây chính là dữ liệu ngoại lai. Ta xác định những trường dữ liệu ta cần xử lý bao gồm: “views”, “likes” và “dislikes”.

Đầu tiên, ta mô hình hóa những trường “views”, “likes” và “dislikes” để tìm ra những giới hạn giá trị có tần xuất xảy ra thấp nhất:

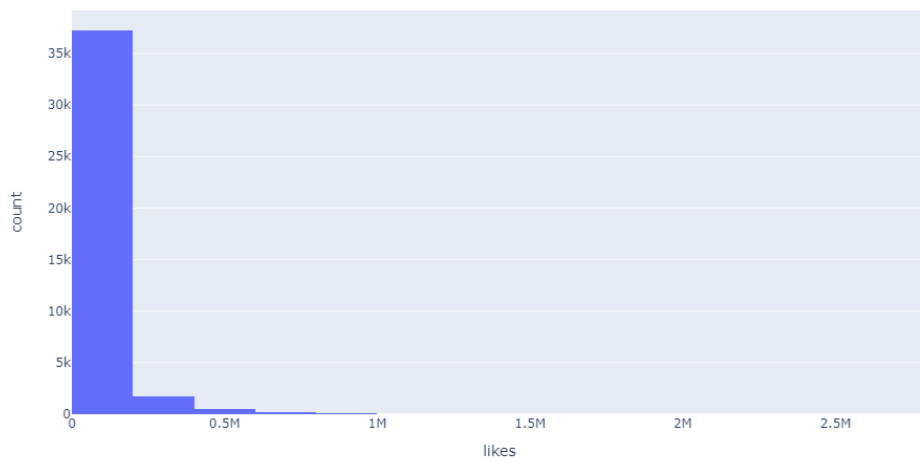
```
fig = px.histogram(df,x="views", nbins=20)
fig = px.histogram(df,x="likes",nbins = 20)
fig = px.histogram(df,x="dislikes", nbins=20)

fig.show()
```

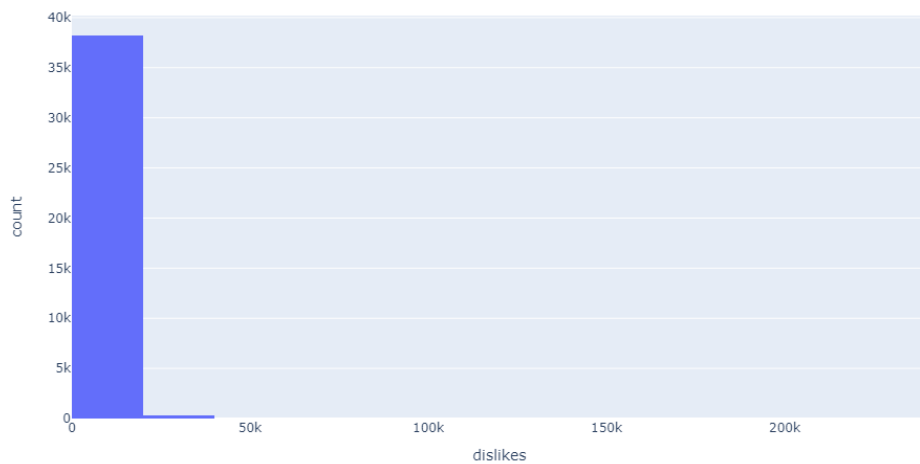
*Bảng 3.3 Vẽ biểu đồ các thuộc tính của dữ liệu*



*Hình 3.8 Biểu đồ số lượng view của bộ dữ liệu*



*Hình 3.9 Biểu đồ số lượng like của bộ dữ liệu*



*Hình 3.10 Biểu đồ số lượng dislike của bộ dữ liệu*

Từ những mô hình trên, ta thấy được giới hạn ngoại lai của từng trường dữ liệu:

- views: 40 000 000
- likes: 400 000
- dislikes: 39 000

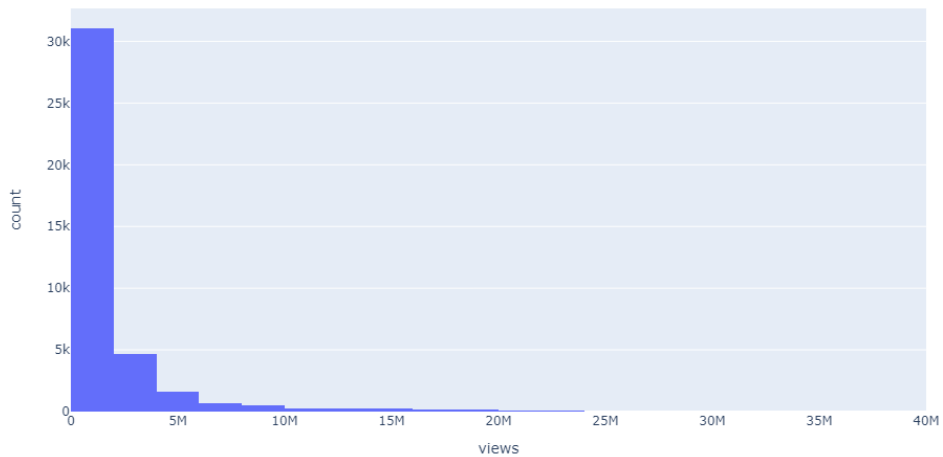
Tiếp theo, ta loại bỏ những dữ liệu nằm ngoài giới hạn ngoại lai:

```
df = df.drop(df[df.views >=40000000].index)
fig = px.histogram(df,x="views", nbins=20)

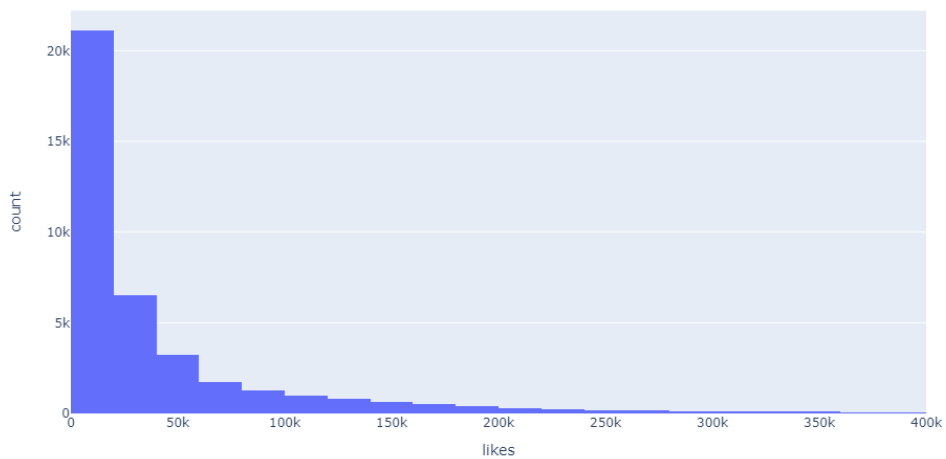
df = df.drop(df[df.likes >=400000].index)
fig = px.histogram(df,x="likes", nbins = 20)

df = df.drop(df[df.dislikes >=39000].index)
fig = px.histogram(df,x="dislikes", nbins = 20)
```

*Bảng 3.4 Vẽ biểu đồ sau khi xử lý giá trị ngoại lai*

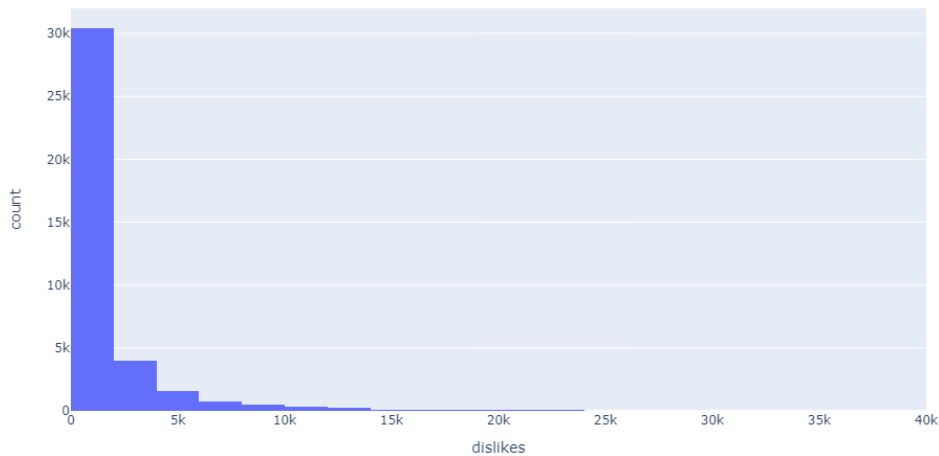


*Hình 3.11 Biểu đồ số lượng view sau khi xử lý ngoại lai*



*Hình 3.12 Biểu đồ số lượng like sau khi xử lý ngoại lai*





Hình 3.13 Biểu đồ số lượng dislike sau khi xử lý ngoại lai

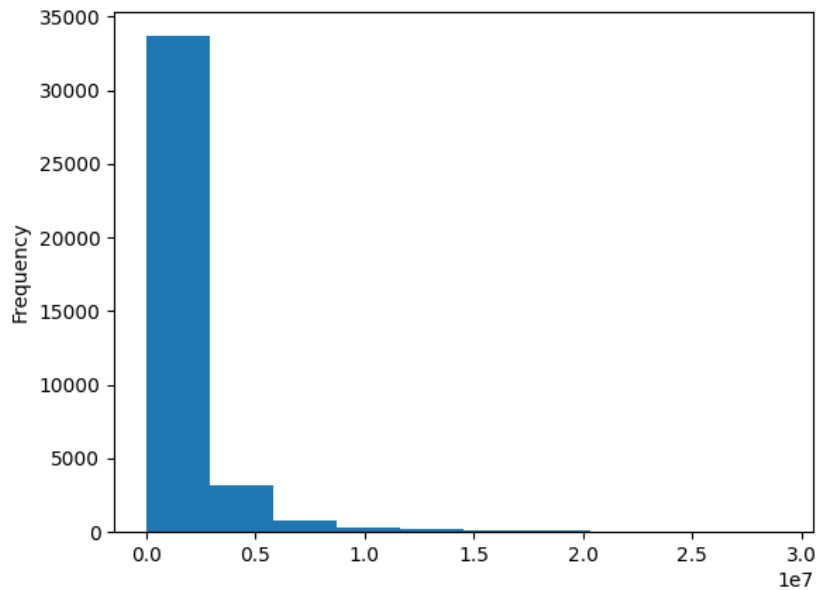
### 3.6. Xử lý dữ liệu nhiều

Phương pháp binning data là một kỹ thuật xử lý dữ liệu nhiều bằng cách chia dữ liệu thành các phân khúc rời rạc và đặt chúng vào các nhóm (bins) tương ứng. Việc này giúp giảm thiểu ảnh hưởng của dữ liệu nhiều lên quá trình phân tích và giúp tăng độ chính xác trong việc đưa ra quyết định.

Có thể áp dụng phương pháp binning data cho các biến số liên tục hoặc phân loại. Để thực hiện phương pháp này, trước tiên cần xác định số lượng các nhóm cần tạo ra. Số lượng nhóm này có thể được chọn dựa trên kinh nghiệm hoặc các phương pháp thống kê như phân phối tần suất của dữ liệu. Sau đó, dữ liệu sẽ được phân bố vào các nhóm tương ứng.

Việc chọn số lượng và kích thước của các nhóm cũng là một yếu tố quan trọng trong phương pháp binning data. Nếu chọn quá nhiều nhóm, dữ liệu sẽ bị chia nhỏ và không còn có tính đại diện, trong khi chọn quá ít nhóm thì sẽ mất mát thông tin và không thể phân tích chi tiết hơn.

Đối với tập dữ liệu này, ta cần phân tích ba biến số: “views”, “likes” và “dislike”. Trước tiên, ta phân tích cột “views”.



Hình 3.14 Biểu đồ tần số cột views

```

maxrange = int(np.ceil(max(df['views'])))
minrange = int(np.floor(min(df['views'])))
ageRange = maxrange - minrange
bins = 10
binwidth = int(np.round(ageRange/bins))

intervals = [ views for views in range(minrange, maxrange + binwidth,
binwidth)]
binlabels = ["bin" + str(i) for i in range (1, int(len(intervals)))]

df['Views_Cut'] = pd.cut(df['views'], bins = intervals, labels = None,
include_lowest=True)

df.groupby('Views_Cut')['views'].count().plot.bar()
plt.xticks(rotation=52)
plt.ylabel('observation count')

```

Bảng 3.5 Xử lý dữ liệu nhiễu của biến views

Ở trên, ta thực hiện tính toán giá trị cận trên của phạm vi "views" bằng cách lấy giá trị tối đa của cột "views" và làm tròn lên đến số nguyên gần nhất. Tương tự, ta tính toán giá trị cận dưới của phạm vi "views" bằng cách lấy giá trị tối thiểu của cột "views" và làm tròn xuống đến số nguyên gần nhất.

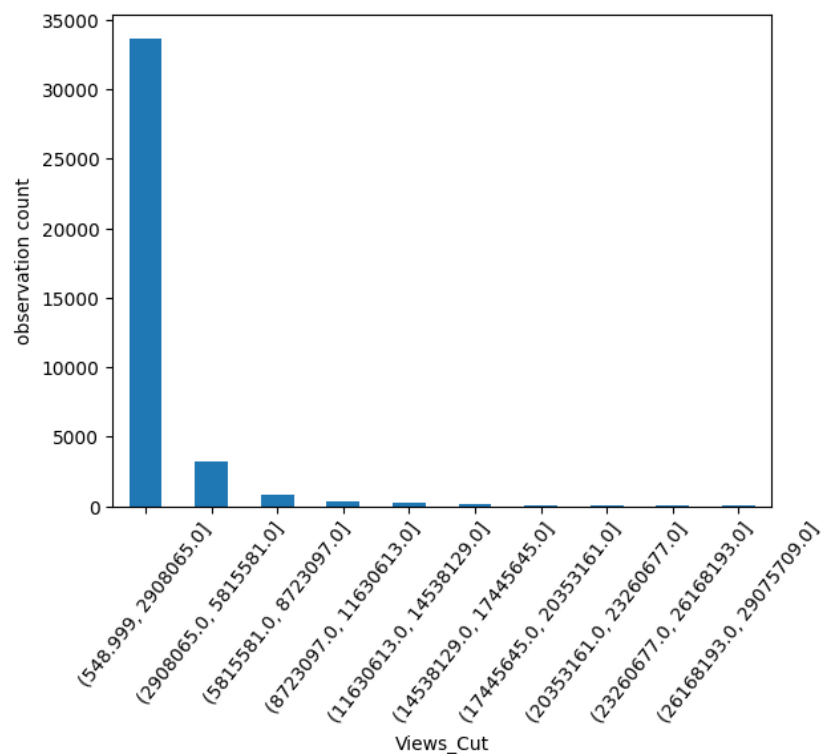
Sau đó, mã tính toán phạm vi "views" bằng cách tính hiệu giữa cận trên và cận dưới của phạm vi. Biến bins được đặt bằng 10 để chỉ định số lượng khoảng chia phạm

vi "views". Mã tính toán độ rộng của mỗi khoảng bằng cách chia phạm vi "views" cho số lượng khoảng và làm tròn đến số nguyên gần nhất.

Sau đó, mã tạo ra một danh sách intervals các khoảng bằng cách sử dụng hàm `range()` từ giá trị cận dưới đến giá trị cận trên của phạm vi "views", với độ rộng của mỗi khoảng là `binwidth`.

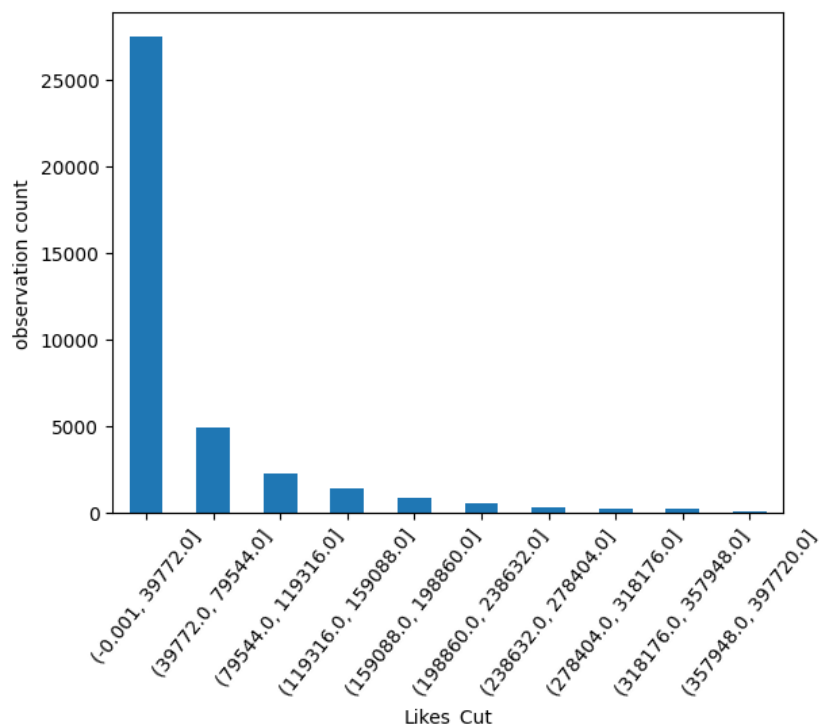
Ta tạo ra một cột mới trong DataFrame df có tên là "Views\_Cut" bằng cách sử dụng hàm `pd.cut()` để chia cột "views" thành các khoảng bằng với các khoảng được xác định trước đó bằng intervals, với các nhãn được xác định bằng binlabels.

Cuối cùng, mã sử dụng hàm `groupby()` để nhóm dữ liệu theo cột "Views\_Cut" và đếm số lượng quan sát trong mỗi khoảng.

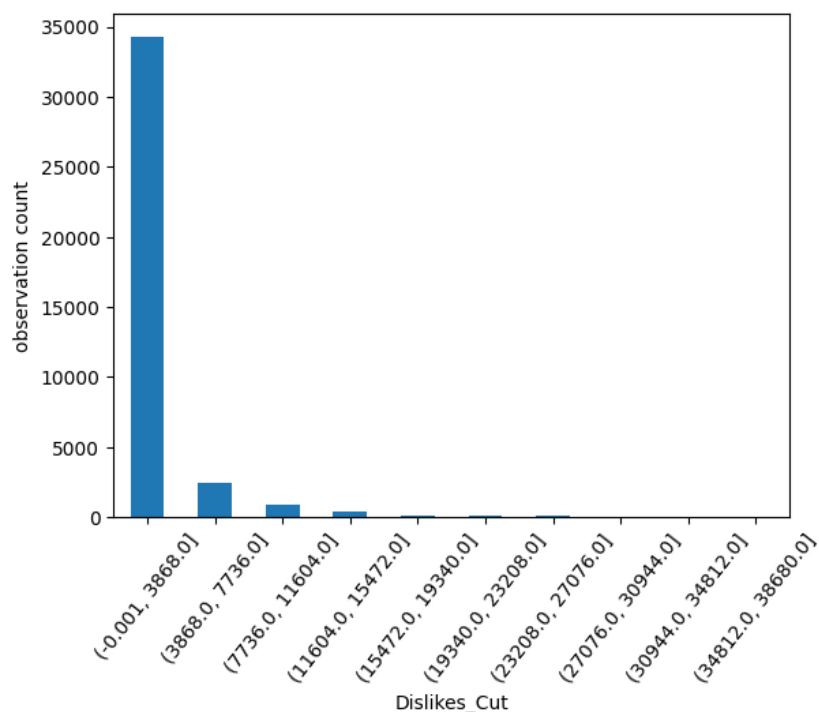


Hình 3.15 Biểu đồ tần số cột views sau khi xử lý dữ liệu nhiễu

Tương tự đối với hai cột “like” và “dislike”, ta tiếp tục phân chia dữ liệu và nhóm chúng theo các bins. Ta được:



Hình 3.16 Biểu đồ tần số cột likes sau khi xử lý dữ liệu nhiễu



Hình 3.17 Biểu đồ tần số cột dislike sau khi xử lý dữ liệu nhiễu

Sau khi đã hoàn thành việc phân nhóm dữ liệu, ta xóa các cột “Views\_Cut”, “Dislikes\_Cut” và “Likes\_Cut” để làm cho DataFrame nhỏ hơn và dễ quản lý hơn.

```
df = df.drop('Dislikes_Cut', axis='columns')
df = df.drop('Views_Cut', axis='columns')
df = df.drop('Likes_Cut', axis='columns')
```

*Bảng 3.6 Xóa cột không cần thiết*

Trên là quá trình tiền xử lý dữ liệu, quá trình này là quá trình quan trọng trong khai phá dữ liệu, giúp chuẩn bị dữ liệu trước khi được áp dụng các phương pháp phân tích. Quá trình này nhằm loại bỏ các giá trị nhiễu, các giá trị bị thiếu, chuẩn hóa các giá trị dữ liệu, rút trích các đặc trưng quan trọng của dữ liệu và chuyển đổi chúng thành các định dạng phù hợp cho việc phân tích.

Việc thực hiện các bước tiền xử lý dữ liệu đảm bảo rằng dữ liệu đã được sạch và phù hợp để sử dụng cho các mô hình phân tích và khai thác dữ liệu.

## CHƯƠNG 4: PHÂN TÍCH DỮ LIỆU KHÁM PHÁ

### 4.1. Phân tích đơn biến

Các đặc tính mô tả của dữ liệu là các thống kê mô tả mà được sử dụng để tóm tắt và mô tả các tính chất của dữ liệu. Các đặc tính mô tả này giúp ta hiểu được cách dữ liệu được phân bố và hình thành trong tập dữ liệu, từ đó có thể rút ra những kết luận và quyết định phù hợp.

Các đặc tính mô tả thường được sử dụng bao gồm:

- Giá trị trung bình.
- Giá trị trung vị.
- Độ lệch chuẩn.
- Phân vị (biểu diễn dưới dạng biểu đồ line).

#### 4.1.1. Thống kê giá trị trung bình

```
[43] df.mean(axis=0, numeric_only=True)

views      1.392780e+06
likes      4.096883e+04
dislikes    1.703419e+03
dtype: float64
```

Hình 4.1 Giá trị trung bình của các cột thuộc tính

Đây là giá trị trung bình của ba thuộc tính trong tập dữ liệu (views, likes, dislikes). Trung bình được sử dụng để biểu thị mức độ trung thành của dữ liệu đối với một số giá trị cụ thể:

- Views:  $1.392780 \times 10^6$
- Likes:  $4.096883 \times 10^4$
- Dislikes:  $1.703419 \times 10^3$

#### 4.1.2. Thống kê giá trị trung vị

```
[44] df.median(axis=0, numeric_only=True)

views      629429.0
likes      16667.0
dislikes     585.0
dtype: float64
```

Hình 4.2 Giá trị trung vị của các cột thuộc tính

Đây là giá trị ở giữa tập dữ liệu, nghĩa là có nửa số giá trị trong tập dữ liệu lớn hơn trung vị và nửa còn lại nhỏ hơn:

- Views: 629429.0

- Likes: 16667.0
- Dislikes: 585.0

#### 4.1.3. Độ lệch chuẩn của các dữ liệu

Là căn bậc hai của phương sai, đây là một độ đo phổ biến để đo sự phân tán của dữ liệu. Sử dụng phương thức `std()` trong thư viện Numpy của Python để tính độ lệch chuẩn.

```
[45] df.std(axis=0, numeric_only=True)

views      2.271627e+06
likes      6.145914e+04
dislikes    3.384368e+03
dtype: float64
```

Hình 4.3 Độ lệch chuẩn của các cột thuộc tính

- Views:  $2.271627 \times 10^6$
- Likes:  $6.145914 \times 10^4$
- Dislikes:  $3.384368 \times 10^3$

#### 4.1.4. Phân vị

Phân vị được sử dụng để phân loại các giá trị trong tập dữ liệu. Nó đại diện cho một phần trăm của các giá trị trong tập dữ liệu, với ý nghĩa là giá trị đó bé hơn hoặc bằng phần trăm đó các giá trị khác trong tập dữ liệu.

index	views	likes	dislikes
count	38467	38467	38467
mean	1392780.363	40968.83407	1703.418853
std	2271627.117	61459.13838	3384.368136
min	549	0	0
25%	232608	5223.5	192
50%	629429	16667	585
75%	1565474	46971	1635
max	29075706	397715	38675

Bảng 4.1 Bảng phân vị của các cột thuộc tính

Thống kê mô tả cho thấy:

- Trung bình lượt views là 1392780.363, số lượt views thấp nhất là 549 và cao nhất là 29075706. Độ lệch chuẩn so với giá trị trung bình là 2271627.117. Đây là chuỗi có xu hướng phân phối lệch về bên trái vì median (mức phân vị 50%) nhỏ hơn mean.
- Trung bình lượt likes là 40968.83407, số lượt views thấp nhất là 0 và cao nhất là 397715. Độ lệch chuẩn so với giá trị trung bình là 61459.13838. Đây là chuỗi

có xu hướng phân phối lệch về bên trái vì median (mức phân vị 50%) nhỏ hơn mean.

- Trung bình lượt dislikes là 1703.418853, số lượt views thấp nhất là 0 và cao nhất là 38675. Độ lệch chuẩn so với giá trị trung bình là 3384.368136. Đây là chuỗi có xu hướng phân phối lệch về bên trái vì median(mức phân vị 50%) nhỏ hơn mean.

## 4.2. Phân tích đa biến

### 4.2.1. Ma trận tương quan

Ma trận tương quan để hiểu tương quan giữa các biến số trong phân tích đa biến. Nó được sử dụng để mô tả mức độ liên quan giữa hai hay nhiều biến số trong tập dữ liệu. Khi hai biến số có tương quan cao, điều đó có nghĩa là khi giá trị của một biến số thay đổi, giá trị của biến số kia sẽ có xu hướng thay đổi theo cùng một hướng.

Ma trận tương quan của tập dữ liệu là một bảng chứa các hệ số tương quan giữa tất cả các cặp biến số trong tập dữ liệu. Các hệ số tương quan này có giá trị từ 0 đến 1, trong đó giá trị 1 biểu thị mối tương quan dương hoàn hảo giữa hai biến và giá trị 0 biểu thị không có tương quan nào giữa hai biến.

Sau khi chuẩn bị dữ liệu, ta tính toán ma trận tương quan bằng cách sử dụng công thức tính toán tương quan giữa hai biến:

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

Trong đó  $\text{cov}(X, Y)$  là hiệp phương sai giữa hai biến X và Y, và  $\sigma_X$  và  $\sigma_Y$  lần lượt là độ lệch chuẩn của X và Y. Ma trận tương quan sẽ có kích thước bằng số lượng biến, trong đó phần tử tại hàng i và cột j sẽ là tương quan giữa biến i và biến j.

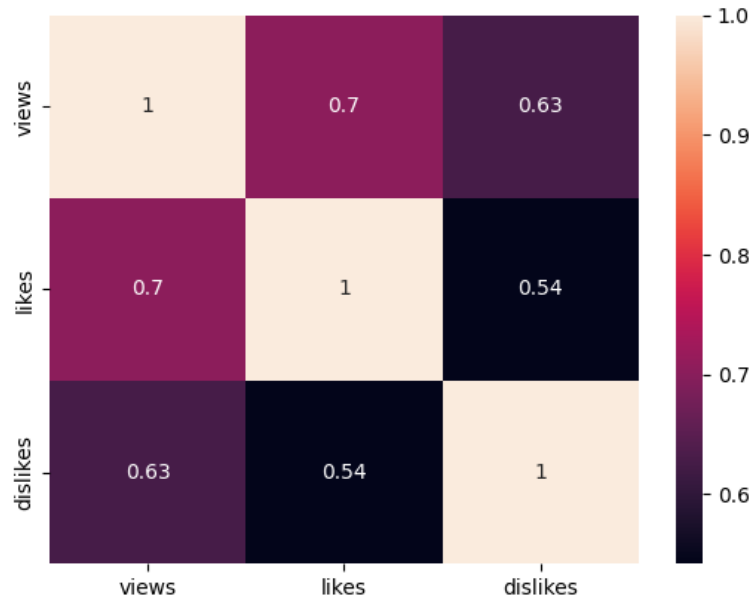
```
correlationMatrix = df[["views", "likes", "dislikes"]].corr()
print(correlationMatrix)
```

	views	likes	dislikes
views	1	0.70437	0.62741
likes	0.70437	1	0.54201
dislikes	0.62741	0.54201	1

Bảng 4.2 Bảng ma trận tương quan

Để trực quan hóa ma trận tương quan, ta có thể sử dụng biểu đồ heatmap (đồ thị nhiệt) để hiển thị mức độ tương quan giữa các biến dưới dạng màu sắc. Các giá trị tương quan lớn sẽ được đại diện bởi màu đậm, còn các giá trị tương quan nhỏ sẽ được đại diện bởi màu nhạt.





Hình 4.4 Đồ thị nhiệt giữa ba thuộc tính chính của tập dữ liệu

Dựa trên ma trận tương quan và biểu đồ heatmap, ta có thể phân tích mối tương quan giữa các biến để tìm hiểu mức độ ảnh hưởng của chúng lên nhau. Các biến có tương quan lớn hơn 0.7 thường được coi là có mối quan hệ mạnh, các biến có tương quan bằng 1 được coi là có mối quan hệ rất mạnh còn các biến có tương quan nhỏ hơn 0.3 thường được coi là không có mối quan hệ đáng kể với nhau:

- Views: số lượng video có lượt views lên top trending có độ tương quan rất cao
- Likes: số lượng video có lượt likes lên top trending có độ tương quan cao
- Dislikes: số lượng video có lượt dislikes lên top trending có độ tương quan trung bình

Hệ số tương quan  $> 0.5$  cho thấy được mối tương quan giữa lượt views, lượt likes và dislikes là mạnh. Số lượng video có lượt views cao sẽ kéo theo lượt likes và dislikes tăng lên.

Phân tích ma trận tương quan giúp chúng ta hiểu được mức độ tương quan giữa các biến số trong tập dữ liệu và giúp chúng ta loại bỏ các biến có tương quan cao hoặc có thể kết hợp chúng lại để giảm số lượng biến và cải thiện hiệu suất của mô hình.

#### 4.2.2. Phân tích chuỗi thời gian

Chúng ta sẽ phân tích khoảng thời gian những ngày hình thành xu hướng của các video trong bộ dữ liệu.

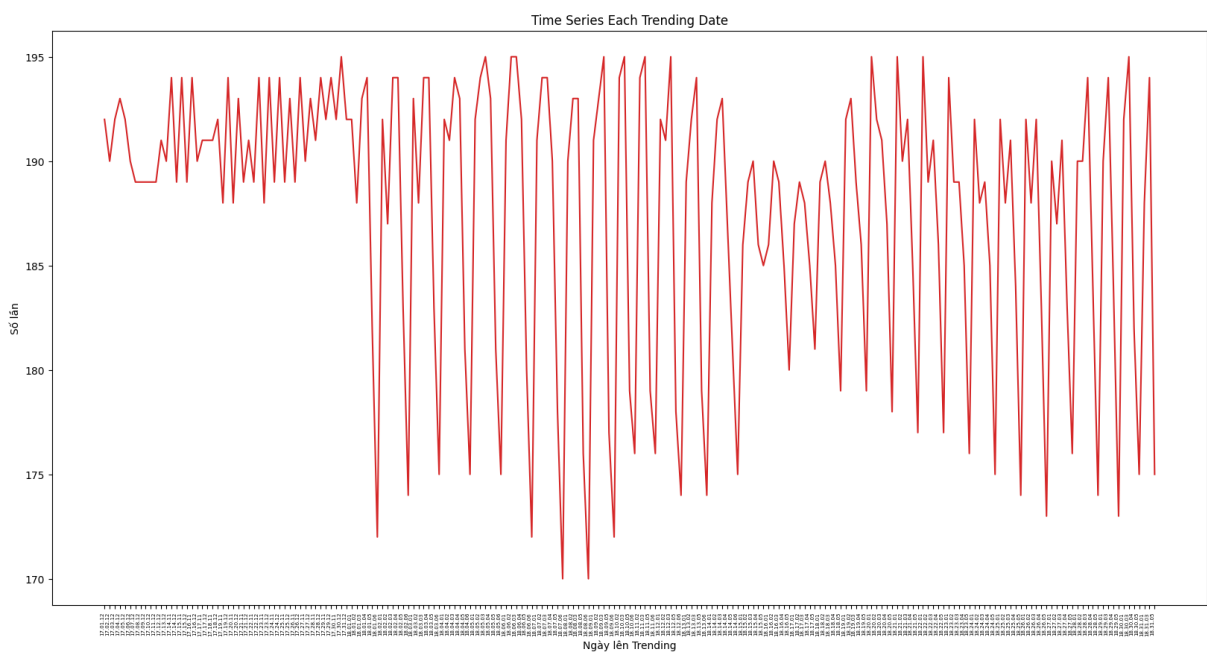
```
time_series = df.groupby(["trending_date"])["trending_date"].count()
time_series = pd.DataFrame({"Trending Date": time_series.index, "Count":
time_series.values})
time_series.head(10)
```

index	Trending Date	Count
0	17.01.12	192
1	17.02.12	190
2	17.03.12	192
3	17.04.12	193
4	17.05.12	192
5	17.06.12	190
6	17.07.12	189
7	17.08.12	189
8	17.09.12	189
9	17.10.12	189

*Bảng 4.3 Bảng phân bố các ngày hình thành xu hướng của video*

Cách phân bố ngày lên trending của các video lên xu hướng là bình thường. Ta thấy số lần video lên top trending trung bình 189 lần.

Sau đó, ta trực quan hóa dữ liệu bằng đồ thị đường để biểu diễn dữ liệu theo chuỗi thời gian:



*Hình 4.5 Biểu đồ phân bố video vào các ngày xu hướng*

## CHƯƠNG 5: KHAI PHÁ DỮ LIỆU

### 5.1. Đặt vấn đề

Dựa vào tập dữ liệu hiện có, với các thuộc tính: ngày thịnh hành, tiêu đề, số lượt xem, số lượt thích, số lượt không thích, ngày phát hành và mô tả. Chúng ta có thể đặt ra những câu hỏi để tiến hành khai phá tập dữ liệu.

Đầu tiên, để hiểu rõ hơn về đa dạng và đặc điểm của các video thịnh hành, chúng ta sẽ phân nhóm chúng dựa trên các thuộc tính như lượt xem, lượt thích, lượt không thích,... Mục tiêu là xác định các nhóm dữ liệu có đặc điểm chung về mức độ quan tâm từ người xem và phản hồi từ cộng đồng YouTube. Vậy chúng ta đặt ra vấn đề, ***“Có bao nhiêu nhóm các video dựa trên các thuộc tính như lượt xem, lượt thích, lượt không thích?”***.

Tiếp theo, số lượt xem, lượt thích và lượt không thích là những chỉ số quan trọng trong việc đánh giá sự thành công của một video trên nền tảng YouTube. Câu hỏi được đưa ở đây là: ***“Có thể dự đoán được số lượt xem của một video dựa trên số lượt thích và không thích không?”***. Câu hỏi này cũng rất quan trọng trong việc giúp các chủ sở hữu video có thể hiểu rõ hơn về yếu tố quan trọng để tạo ra một video thành công và thu hút được nhiều lượt xem trên YouTube. Chính vì vậy, việc tìm ra câu trả lời cho câu hỏi này sẽ mang lại nhiều giá trị và hỗ trợ cho các chủ sở hữu video đưa ra các chiến lược quảng cáo và tiếp cận được nhiều khán giả hơn trên nền tảng YouTube.

Vì vậy, chúng ta sẽ tiến hành phân tích tập dữ liệu và tìm câu trả lời cho hai câu hỏi:

- ***Có bao nhiêu nhóm các video dựa trên các thuộc tính như lượt xem, lượt thích, lượt không thích?***
- ***Có thể dự đoán được số lượt xem của một video dựa trên số lượt thích và không thích không?***

### 5.2. Thuật toán sử dụng

#### 5.2.1. Phân nhóm các video dựa trên các thuộc tính của chúng

Để giải quyết vấn đề này, phân nhóm các video dựa trên các thuộc tính như lượt xem, lượt thích, lượt không thích, một lựa chọn phù hợp là sử dụng thuật toán K-means Clustering.

K-means Clustering là một thuật toán phân cụm phổ biến và đơn giản, phân chia dữ liệu thành các nhóm dựa trên sự tương đồng về khoảng cách giữa các điểm dữ liệu. Trong trường hợp này, chúng ta muốn phân nhóm các video dựa trên các thuộc tính như lượt xem, lượt thích, lượt không thích. K-means Clustering là một lựa chọn hợp lý để xác định các nhóm video tương tự về các thuộc tính này. K-means Clustering là một

thuật toán đơn giản và hiệu quả trong việc phân cụm dữ liệu. Nó có khả năng xử lý dữ liệu lớn và hoạt động tốt với các thuộc tính số.

Ta có thể sử dụng thư viện scikit-learn trong Python để áp dụng thuật toán K-means Clustering. Đầu tiên, tiền xử lý dữ liệu bằng cách chọn các thuộc tính cần thiết và chuẩn hóa dữ liệu nếu cần. Sau đó, xác định số lượng nhóm mong muốn và áp dụng thuật toán K-means Clustering để phân nhóm các video. Bạn có thể sử dụng khoảng cách Euclidean hoặc khoảng cách cosine để đo độ tương đồng giữa các điểm dữ liệu.

Trong quá trình áp dụng, bạn cần chú ý đến số lượng nhóm ( $k$ ) mà bạn muốn tạo ra và các phương pháp khởi tạo khác nhau (như ngẫu nhiên hoặc chọn ngẫu nhiên các điểm dữ liệu làm trung tâm ban đầu) có thể ảnh hưởng đến kết quả của thuật toán. Để đánh giá hiệu suất của phân cụm, bạn có thể sử dụng các độ đo như Silhouette Score hoặc Calinski-Harabasz Index.

### 5.2.2. Dự đoán được số lượt xem của một video

Để giải quyết vấn đề này, dự đoán được số lượt xem của một video dựa trên số lượt thích và không thích không, một lựa chọn phù hợp là sử dụng thuật toán Support Vector Machines (SVM).

SVM có khả năng xử lý cả dữ liệu có tính chất tuyến tính và phi tuyến tính. Trong trường hợp các thuộc tính lượt thích và không thích không có mối quan hệ tuyến tính trực tiếp với lượt xem, SVM có thể tìm ra các siêu phẳng phi tuyến tính để tạo ra mô hình dự đoán chính xác. SVM có khả năng xử lý các tập dữ liệu có số chiều cao, tức là nhiều thuộc tính đầu vào. Trong trường hợp chúng ta có nhiều thuộc tính khác nhau để dự đoán số lượt xem, SVM có thể xử lý một số lượng lớn các biến đầu vào.

SVM là một thuật toán học máy mạnh mẽ với khả năng tìm ra các giải pháp tối ưu cho bài toán phân loại và hồi quy. SVM có khả năng xử lý các tập dữ liệu lớn và đa dạng. Thuật toán sử dụng một tập hợp con của các điểm dữ liệu để xây dựng các đường ranh giới phân chia, giúp mô hình linh hoạt và chính xác.

Để áp dụng thuật toán, chúng ta cần chuẩn bị dữ liệu bằng việc tách các thuộc tính lượt thích và không thích làm đầu vào ( $X$ ) và số lượt xem làm đầu ra ( $y$ ). Tiếp theo, chúng ta áp dụng thuật toán SVM trên dữ liệu huấn luyện để tạo ra mô hình. Sau khi mô hình được huấn luyện, chúng ta sử dụng nó để dự đoán số lượt xem của các video mới dựa trên số lượt thích và không thích. Cuối cùng, chúng ta đánh giá hiệu suất của mô hình dựa trên các độ đo như độ chính xác, độ sai số và độ phân loại đúng.

## CHƯƠNG 6: ĐÁNH GIÁ VÀ CHỌN THUẬT TOÁN

### 6.1. Áp dụng thuật toán

#### 6.1.1. Thuật toán K-means Clustering

Trước tiên, ta tạo tập dữ liệu chứa các trường cần thiết: views, likes, dislikes

```
features = ["views", "likes", "dislikes"]  
data = df[features].copy()
```

Bảng 6.1 Tạo tập data chứa dữ liệu cần thiết

```
[ ] data
```

	views	likes	dislikes
video_id			
2kyS6SvSYSE	748374	57527	2966
1ZAPwfrtAFY	2418783	97185	6146
5qpjK5DgCt4	3191434	146033	5339
puqaWrEC7tY	343168	10172	666
d380meD0W0M	2095731	132235	1989
...	...	...	...
_QWZvU7VCn8	5564576	46351	2295
7UoP9ABJXGE	5534278	45128	1591
BZt0qjTWNhw	1685609	38160	1385
D6Oy4LfoqsU	1066451	48068	1032
oV0zkMe1K8s	5660813	192957	2846

38467 rows x 3 columns

Hình 6.1 Dữ liệu được trả về

Nhằm dễ dàng xử lý dữ liệu, ta có thể đưa các giá trị dữ liệu về một khoảng giá trị cụ thể để dễ dàng so sánh và xử lý. Ta thực hiện bằng cách chuẩn hóa dữ liệu theo phạm vi từ 1 đến 10.

```
data = ((data - data.min()) / (data.max() - data.min())) * 9 + 1
```

Bảng 6.2 Chuẩn hóa giá trị trong miền từ 1 đến 10

```
[209] data.describe()
```

	views	likes	dislikes
count	38467.000000	38467.000000	38467.000000
mean	1.430955	1.927095	1.396400
std	0.703165	1.390775	0.787571
min	1.000000	1.000000	1.000000
25%	1.071832	1.118204	1.044680
50%	1.194665	1.377162	1.136134
75%	1.484411	2.062919	1.380478
max	10.000000	10.000000	10.000000

Hình 6.2 Giá trị trả về

Tiếp theo là các hàm cần thiết sử dụng trong thuật toán K-means Clustering:

- Tạo tâm ngẫu nhiên cho tập dữ liệu:

```
def random_centroids(data, k):  
    centroids = []  
    for i in range(k):  
        centroid = data.apply(lambda x: float(x.sample()))  
        centroids.append(centroid)  
    return pd.concat(centroids, axis=1)
```

Bảng 6.3 Hàm random\_centroids

- Tạo labels cho tập dữ liệu:

```
def get_labels(data, centroids):  
    distances = centroids.apply(lambda x: np.sqrt(((data - x) **  
2).sum(axis=1))))  
    return distances.idxmin(axis=1)
```

Bảng 6.4 Hàm get\_labels

- Tạo tâm mới:

```
def new_centroids(data, labels, k):  
    centroids = data.groupby(labels).apply(lambda x:  
np.exp(np.log(x).mean())) .T  
    return centroids
```

Bảng 6.5 Hàm new\_centroids

- Vẽ đồ thị biểu diễn cho các cụm:

```
def plot_clusters(data, labels, centroids, iteration):
    pca = PCA(n_components=2)
    data_2d = pca.fit_transform(data)
    centroids_2d = pca.transform(centroids.T)
    clear_output(wait=True)
    plt.title(f'Iteration {iteration}')
    plt.scatter(x=data_2d[:,0], y=data_2d[:,1], c=labels, s=1)
    plt.scatter(x=centroids_2d[:,0], y=centroids_2d[:,1], s=10,
marker='x', color = 'r')
    plt.show()
```

*Bảng 6.6 Hàm plot\_clusters*

Phần chính của thuật toán K-means Clustering:

```
max_iterations = 100
centroid_count = 5

centroids = random_centroids(data, centroid_count)
old_centroids = pd.DataFrame()
iteration = 1

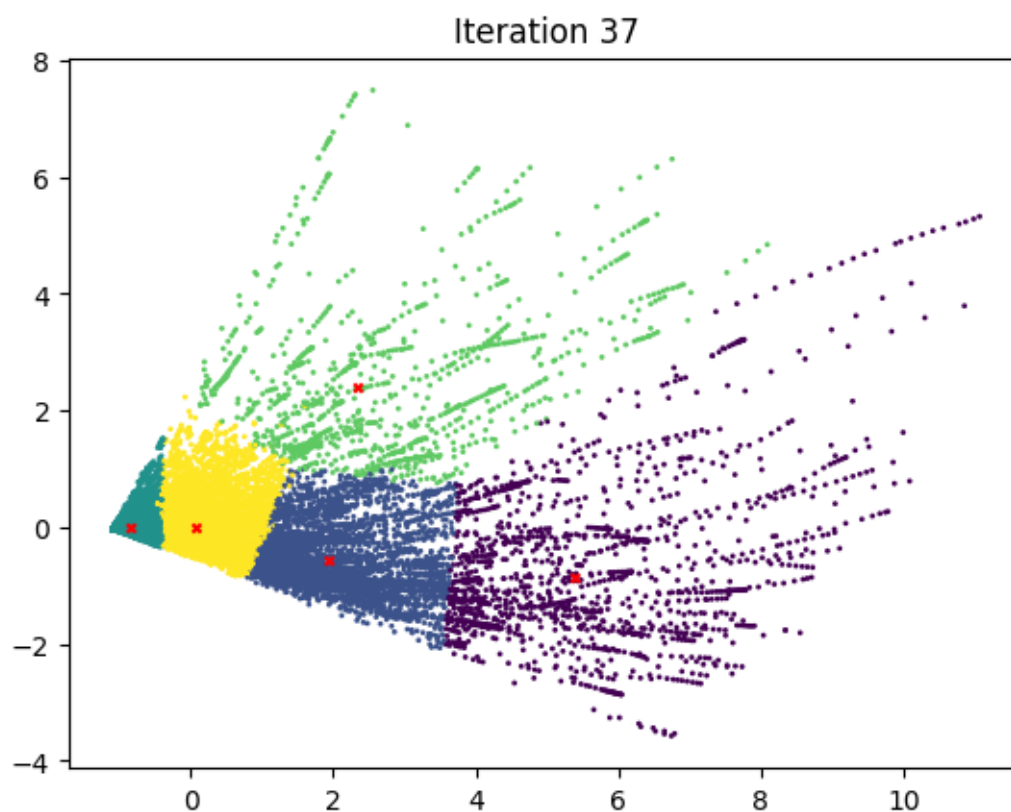
while iteration < max_iterations and not
centroids.equals(old_centroids):
    old_centroids = centroids

    labels = get_labels(data, centroids)
    centroids = new_centroids(data, labels, centroid_count)
    plot_clusters(data, labels, centroids, iteration)
    iteration += 1
```

*Bảng 6.7 Kết hợp các hàm để phân nhóm tập dữ liệu*

Để tiến hành phân nhóm, ta thực hiện các bước sau:

- Khởi tạo các centroid ban đầu: Trước khi bắt đầu thuật toán, cần khởi tạo ngẫu nhiên các centroid ban đầu.
- Gán nhãn cho các mẫu dữ liệu: Với các centroid đã được khởi tạo, cần gán nhãn cho từng mẫu dữ liệu dựa trên khoảng cách từ mẫu đó đến các centroid.
- Cập nhật centroid mới: Dựa trên nhãn đã được gán, cần tính toán lại vị trí của các centroid mới dựa trên trung bình của các mẫu dữ liệu trong cùng một nhóm.
- Lặp lại quá trình cho đến khi tiêu chí dừng được đáp ứng: Quá trình gán nhãn và cập nhật centroid mới sẽ được lặp lại cho đến khi các centroid không thay đổi hoặc đạt tới số lần lặp tối đa.
- Vẽ đồ thị các nhóm dữ liệu bằng hàm plot\_clusters.



Hình 6.3 Đồ thị phân nhóm bằng K-means Clustering

index	0	1	2	3	4
views	3.122072	1.872732	1.118228	2.57473	1.502321
likes	6.972804	3.877585	1.192435	2.940626	1.969856
dislikes	2.555887	1.642418	1.096252	4.399912	1.400011

Bảng 6.8 Dữ liệu các centroid được tạo

Sau khi hoàn tất việc phân nhóm, các video sẽ được gắn nhãn đúng với nhóm mà chúng thuộc về. Qua đó, chúng ta nhận biết sự tương đồng hoặc khác biệt giữa các video trong tập dữ liệu. Các nhóm video có thể được hiểu là các tập hợp video có đặc điểm tương tự nhau.



```
[222] df[labels == 3][["title"] + features]
```

	title	views	likes	dislikes
video_id				
1ZAPwfrtAFY	The Trump Presidency: Last Week Tonight with J...	2418783	97185	6146
5qpjK5DgCt4	Racist Superman   Rudy Mancuso, King Bach & Le...	3191434	146033	5339
d380meD0W0M	I Dare You: GOING BALDI?	2095731	132235	1989
5E4ZBSInqUU	Marshmello - Blocks (Official Music Video)	687582	114188	1333
ujyTQNNjjDU	G-Eazy - The Plan (Official Video)	2642930	115795	3055
...	...	...	...	...
99t4EBwIAt8	Shawn Mendes Answers the Web's Most Searched Q...	2310794	105783	558
SQsPvrev_bQ	435	2252933	129865	1550
oLDbO545aKQ	Terrible Magicians   Rudy Mancuso & Juanpa Zurita	3825440	196635	4514
Xr2rgT9uEnA	LIE DETECTOR TEST WITH MY GIRLFRIEND!	3229540	109945	3062
oV0zkMe1K8s	How Black Panther Should Have Ended	5660813	192957	2846

3338 rows x 4 columns

Hình 6.4 Ví dụ minh họa việc phân nhóm thành công

### 6.1.2. Thuật toán Support Vector Machines

Đầu tiên, ta khai báo các thư viện được sử dụng trong thuật toán:

```
from sklearn.svm import SVR
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
```

Bảng 6.9 Khai báo các thư viện

Tiếp theo, ta tách các thuộc tính được sử dụng trong thuật toán để dễ dàng trong việc xuất trích và khai báo:

```
X = df[['likes', 'dislikes']].copy()
y = df['views'].copy()
```

Bảng 6.10 Tạo biến X và y

Ta tách tập dữ liệu (X) và (y) thành hai phần:

- Train: sử dụng để huấn luyện cho mô hình
- Test: sử dụng để kiểm tra độ chính xác của mô hình

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)
```

Bảng 6.11 Tách tập training và testing

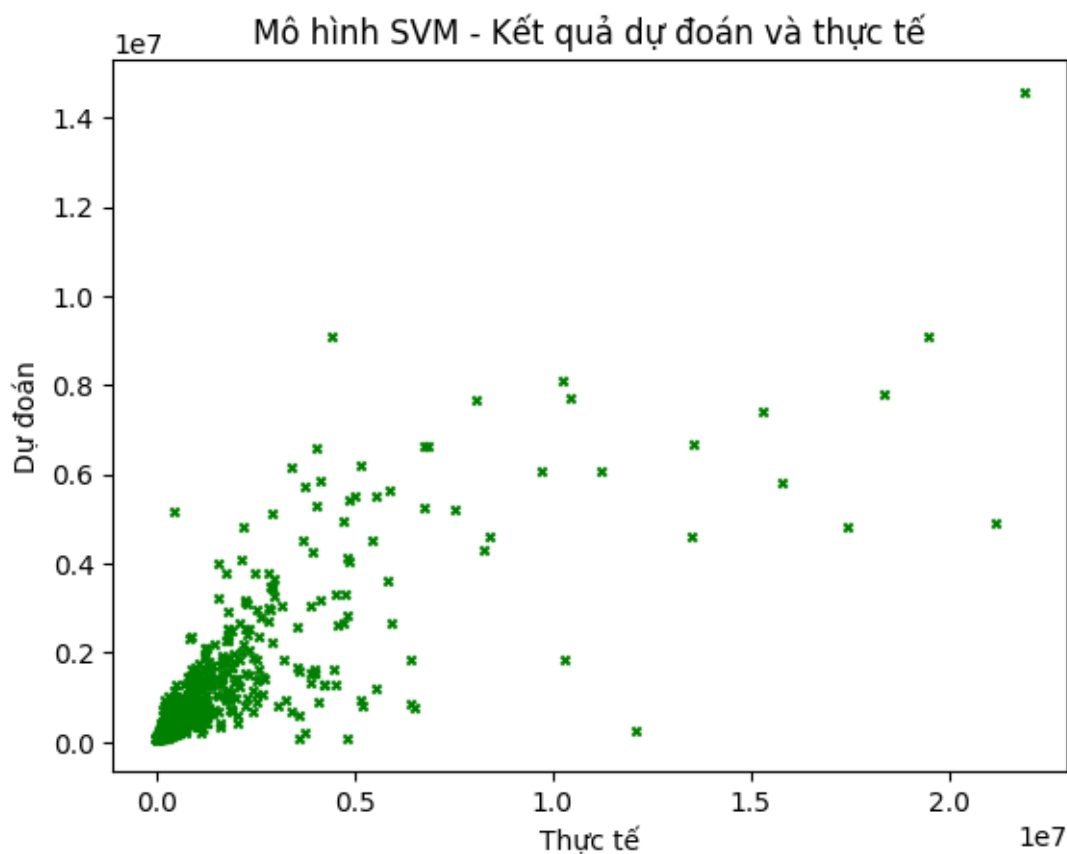
Phần chính của thuật toán bao gồm việc tạo mô hình và tạo tập dữ liệu dự đoán:

```
# Tạo mô hình SVM
svm = SVR(kernel='linear')
# Huấn luyện mô hình
svm.fit(X_train, y_train)
# Dự đoán số lượt xem cho tập kiểm tra
y_pred = svm.predict(X_test)
```

Bảng 6.12 Mô hình Support Vector Machines

Ta thực hiện các bước sau:

- Khởi tạo mô hình SVM với các tham số tùy chọn (ví dụ: linear SVM, kernel SVM).
- Sử dụng phương thức fit để huấn luyện mô hình trên dữ liệu huấn luyện.
- Sử dụng phương thức predict để dự đoán nhãn của các điểm dữ liệu kiểm tra dựa trên mô hình đã được huấn luyện.



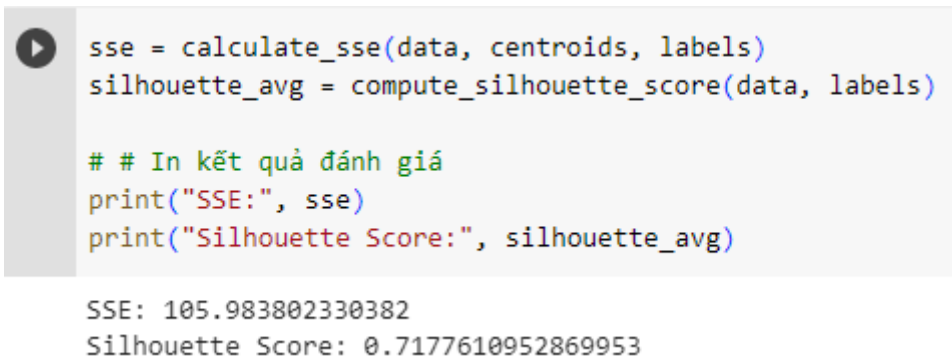
Hình 6.5 Mô hình kết quả của thuật toán Support Vector Machines

Sau khi mô hình đã được huấn luyện và điều chỉnh, bạn có thể sử dụng nó để dự đoán nhãn cho dữ liệu mới, bằng cách sử dụng phương thức predict.

## 6.2. Tiến hành đánh giá kết quả

### 6.2.1. Thuật toán K-means Clustering

- *SSE* (Sum of Squared Errors): Đây là độ đo thường được sử dụng để đánh giá hiệu suất của thuật toán K-means. Nó tính tổng bình phương khoảng cách từ mỗi điểm dữ liệu tới centroid tương ứng của nó. Độ giảm SSE càng lớn, thuật toán càng tốt.
- *Silhouette Score*: Đây là một độ đo đánh giá chất lượng phân nhóm trong thuật toán K-means. Nó tính toán độ tách biệt giữa các cụm và độ liên kết của các điểm dữ liệu trong cùng một cụm. Giá trị Silhouette Score nằm trong khoảng  $[-1, 1]$ , với giá trị càng gần 1 cho thấy phân nhóm tốt, giá trị gần 0 cho thấy phân nhóm không rõ ràng, và giá trị gần -1 cho thấy phân nhóm không tốt.



```
sse = calculate_sse(data, centroids, labels)
silhouette_avg = compute_silhouette_score(data, labels)

## In kết quả đánh giá
print("SSE:", sse)
print("Silhouette Score:", silhouette_avg)
```

SSE: 105.983802330382  
Silhouette Score: 0.7177610952869953

Hình 6.6 Giá trị đánh giá của mô hình K-means Clustering

Dựa vào các giá trị trên, ta có thể đánh giá như sau:

- Giá trị SSE nhỏ (105.983802330382) cho thấy tổng bình phương sai số giữa các điểm dữ liệu và các trung tâm cụm là khá thấp, cho thấy mức độ tổ chức của phân cụm khá tốt.
- Giá trị Silhouette Score (0.7177610952869953) lớn hơn 0.5, đây là một giá trị tương đối tốt. Nó cho thấy các điểm dữ liệu trong cùng một cụm gần nhau hơn và xa các cụm khác.

### 6.2.2. Thuật toán Support Vector Machines

- *mean\_squared\_error*: Đây là một metric đo lường sự sai khác trung bình giữa các giá trị dự đoán và giá trị thực tế. Nó tính toán giá trị trung bình của bình phương sai khác giữa các điểm dữ liệu. Giá trị MSE càng thấp thì mô hình càng tốt, với giá trị 0 nghĩa là mô hình dự đoán hoàn toàn chính xác.
- *r2\_score* (hay còn gọi là R-squared score): Đây là một metric đo lường khả năng giải thích của mô hình hồi quy trên dữ liệu. Nó đo lường tỉ lệ phương sai của biến mục tiêu mà mô hình có thể giải thích được. Giá trị R-squared càng gần 1 thì mô

hình càng tốt, với giá trị 1 nghĩa là mô hình giải thích được toàn bộ phương sai của biến mục tiêu.

Cả hai metric này đều sử dụng các giá trị dự đoán và giá trị thực tế để đánh giá hiệu suất của mô hình hồi quy. MSE tập trung vào đo lường sự chính xác về mức độ sai lệch dự đoán, trong khi R-squared đo lường khả năng giải thích của mô hình trên dữ liệu.

```
[74] from sklearn.metrics import mean_squared_error, r2_score

# Evaluate the model
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print("Mean Squared Error:", mse)
print("R-squared Score:", r2)

Mean Squared Error: 180.107
R-squared Score: 0.7688750018867881
```

*Hình 6.7 Giá trị đánh giá của mô hình Support Vector Machines*

Thuật toán của ta trả về giá trị:

- MSE: Giá trị MSE là 180.107, điều này cho thấy độ lệch bình phương trung bình giữa dự đoán và giá trị thực tế là khá nhỏ và chấp nhận được. Mô hình có độ chính xác tương đối cao trong việc dự đoán giá trị.
- R-squared score: Giá trị R-squared score là 0.7688750018867881, cho thấy mô hình giải thích được khoảng 76.89% phương sai của biến mục tiêu. Điều này cho thấy mô hình có khả năng dự đoán tốt và giải thích một phần lớn sự biến thiên của biến mục tiêu.

## CHƯƠNG 7: KẾT QUẢ VÀ THẢO LUẬN

### 7.1. Đánh giá kết quả

Để đánh giá kết quả từ hai phương pháp áp dụng, ta tiến hành kiểm tra:

- Số lượng nhóm video: Dựa vào thuật toán K-means Clustering, ta phân tập dữ liệu thành 5 nhóm với các tâm khác nhau. Thuật toán đã hoàn thành việc phân nhóm theo các thuộc tính yêu cầu (likes, views, dislikes).
- Dự đoán số lượt xem: Dựa vào thuật toán Support Vector Machines, ta đã dự đoán số lượt xem của một video dựa trên số lượt thích và không thích. Sau khi đánh giá kết quả bằng cách kiểm tra độ chính xác, độ phù hợp và độ tin cậy của mô hình dự đoán, so sánh kết quả dự đoán với dữ liệu thực tế ta đã xác định khả năng của mô hình trong việc dự đoán số lượt xem.
- Độ tin cậy và độ ổn định: Việc các thông số đánh giá trả về kết quả không khả quan, ta thấy được độ tin cậy và tính ổn định của mô hình chỉ mang tính chất tương đối. Việc bổ sung và cải thiện mô hình là điều cần thiết để có thể đưa mô hình vào việc ứng dụng thực tế.

### 7.2. Điểm mạnh của nghiên cứu

Điểm mạnh của nghiên cứu như sau:

- Tính ứng dụng: Đề tài tập trung vào việc áp dụng các thuật toán và phương pháp machine learning để giải quyết vấn đề thực tế trong lĩnh vực truyền thông và xã hội. Việc dự đoán số lượt xem của video và phân nhóm video dựa trên các thuộc tính như lượt xem, lượt thích và lượt không thích có tính ứng dụng cao và có thể hỗ trợ các quyết định quan trọng trong việc tối ưu hóa hiệu quả và tăng cường tương tác của video trên các nền tảng truyền thông xã hội.
- Phân tích đa chiều: Đề tài xem xét các yếu tố quan trọng như lượt xem, lượt thích và lượt không thích để hiểu và phân loại video theo các nhóm tương ứng. Điều này cho phép hiểu rõ hơn về mức độ thu hút và tương tác của video với khán giả. Đồng thời, phương pháp phân nhóm giúp tạo ra cái nhìn tổng quan về sự đa dạng và phân phối của các video trong tập dữ liệu.
- Sử dụng các thuật toán machine learning: Đề tài sử dụng các thuật toán phân loại Support Vector Machines và K-means Clustering để xây dựng các mô hình và thực hiện phân nhóm và dự đoán số lượt xem của video. Các thuật toán này đã được chứng minh hiệu quả trong nhiều bài toán và có khả năng xử lý dữ liệu phức tạp và tạo ra kết quả đáng tin cậy.
- Tiềm năng phát triển: Nghiên cứu có tiềm năng phát triển và mở rộng để bao gồm các thuộc tính khác và áp dụng trên các tập dữ liệu lớn hơn.

Tuy nhiên, để đánh giá toàn diện, cần xem xét cả các điểm yếu của nghiên cứu và giới hạn của phương pháp sử dụng để đảm bảo tính khách quan và đáng tin cậy.

### 7.3. Điểm yếu của nghiên cứu

Một số điểm yếu có thể của đề tài là:

- Hạn chế dữ liệu: Điểm yếu chung của các dự án liên quan đến dữ liệu là có thể gặp hạn chế về tập dữ liệu. Dữ liệu có thể không đủ lớn, không đại diện hoặc không đầy đủ để phân tích một cách toàn diện. Điều này có thể ảnh hưởng đến khả năng đưa ra dự đoán chính xác và phân nhóm một cách tốt nhất.
- Thiếu yếu tố quan trọng: Trong đề tài này, chúng ta chỉ xem xét lượt xem, lượt thích và lượt không thích làm các yếu tố chính để dự đoán và phân nhóm video. Thiếu yếu tố này có thể làm giảm độ chính xác và độ tin cậy của kết quả dự đoán và phân nhóm.
- Sự giới hạn của các thuật toán: Mặc dù các thuật toán được sử dụng trong đề tài đã được chứng minh là hiệu quả trong nhiều bài toán, nhưng chúng có thể có những giới hạn riêng. Ví dụ, thuật toán K-means Clustering yêu cầu xác định số lượng nhóm trước đó, điều này có thể làm giới hạn khả năng phân nhóm một cách chính xác. Các thuật toán machine learning khác cũng có những giới hạn và điều kiện giả định riêng, và việc không tuân thủ các điều kiện này có thể ảnh hưởng đến hiệu quả của mô hình.
- Độ chính xác và độ tin cậy: Mặc dù chúng ta sử dụng các độ đo như Mean Squared Error, R-squared score và Silhouette Score để đánh giá kết quả, nhưng điều này không đảm bảo độ chính xác và độ tin cậy tuyệt đối của kết quả. Các độ đo này chỉ cung cấp một cái nhìn tổng quan và phụ thuộc vào giả định và giới hạn của chúng. Việc đánh giá kết quả cần được thực hiện cẩn thận và kết hợp với phản hồi từ người dùng.

## CHƯƠNG 8: KẾT LUẬN

### 8.1. Tổng kết kết quả

Trong đề tài này, chúng ta đã thực hiện các công việc sau và đạt được những kết quả sau đây:

- Thu thập dữ liệu: Chúng ta đã thu thập dữ liệu về các video trên các nền tảng truyền thông xã hội, bao gồm thông tin về lượt xem, lượt thích và lượt không thích.
- Tiền xử lý dữ liệu: Chúng ta đã thực hiện các bước tiền xử lý dữ liệu như xóa dữ liệu trùng lặp, xử lý dữ liệu thiếu và chuẩn hóa dữ liệu để chuẩn bị cho việc phân tích.
- Phân nhóm video: Sử dụng thuật toán K-means Clustering, chúng ta đã phân nhóm các video dựa trên các thuộc tính lượt xem, lượt thích và lượt không thích. Kết quả phân nhóm giúp chúng ta hiểu rõ hơn về sự đa dạng và phân phối của các video trong tập dữ liệu.
- Dự đoán số lượt xem: Sử dụng thuật toán machine learning Support Vector Machines, chúng ta đã xây dựng mô hình dự đoán số lượt xem của một video dựa trên số lượt thích và không thích. Kết quả dự đoán giúp chúng ta ước lượng và đánh giá tiềm năng của một video trong việc thu hút sự quan tâm và tương tác từ khán giả.
- Đánh giá kết quả: Chúng ta đã sử dụng các độ đo như Mean Squared Error, R-squared score và Silhouette Score để đánh giá hiệu quả của các mô hình và phương pháp được sử dụng. Các độ đo này cho phép chúng ta đánh giá mức độ chính xác và phù hợp của kết quả dự đoán và phân nhóm.

Tổng kết lại, đề tài đã thành công trong việc phân nhóm và dự đoán số lượt xem của video dựa trên các thuộc tính như lượt xem, lượt thích và lượt không thích. Kết quả này có thể hỗ trợ quyết định và tối ưu hóa hiệu quả và tương tác của video trên các nền tảng truyền thông xã hội.

### 8.2. Kết luận hiệu quả

Đề tài đã đạt được hiệu quả đáng kể trong việc phân nhóm và dự đoán số lượt xem của các video dựa trên các thuộc tính lượt xem, lượt thích và lượt không thích. Các phương pháp và mô hình machine learning được áp dụng đã mang lại kết quả đáng tin cậy và hữu ích trong việc hiểu và tối ưu hóa hoạt động của các video trên nền tảng truyền thông xã hội.

Phân nhóm video giúp chúng ta nhìn nhận rõ hơn về sự đa dạng và phân bố của các video trong tập dữ liệu. Điều này có thể giúp chúng ta tối ưu hóa chiến lược phân phối video, tăng cường sự tương tác và quan tâm từ khán giả.

Mô hình dự đoán số lượt xem của video dựa trên lượt thích và không thích cung cấp một công cụ mạnh để ước lượng và đánh giá tiềm năng của một video. Kết quả dự đoán có thể hỗ trợ quyết định về nội dung, quảng cáo và thời gian phát hành video để tối đa hóa sự quan tâm và tương tác từ khán giả.

Tổng cộng, đề tài đã đạt được hiệu quả cao và mang lại những kết quả quan trọng và hữu ích trong việc phân nhóm và dự đoán số lượt xem của video. Các phương pháp và kỹ thuật đã được áp dụng một cách khéo léo và đáng tin cậy, mang lại giá trị và tiềm năng trong việc tối ưu hóa hiệu quả và tương tác của video trên các nền tảng truyền thông xã hội.

## **TÀI LIỆU THAM KHẢO**

- [1] Wikipedia (2023). “K-means Clustering”. Truy cập tại: [https://en.wikipedia.org/wiki/K-means\\_clustering](https://en.wikipedia.org/wiki/K-means_clustering) (truy cập ngày 10/05/2023)
- [2] Wikipedia (2023). “Support Vector Machines”. Truy cập tại: [https://en.wikipedia.org/wiki/Support-vector\\_machine](https://en.wikipedia.org/wiki/Support-vector_machine) (truy cập ngày 10/05/2023)
- [3] Scikit-learn documentation (2023). “Evaluation Metrics for Classification”. Truy cập tại: [https://scikit-learn.org/stable/modules/model\\_evaluation.html#classification-metrics](https://scikit-learn.org/stable/modules/model_evaluation.html#classification-metrics) (truy cập ngày 10/05/2023)
- [4] Scikit-learn documentation (2023). “Silhouette Score”. Truy cập tại: [https://scikit-learn.org/stable/modules/generated/sklearn.metrics.silhouette\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.silhouette_score.html) (truy cập ngày 10/05/2023)

**--HẾT--**