

# CLUSTERING AND FITTING:

## NETFLIX MOVIES AND TV SHOWS ANALYSIS

Name: Fawas Afsal  
Student ID: 22080258

Dataset: <https://www.kaggle.com/code/onyonixch/netflix-movies-tv-shows-eda-and-clustering/input>

### Introduction

- This dataset includes a wide variety of movie and TV show information, including genre, title, director, cast member names country of origin date released by clustering and fitting, the analysis aims to reveal patterns and connections in this multi-media dataset.
- Through aggregating related titles by commonalities, can analyze patterns in the development of content, consumer tastes and demands, as well as production features.
- Furthermore, the variety of explored genres and themes allows for a better understanding of the entertainment industry.

### Aims and Objectives

#### Aim

The project aim is to utilize autoencoders to analyze and cluster Netflix movies and TV shows, revealing hidden designs and giving experiences into content conveyance trends.

#### Objectives

- To import and research the Netflix dataset, and find and resolve data disparities.
- To develop an autoencoder system for the extraction of elements and clustering.
- To execute a customized clustering layer which will result in more successful cluster designations.
- To prepare the model in two stages, we upgrade for the portrayal of highlights and clustering precision.
- To show the clustering discoveries through static and intuitive diagrams, taking into consideration better examination of content examples.

### Background

- The research centers around doing an inside-and-out investigation of Netflix network shows and films utilizing cutting-edge machine-learning methods.
- This involves finding and remedying data anomalies that incorporate copies, missing numbers, and interesting classes to guarantee the dataset's trustworthiness.
- The following center moves to utilize autoencoders, a kind of Artificial Neural Network (ANN), empowering unattended clustering (Camarrone and Van Hulle, 2019).
- Moreover, a custom clustering part is added to the model, working on its ability to relegate significant groupings in light of learned portrayals.

### Methods

#### Data Loading and Investigation:

- The review started bringing in the Netflix dataset utilizing the 'netflix\_titles.csv' record and dissecting its aspects, types, and plausible missing qualities.
- Copies were erased to keep up with data honesty. Pivotal bits of knowledge, for example, quarterly delivery examples and top substance suppliers, were outwardly portrayed utilizing diagrams.

#### Clustering Using Autoencoders:

- Autoencoders, a sort of neural organization, have been utilized for the extraction of elements and gathering.
- A remarkable clustering layer was executed to dispense bunches given learned portrayals. K-implies clustering made it simpler to assign starting clusters.

#### Results and Representation:

- After model preparation, this Netflix dataset was expanded with cluster tasks. Clustering results were shown utilizing static and intelligent diagrams (Eklund and Jong-Min, 2022).
- This considered an intensive understanding of content patterns and gave experiences about cluster scattering across the dataset.

### Result Implementation

#### GOAL 1: Data Loading and Exploration

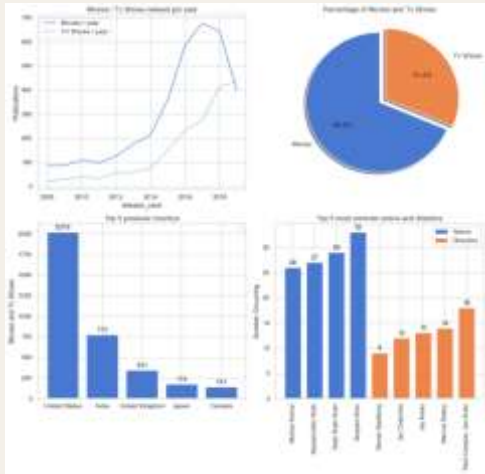
- Specifically, the 'director,' 'cast,' 'country,' 'date\_added,' and 'rating' columns have missing numbers, including 'director' having 1969, 'cast' 570, 'country' 476, 'date\_added' 11, and 'rating' 10 invalid qualities.

```
1.) The shape of the Dataset is (8028, 12), 8028 rows and 12 columns
2.) The dataset columns contain the following datatypes:
show_id    int64
type       object
title      object
director   object
cast       object
country    object
date_added object
release_year int64
rating     object
duration   object
listed_in  object
description object
other_text object
3.) Non nulls in the dataset:
show_id    8028
type       8028
title      8028
director   8028
cast       8028
country    476
date_added 11
release_year 8028
rating     10
duration   8028
listed_in  8028
description 8028
other_text 8028
```

- Replication in the dataset was found by title and purposefully eliminated, bringing about the deficiency of 62 columns.

```
4.) Check if there are duplicat titles in the dataset and remove the duplicates:
62 rows of duplicat titles have been removed
5.) Count number of unique genres:
There are 698 unique categories / genres in this dataset
```

- In this manner, a type evaluation affirmed the presence of 608 unmistakable classes or kinds inside the Netflix dataset, demonstrating the organization's differentiated substance contributions.



- The code area involves Python's Counter for investigating Netflix data, with an accentuation on happy examples (Jha *et al.* 2022). It isolates motion pictures and TV series, ascertains yearly counts, and perceives top makers, and entertainers, including directors.

#### GOAL 2: Clustering Section for Imports and text tokenizing

- The code makes a tokenizer utilizing the Keras Tokenizer class for getting ready portrayals of text given by the Netflix dataset.
- It confines the rundown of words to the 10,000 ordinarily noticed terms and arranges every depiction.

```
Epoch 1/8
WARNING:tensorflow:From C:\Users\Tech Assignment\AppData\Local\Microsoft\Windows\Apps\PythonSoftwareFoundation.Python.3.10.charmap01\python.exe: The name tf.nn.rnn_cell.LSTMCell is deprecated. Please use tf.nn.rnn_cell.LSTMCell instead.
40/40 [=====] - 4s 43ms/step - loss: 0.0100
Epoch 2/8
40/40 [=====] - 2s 42ms/step - loss: 0.0087
Epoch 3/8
40/40 [=====] - 2s 42ms/step - loss: 0.0085
Epoch 4/8
40/40 [=====] - 2s 43ms/step - loss: 0.0085
Epoch 5/8
40/40 [=====] - 2s 43ms/step - loss: 0.0085
Epoch 6/8
40/40 [=====] - 2s 43ms/step - loss: 0.0085
Epoch 7/8
40/40 [=====] - 2s 43ms/step - loss: 0.0085
Epoch 8/8
40/40 [=====] - 2s 43ms/step - loss: 0.0085
```

- Autoencoder engineering, which incorporates both interpreting and encoding parts, works with a specific number of hidden layers and numerous layers, which ultimately brings about a remaining hidden layer (Shayna *et al.* 2022).

```
103/103 [=====] - 1s 3ms/step
```

- Following approximately 0.0100, the misfortune dropped to 0.0085, showing that the info data was effectively remade and the autoencoder was pre-prepared.

- The review utilized an autoencoder towards clustering, with 20 clusters prepared of more than 8 ages and a bunch size is approximately 128. The autoencoder engineering has layers with aspects of [x.shape[1], 500, 500, 1000, 18].

- The statement cycle embraced the 'fan\_in' mode with uniform dissemination and the underlying preparation enhancer was 'rmsprop', with energy set at 0.9.

- The 'ClusteringLayer' class is an extraordinary part of the neural organization's plan that considers clustering utilizing K-implies with predefined clusters.

- It produces the probability dissemination (q) of clustered data focuses while constantly adjusting alpha qualities for ideal portrayal learning.

#### GOAL 3: Cluster prediction and Model Fitting

- The clustering model iteratively changed its boundaries with refreshing objective dispersions. Preparing utilized small-scale bunches to increment effectiveness.
- Moderate appraisals empower versatile learning. The last model has been put something aside for organization.

```
data_all['cluster'] = y_pred
data_all.head()
```

	show_id	type	title	director	cast	country	date_added	release_year	rating	duration	listed_in	description	cluster
9	81145628	Movie	Naam of the South King	Richard Yoo, Tim Voth	Joan Marrett, Andrew Toff, Brian Dobson, Cole	United States, India, South Korea, China	September 9, 2019	2018	TV-PG	90 min	Children & Family Movies, Canadian	Before planning an awesome wedding for his gr...	1
1	80117401	Movie	Jarhead: Whatever it Takes	Naal	Jarrod Argersoll	United Kingdom	September 9, 2019	2016	TV-MA	94 min	Stand-Up Comedy	Jarrod Argersoll shares his on-the-challenges of life...	8
2	70234438	TV Show	Transformers: Prime	Naal	Peter Cullen, Sumner, Mark Rolston, Frank	United States	September 9, 2019	2013	TV-PG	Season 1	Kids TV	With the help of three human allies, the Autob...	8
3	80858054	TV Show	Transformers: Robots in Disguise	Naal	Will Friedle, Darren Criss, Connor	United States	September 9, 2019	2018	TV-PG	Season 1	Kids TV	When a prison ship crash-lands hundreds of...	18
4	80125078	Movie	Awesetup	Forrest	Naal	United States	September 9, 2017	2017	TV-14	99 min	Canadian	When nearly high schooler can't really attract...	7

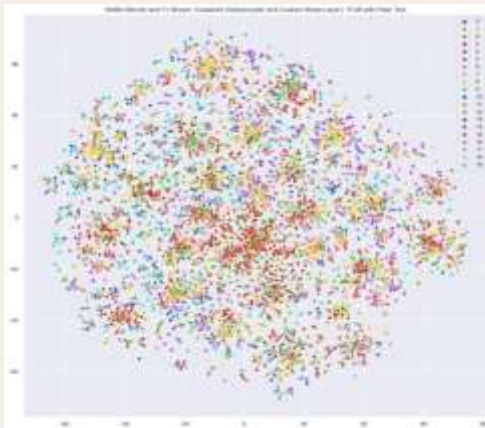
- The dataset incorporates a different scope of Netflix material, like motion pictures and television series.

- Clustering utilizing autoencoders uncovered unmistakable substance designs. For instance, Cluster 1 gives family-accommodating movies, Cluster 9 spotlights stand-up satire and an assortment of television series, while Cluster 7 incorporates teen-oriented films.

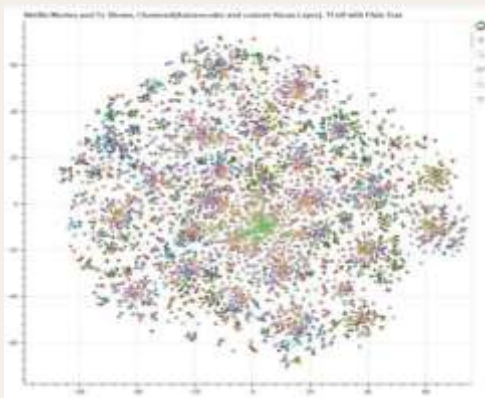
cluster	
12	691
1	583
17	513
8	420
8	417
13	378
16	366
18	322
4	298
7	257
10	274
15	272
14	256
19	230
0	210
2	208
5	206
11	176
6	85
0	22

Name: count, dtype: int64

- The Netflix dataset, which had 6,172 instances, has been streamlined into two aspects by utilizing t-distributed Sequential Neighbor Embedding (t-SNE).



- This modification took into account visual depiction and understanding of the dataset's fundamental cluster structure.



- A scatter plot for gathered Netflix films and TV series was made utilizing Bokeh, a Python electronic visualization bundle, as a result of autoencoder embedding.

### Conclusion

Finally, the classification and regression analyses provide insights into the complex world of movies and TV series.

The created clusters help to understand the similarities of titles what allows for content categorization and recommendation systems.

The nature of this work allows for the expression of some depth in highlighting factors that affect performance and reception by entertainment content providers as well as consumers.

The diversity of the dataset, covering a wide range of genres and regions, demonstrates that entertainment is a global phenomenon.

With the changing technology and viewer's preferences, the findings from this analysis will continue to add value and contribute to a more practical approach in content development, distribution, and audience engagement.