

Multi-Modal Genre Classification for Movies

By – Navneet Gupta ,Sahil Singhavi, Abhay Bhardwaj

Abstract—Using Machine learning models and algorithms to predict the genres of a particular movie. In this paper, we addressed the multi-label classification of the movie genres in a multimodal way. Machine learning is a branch of artificial intelligence (AI) and computer science which focuses on the use of data and algorithms to imitate the way that humans learn gradually improving its accuracy. Some real-world applications of machine learning are speech recognition, customer service like online chat bots, Alexa, google assistance, recommendation engines on YouTube etc. To elaborate further let's take movies. Humans can identify the genres of movie very easily. We only require bare minimum information such as posters or name of a movie or dialogs etc. to guess the genre of a movie and in majority of cases we get it right. The idea is to be able to the same with the computer using machine learning, but movie genre classification is a challenging task and is gaining attention a lot.

Index Terms— multi-modal, tensor-flow



1 INTRODUCTION

In today's world computer science has improved the standards of living significantly. We have smart phones laptops and other electronic gadgets combined with the Internet make our lives easier and efficient. One of the main technologies which makes these gadgets more user friendly and consumer satisfactory is Artificial intelligence (AI) and Machine learning.

Streaming media services have grown steadily over the past decade as a practical and comfortable way of allowing consumers to have access to films, series, documentaries, and so on. Big companies like Netflix, Amazon prime video, YouTube have hugely invested in this area for making it more user friendly and understanding the demands of the user like what type of genres they prefer and what are the best recommendations for a particular user. Etc. thus, for such type of features ML is the answer. It is an attempt to imitate what the human mind does. We can guess the genre of the movie just by looking at the movie title or poster dialogs etc. thus our perception is very high compared a machine. making a computer do all that work is the aim of ML. We need to test the accuracy of various ML models which predict the genre of a particular movie. For this here there are two types of data available one is textual that is subtitles and the other is visual that is posters. We scrap the data according to our requirements after that use the non-deep conventional model for the textual data.

To identify the genres from the posters we use deep learning.

Project Outline

- iii) Scraping dataset
- iv) Pre-processing dataset
- v) Introduction Deep Learning
- vi) ML Models for poster
- vii) ML Models for Textual data

2 METHODOLOGY

In ML, task of classification means to use available data to learn function which can assign category to data point.

For e.g, assign a genre to a movie, like "Romantic Comedy", "Action", "Thriller". Another example could be automatically assigning category to news articles, like "Sports" & "Politics".

Given:

- i) data point x_i
- ii) set of categories y_1, y_2, \dots, y_n that x_i belong to.

Task:

Predicting correct category y_k for new data point x_k not present in given dataset.

Problem:

don't know how the x and y are related mathematically.

Assumption:

Assuming exists a function f relating x and y i.e., $f(x_i) = y_i$

Approach:

f not known, learn a function g which approximates f .

Important consideration:

$f(x_i) = g(x_i) = y_i$ for all x_i , after that two functions f and g are exactly equal. This wont realistically ever happen, and we will only be able to approximate true function f using g . means, sometimes the prediction $g(x_i)$ will not be correct. And essentially, our whole goal is to find g which makes really low number of such errors. .

should mention that this is a specific kind of learning problem which we call "Supervised Learning". the idea that g approximates f well for data do not present in our dataset is called "Generalization". It is paramount that our model generalizes, or else all our claims will only be true about data we already have, and our predictions will not be correct.

There are many other kinds, but supervised learning is most popular and well-studied kind.

In ML network, Multi-Modal is refer to multiple kinds of data. For example, consider a YouTube video. It can be thought to contain 3 different modalities -

- i) video modality
- ii) audio modality
- iii) textual modality

3 LITERATURE REVIEW

We will scratch information from 2 unique film sources - IMDB and TMDB "IMDB:http://www.imdb.com/". For those uninformed, IMDB is the essential wellspring of data about motion pictures on the web. It is massively rich with banners, audits, rundown, appraisals, and numerous other data on each film. We will utilize this as our essential information source. "TMDB: https://www.themoviedb.org/". TMDB, or The Movie Database, is an open-source rendition of IMDB, with an allowed to utilize API that can be utilized to gather data. do require an API key, however it tends to be gotten free of charge simply by making a solicitation in the wake of making a free record.

```
!pip install tmdbsimple
!pip install wget
!pip install IMDbPY
```

The data is called from the website mention and how to setup the API and calling our running time data we have described in [Appendix1.1]

A comparison of IMDB and TMDB data

Now that we have both systems running, let's do a very short comparison for the same movie.

```
print ("The genres for The Matrix pulled from IMDB are -",movie['genres'])
print ("The genres for The Matrix pulled from TMDB are -",get_movie_genres_tmdb("The Matrix"))
```

[Appendix1.2]

Building a dataset to work

Utilizing similar API , we will simply pull results from the main 50 pages. As referenced before, the "page" characteristic of the order top movies=all_movies.popular() can be utilized for this reason.

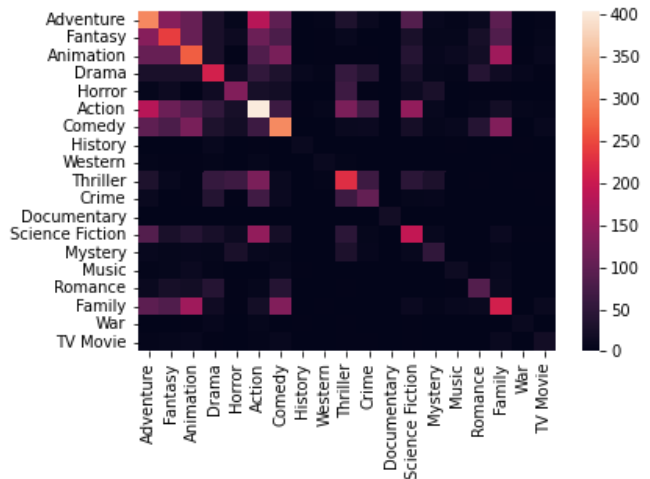
As the data is excessively huge with more properties in a portion of the code underneath will store the data into python "pickle" documents so it tends to be prepared straightforwardly from memory, instead of being downloaded without fail. Once done, should remark out any code which produced an article that was cured and is not generally required.

For reference check the code [Appendix 2.1]

Analysis of Movie Genres

Pairwise Analysis As our dataset is multi label, certainly searching at the distribution of genres is not enough. It might be beneficial to look which genres co-arise, as it'd shed a few light on inherent biases in our dataset. So, for the top one thousand movies allow's do a little pairwise evaluation for genre 'distributions'. Our fundamental purpose is to look which genres arise collectively within the equal movie. So, we first define a characteristic which

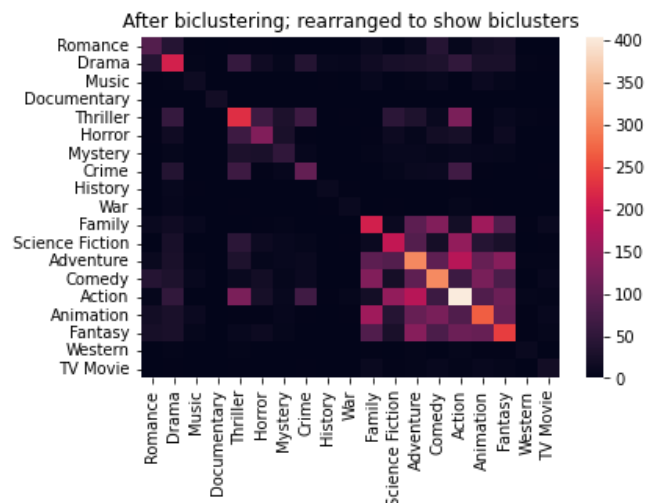
takes a listing and makes all viable pairs from it. After that, we pull the list of genres for a movie and run this function at the listing of genres to get all pairs of genres which arise together. We made the shape that's a 19X19 shape which shows 19 Genres. This shape counts the wide variety of simultaneous occurrences of genres in same film.



The above image indicates how frequently the genres arise collectively, as a heatmap essential factor to word within the above plot is the diagonal. The diagonal corresponds to self-pairs, i.e., quantity of instances a genre, say Drama took place with Drama. Which is basically just a be counted of the entire instances that style befell! As we are able to see there are a number of dramas within the records set, it is also a very unspecific label. There are nearly no documentaries or tv films. Horror is a distinct label, and romance is also now not too widely unfold.

To account for this unbalanced data, multiple things we can try to explore what interesting relationships can be found. Delving Deeper into co-occurrence of genres

we want to look for nice groups of genres that co-occur. while the information might not display that directly, we will play with the numbers to peer if that's feasible.



The method used for this is known as biclustering. Searching at the above figure, "boxes" or groups of movie

genres automatically emerge! Intuitively - Crime, Sci-Fi, thriller, action, Horror, Drama, thriller, and so forth. co-occur. AND, Romance, myth, circle of relatives, song, adventure, and many others. co-occur.

For code of above check [Appendix 2.4]

One assignment is the large variety of the drama style. It makes the two clusters exceptionally overlapping. If we merge it together with motion thriller, etc. we are able to emerge as with nearly all films just having that label.

based on playing round with the stuff above, we are able to sort the information into the subsequent style categories - "Drama, movement, technological know-how Fiction, thrilling (thriller, crime, mystery), uplifting (adventure, delusion, animation, comedy, romance, circle of relatives), Horror, datasets"

This categorization is subjective and in no way the simplest proper solution. One can also simply live with the original labels and best exclude the ones with now not enough facts. Such hints are vital to balance the dataset, it allows us to increase or decrease the strength of certain indicators, making it viable to improve our inferences.

The discussion or the things which comes in our mind after seeing the graph

- Is length correlated with film style?
- Are film posters darker for horror than for love end comedy?
- Are some genres in particular released extra frequently at a sure time of year?
- Is the RPG score correlated with the genre?

Based on this new category set, we will now pull posters from TMDB as our training data!

Check [Appendix 1.3]

Scraping done!

Dataset Building from scraped information!

Because the data is scraped now, we want it to apply it however nevertheless in pieces so we want to combine and clean the data. This task is easy, however extraordinarily vital. it's basically what's going to set the degree for the entire project. given that have the freedom to cast their own challenge in the framework supplying, there are numerous decisions that have to make to finalize very own model of the task.

As we are working on a classification problem, we need to make two decisions given the data at hand -

- viii) What do we want to predict, i.e., what's our Y?
- ix) What features to use for predicting this Y, i.e., what X should we use?

There are many different options possible, and it comes down to decide what's most exciting. I will be picking my own version for the example, but it is imperative that think this through, and come up with a version which excites!

As an example, here are some possible ways to frame Y, while still sticking to the problem of genre prediction -

x) Assume every movie can have multiple genres, and after that it becomes a multi-label classification problem. For example, a movie can be Action, Horror and Adventure simultaneously. Thus, every movie can be more than

one genre.

xi) Make clusters of genres as we did in Milestone 1 using biclustering, and after that every movie can have only 1 genre. This way, the problem becomes a simpler, multi-class problem. For example, a movie could have the class - Uplifting (refer Milestone 1), or Horror or History. No movie gets more than one class.

For the purposes of this implementation, I'm going with the first case explained above - i.e., a multi-label classification problem.

Similarly, for designing our input features i.e., X, may pick any features think make sense, for example, the Director of a movie may be a good predictor for genre. OR they may choose any features they design using algorithms like PCA. Given the richness of IMDB, TMDB and alternate sources like Wikipedia, there is a plethora of options available. Another important thing is that in doing so, we must also make many more small implementation decisions on the way. For example, what genres are we going to include? what movies are we going to include? All these are open ended!

As we on a classification problem, we want to make two choices given the records at hand -

- What do we want to predict, i.e., what's our Y?
- Features to apply for predicting the Y, i.e., what X to use?

there are numerous distinct options feasible, and it comes all the way down to determine what is most thrilling. I can be choosing my own version for the instance, but it is vital that suppose this via, and give you a version which excites!

As an example, right here are a few viable approaches to frame Y, whilst nonetheless sticking to the problem of genre prediction -

- count on every movie may have more than one genres, after which it will become a multi-label classification problem. for instance, a film can be motion, Horror and adventure concurrently. hence, each movie may be a couple of style.

• Make clusters of genres as we did in Milestone 1 the use of biclustering, after which every movie can have most effective 1 genre. This manner, the problem becomes a simpler, multi-class problem. as an example, a movie ought to have the class - Uplifting (refer Milestone 1), or Horror or history. No film gets more than one class.

For the purposes of this implementation, we going with the first case defined above - i.e., a multi-label type problem.

similarly, for designing our enter functions i.e., X, may pick any feature think make sense,

Implementation

Implementation decisions made -

1. The problem is framed here as a multi-label problem explained above.
2. We will try to predict multiple genres associated with a movie. This will be our Y.
3. We will use 2 different kinds of X - text and images.
4. For the text part - Input features being used to

predict the genre is a form of the movie's plot available from TMDb using the property 'overview'. This will be our X.

5. For the image part - we will use the scraped poster images as our X.
- 6.

We will first look at some conventional machine learning models, which were popular before the recent rise of neural networks and deep learning. For the poster image to genre prediction, I have avoided using this because conventional ML models are simply not used anymore without using deep learning for feature extraction (all discussed in detail ahead, don't be scared by the jargon). For the movie overview to genre prediction problem, we will look at both conventional models and deep learning models.

let's build our X and Y!

First, let's identify movies that have overviews. Next few steps are going to be a good example on why data cleaning is important!

Now let's store the genres for these movies in a list that we will later transform into a binarized vector.

Binarized vector representation is a very common and important way data is stored/represented in ML. Essentially, it's a way to reduce a categorical variable with n possible values to n binary indicator variables. What does that mean? For example, let [(1,3), (4)] be the list saying that sample A has two labels 1 and 3, and sample B has one label 4. For every sample, for every possible label, the representation is simply 1 if it has that label, and 0 if it doesn't have that label. So, the binarized version of the above list will be -

Implementation decisions made -

- 1) The problem is framed right here as a multi-label problem explained above.
- 2) we can try and are predicting multiple genres associated with a film. this could be our Y.
- 3) we will use 2 exclusive kind of X - text and pics.
- 4) For the text part - entered features being used to predicting the genre is a shape of the movie's plot available from TMDb using the property 'review'. this will be our X.

- For image part - we can use the scraped poster image as our X.

we are able to first look at some traditional machine learning models, which were famous before the recent upward push of neural networks and deep learning. For the poster photo to style prediction, we have avoided the usage of this due to the fact conventional ML models are clearly not used anymore without the use of deep learning for feature extraction. For the movie overview to genre prediction problem, we can have a look at both conventional and deep learning models.

let's construct our X and Y!

First, let's discover films which have overviews.

Now allow's keep the genres for those films in a listing that we are able to later rework into a binarized vector.

Binarized vector illustration is a common and essential manner data is saved/represented in ML. basically, its

wayr to lessen a express variable with n feasible values to n binary indicator variables.

[(1,0,1,0)], (0,0,0,1)]

See [Appendix 1.4]

Storing movie overview in a matrix

The way we will be storing the X matrix is called a "Bag of words" representation. The basic idea of this representation in our context is that we can think of all the distinct words that are possible in the movies' reviews as a distinct object. And after that every movie overview can be thought as a "Bag" containing a bunch of these possible objects.

The simple concept of this representation in our context is that we are able to think of all the distinct words which can be viable inside the films' reviews as a awesome object. and after that every movie assessment may be idea as a "Bag" containing a gaggle of these viable objects.

What this means is that, for all the movies that we have the data on, we will first count all the unique words. Say, there's 30,000 unique words. After that we can represent every movie overview as a 30000x1 vector, where each position in the vector corresponds to the presence or absence of a particular word. If the word corresponding to that position is present in the overview, that position will have 1, otherwise it will be 0.

What this indicates is that, for all the movies that we've the statistics on, we will first remember all of the particular phrases. Say, there may be 30,000 particular phrases. After that we will constitute every film overview as a 30000x1 vector, wherein every role in the vector corresponds to the presence or absence of a specific phrase. If the word corresponding to that role is present in overview, that position will have 1, in any other case it will likely be zero.

[appendix 2.5]

For example, let's consider the overview for Matrix Movie -

```
get_movie_info_tmdb('The Matrix')['overview']
```

For "The Matrix" a phrase like "computer" is a more potent indicator of it being a Sci-Fi film, than phrases like "who" or "effective" or "huge". One way computer scientists operating with natural language tackled this trouble inside the beyond (and it is still used very popularly) is what we name TF-IDF i.e., time period Frequency, Inverse file Frequency. The fundamental concept here is that phrases which are strongly indicative of the content of a single document (every film assessment is a record in our case) are words that occur very frequently in that document, and really infrequently in all other documents.

For three variables, it'd be one hundred thousand, and we might want to pattern at 500,000 points. that is already extra than the range of information points we've for maximum training problem we are able to ever stumble upon.

essentially, as the dimensionality (number of functions) of the examples grows, due to the fact a hard and fast-size training set covers a dwindling fraction of the input space. even with a mild size of 100 and a large training set

of a thousand billion examples, the latter covers handiest a fragment of approximately 10–18 of the input area. this is what makes device studying both essential and difficult.

So, yes, if a few words are unimportant, we need to do away with them and reduce the dimensionality of our X matrix. And the manner we will do it is the use of TF-IDF to become aware of un-crucial words. Python shall we us try this with simply one line of code

```
vectorizer=CountVectorizer(max_df=0.95, min_df=0.005)
X=vectorizer.fit_transform(content)
```

We are excluding all words that occur in too many or too few documents, as these are very unlikely to be discriminative. Words that only occur in one document most probably are names, and words that occur in nearly all documents are probably stop words. Note that the values here were not tuned using a validation set. They are just guessing. It is ok to do because we didn't evaluate the performance of these parameters. In a strict case, for example for a publication, it would be better to tune these as well.

So, each movie's overview gets represented by a 1x1201 dimensional vector.

Now, we are ready for the kill. Our data is cleaned, hypothesis is set (Overview can predict movie genre), and the feature/output vectors are prepped. Let's train some models!

we are apart from all phrases that arise in too many or too few documents, as these are impossible to be discriminative. words that most effective arise in a single file maximum probably are names, and words that arise in almost all documents are stop words. note that the values here were not tuned the usage of a validation set. they're simply guessing. it's miles ok to do due to the fact we didn't evaluate the performance of those parameters. In a strict case, as an example for a booklet, it would be better to tune those as properly.

So, each film's overview gets represented through a 1x1201 dimensional vector.

we're ready for the kill. data is cleaned, hypothesis is about (overview predicting movie genre), and the feature/output vectors are prepped. let's train!

[Appendix 2.6]

Data set ready

we are building our own dataset, and didn't want to spend all time waiting for poster image downloads to finish, we are working with extremely small dataset. That is why, the results we will will not be spectacular as compared to conventional machine learning methods.

Non-deep, Conventional ML models with above data

Here is what we will be doing -

1,implementing two different models

2,Deciding a performance metric

3. Difference between the models, their strengths, weaknesses, etc.(Discussion)

There are many implementation decisions . Between feature engineering, hyper-parameter tuning, model selection and how interpretable do want model to be (Read:

Bayesian vs Non-Bayesian approaches) a lot is to be decided.

For our purpose , showing the example of 2 very simple models,

i)SVM

ii)Multinomial Naive Bayes

whole pipeline below

i)A little bit of feature engineering

ii) different Models

iii)Evaluation Metrics chosen

iv)Model comparisons

Second, can only represent based on the data at hand.

A nice way to think of it is to think that start with the problem at hand, but design features constrained by the data have available. If have many independent features that each correlate well with the class, learning is easy. On the other hand, if the class is a very complex function of the features, may not be able to learn it.

in the context of the problem, we would like to predict genre of a film. what we have get admission to to - movie overviews, which might be textual content descriptions of the film plot. The speculation makes sense, review is a quick description of the story and the story is certainly essential in assigning genres to movies.

So, lets enhance our functions by way of gambling with the words within the overviews in our datas. One interesting way to head returned to what we mentioned earlier - TF-IDF. We originally used it to clear out word, however we can also assign the tied values as "significance" values to words, as opposed to treating all of them equally. Tied truly attempts to become aware of the assign a weightage to every word in the bag of phrases.

yet again, the manner it really works is - most film descriptions have the word "The" in it. obviously, it would not tell whatever unique about it. So, weightage have to be inversely proportional to how many movies have the word of their description. that is the IDF element.

However, for the film interstellar, if the description has the phrase space 5 instances, and wormhole 2 instances, after that it is likely greater about space than about wormhole. thus, space must have a excessive weightage. that is TF element.

We genuinely use TF-If to assign weightage to every word inside the bag of phrases

```
from sklearn.feature_extraction.text import TfidfTransformer
tfidf_transformer = TfidfTransformer()
X_tfidf = tfidf_transformer.fit_transform(X)
X_tfidf.shape
```

Output

(1693, 1201)

allow's divide our X and Y matrices into train and test split. Training the model on train split and report the performance at the test slip . consider this like the questions do inside the prolem units v/s the exam. Of direction, they are both (assumed to be) from the same populace of

questions. And doing well on problem units is a good indicator that 'do well in tests, however simply, ought to check before can make any claims approximately knowing the subject.

[Appendix 1.5]

inside the output we observe the performance is via and big poorer for movies which can be less represented like battle and animation, and better for classes like Drama.

Numbers apart, permit's have a look at our model predictions for a small pattern of movies from our test set.

As seen above, the outcomes appear promising, however how do we virtually compare the 2 models? We need to quantify our overall performance so that we are able to say which one's better. Takes us again to what we mentioned right within the beginning – we learning the function g which could approximate the unknown function f . For a few values of x_i , the predictions can be wrong for certain, and we need to reduce it.

For multi label structures, we often keep tune of overall performance using "Precision" and "recall". these are popular metrics and may google to study up greater approximately them if there are new to these words.

[Appendix 1.6]

Here is Evaluation Matrix

Using standard precision recall metrics for evaluating systems average precision and recall score for samples are pretty good! Model working!

[appendix1.6]

Deep Learning(Overview)

Short concept approximately what is deep learning.

As described above, two most crucial principles in doing excellent classification (or regression) are to

- 1) Use the proper illustration which captures the info about the data which is which is relevant to the problem at hand
- 2) the use of the right model which has the capability of making sense of the representation fed to it.

To simply emphasize the importance of illustration in the complex tasks we generally try with Deep Learning, let us talk the original problem which made it famous. The paper is frequently called the "ImageNet task Paper", and it was essentially working on object recognition in pictures.

The "Deep" within the deep learning comes from the fact that it turned into conventionally carried out to Neural Networks. Neural Networks, as we all recognise, are structures organized in layers. Layers of computations. Why do we need layers? due to the fact those layers be sub-tasks that we do within the complicated task of recognizing out a chair. it is able to be thought as a hierarchical split down of a complex task into smelled sub-tasks.

Mathematically, each layer acts like space transformation which takes the pixel value to a excessive dimensional space. whilst we start off, every pixel inside the image is given identical importance in our matrix. With each layer, convolution operations deliver a few elements greater importance, and a few lesser significance. In doing

so, we rework our images to space in which similar searching objects/object parts are closer

What precisely became learnt by those neural networks is tough to know, and an active vicinity of research. but one very crude way to visualize what it does is to assume like - It starts offevolved by generic features within the first layer. something as simple as vertical and horizontal traces. within the next layer, it learns that if integrate the vectors representing vertical and horizontal vectors in different ratios, can make all feasible slanted strains. subsequent layer learns to combine lines to form curves - Say, some thing like the define of a face. those curves come together to form 3-D objects. and so forth. building sub-modules, combining them within the proper manner that deliver it semantics.

So, in a nutshell, the primary few layers of a "Deep" community examine the proper illustration of the facts, given the problem (that is mathematically defined with the aid of goal characteristic seeking to decrease difference between ground truth and predicted labels). The last layer absolutely seems how near or far apart things are in this high dimensional space.

subsequently, we will supply any sort of data a high dimensional illustration the usage of neural networks. beneath we can see high dimensional representations of each words in overviews (text) and posters (photograph). allow's get commenced with the posters i.e., extracting visual features from posters using deep learning

Predicting genre from poster

once again, we have to make an implementation decision. This time, it has greater to do with how tons time are we willing to spend in return for brought accuracy. we are going to use here a way this is generally called Pre-Training in machine Learning.

in place of seeking to re-invent the wheel here, going to borrow this short segment on pre-training from Stanford college's lecture on CNN's. to quote -

"In practice, very few people train an entire Convolutional Network from scratch (with random initialization), because it is relatively rare to have a dataset of sufficient size. Instead, it is common to pretrain a Convent on a very large dataset (e.g., ImageNet, which contains 1.2 million images with 1000 categories), and after that use the Convent either as an initialization or a fixed feature extractor for the task of interest. "

What this indicates, is that within the area spacestack transforms the snap shots to, all photographs which include a "dogse" are closer to each other, and all photos containing a "cat" are closer. consequently, it is a significant area wherein photographs with similar objects are closer.

think about it, now if we pump our posters through this stack, it will embed them in space where posters which include similar objects are closer While a smiling couple could point to romance or drama. The opportunity might be to train the CNN from scratch that is computationally intensive and entails a whole lot of tricks to get the CNN training to converge to the most efficient area transformation.

This manner, we can start off with something robust, after which build on top. We pump our pictures via the pre-trained network to extract the visual features from the posters. After that, the use of those functions as descriptors for the picture, and genres as labels, we train a less difficult neural network from scratch which learns to classification in this dataset. these 2 steps are precisely what we're going to do for predicting genres from film posters..

Extracting visual features from posters

The problem here we answering is can we use poster to predict genre. is hypothesis make sense? Yes. This what graphic designers do g. leaving visual cues to semantics. They make sure that when we look poster of horror movie, we can know it's not smiling image. Can our deep learning system infer such subtleties?

either we train deep neural network ourselves from scratch, or we can use a pre-trained made available to us from Visual Geometry Group at Oxford University, one of the most popular method. called the VGG-net. , it's just space transformation in form of layers.

Kera's is library that makes it easy for us. Some other are TensorFlow and PyTorch. While the Kera's makes it easy for prototype by keeping the syntax simple. Some common ways people refer to step are - "Getting the VGG features of an image", or "Forward Propagating the image through VGG and chopping off the last layer".

[Appendix2.7]

Training a simple neural network model using these VGG features.

let's first get labels on 1570 samples As picture download fails on a few instances, the high-quality way to work with the right model is to examine the poster names downloaded and running from there. these posters can't be uploaded to GitHub as they're too large size

[Appendix 1.7]

The final movie poster set for which we have all the information we need, is movies.

In above code we're making an X NumPy array containing the visual feature of 1 picture in step with row. VGG function are reshaped to be in form (1,25088) and we obtain a matrix of shape (1553,25088).

Our binarized Y NumPy array includes binarized labels corresponding to genre IDs of 1553 films.

Now, we create our own koras neural network to use the VGG functions and after that classify movie genres.

Neural network architectures have gotten complicated over the years. however the simplest ones comprise very wellknown computations organized in layers, as described above. Given the popularity of a number of these, Kera's makes it as smooth as writing out the names of those operations in a sequential order. This manner can make network while completely keeping off Math

Sequential () lets in us to make models the follow this sequential order of layers. unique forms of layers like Dense, Conv2D etc. can be used, and many activation features like RELU, Linear etc. also are to be had.

Let's train our model after that from features extracted from VGG net

Just one hidden layer between the VGG features and the final output layer. The simplest neural network can get. An image goes into this network with the dimensions (1,25088), the first layer's output is 1024 dimensional. This hidden layer output undergoes a pointwise RELU activation. output gets transformed into the output layer of 20 dimensions. It goes through a sigmoid.

The sigmoid, is a function which squashes numbers between 0 and 1.

By squashing score of each 20 output labels between 0 and 1, sigmoid lets us interpret their scores as probabilities. After that we can pick the classes with top 3 or 5 probability scores as predicted genres for movie poster

[Appendix 1.8]

Trained the model using the fit () function. It takes these parameters –

training features and training labels, epochs, batch size and verbose.

verbose. 0="don't print anything as work", 1="Inform me as go".

data set is too large to be loaded into RAM. we load data in batches. For batch size=32 and epochs=10, model loading starts rows from X in batches of 32 every time it calculates the loss and updates the model.and keeps going until it has covered all the samples 10 times.

no. of times model updated = (Total Samples/Batch Size) x (Epochs)

```
mod-
el_visual.fit(X_train, Y_train, epochs=10, batch_size=64,
verbose=1)
```

output

```
Epoch 1/10
20/20 [=====] - 7s
270ms/step - loss: 2.4439 - accuracy: 0.1200
Epoch 2/10
20/20 [=====] - 5s
267ms/step - loss: 0.3075 - accuracy: 0.3090
Epoch 3/10
20/20 [=====] - 5s
267ms/step - loss: 0.0958 - accuracy: 0.3558
Epoch 4/10
20/20 [=====] - 5s
266ms/step - loss: 0.0776 - accuracy: 0.3739
Epoch 5/10
20/20 [=====] - 5s
274ms/step - loss: 0.0662 - accuracy: 0.3591
Epoch 6/10
20/20 [=====] - 5s
271ms/step - loss: 0.0497 - accuracy: 0.3722
Epoch 7/10
20/20 [=====] - 6s
275ms/step - loss: 0.0419 - accuracy: 0.3895
Epoch 8/10
```

```
20/20 [=====] - 5s
269ms/step - loss: 0.0530 - accuracy: 0.3509
Epoch 9/10
20/20 [=====] - 5s
266ms/step - loss: 0.0358 - accuracy: 0.3615
Epoch 10/10
20/20 [=====] - 5s
266ms/step - loss: 0.0333 - accuracy: 0.3541
<keras.callbacks.History at 0x7fbc72f27f10>
```

```
mod-
el_visual.fit(X_train, Y_train, epochs=50, batch_size=64,
verbose=0)
```

For top 10 epochs trained model in a verbose fashion for telling what's happening

After I turned off the verbosity which keeps the code cleaner.

```
Y_preds=model_visual.predict(X_test)
sum(sum(Y_preds))
output
663.1564311981201
```

4 PREDICTIONS

Some Predictions

So, even with just the poster of course text outperforms the visual features, but the thing is that it still works. In more complicated models, we can combine the two to make even better predictions.

```
Predicted: Adventure,Comedy,Fantasy Actual:
Adventure,Fantasy,Action
Predicted: Animation,Adventure,Action Actual:
Animation,Action,Science Fiction,Family
Predicted: Animation,Action,Adventure Actual:
Adventure,Fantasy,Action
Predicted: Drama,Action,Thriller Actual: Ac-
tion,Western,Science Fiction
Predicted: Mystery,Music,Documentary Actual:
Documentary
Predicted: Comedy,Crime,Music Actual: Docu-
mentary,Music
Predicted: Family,Science Fiction,Comedy Actual:
Adventure,Comedy,Science Fiction,Family
```

These models were trained in CPU's, and a simple 1-layer model was used to show because there is a lot of information in this data that the models can extract. With a larger dataset, and more training we are able to bring these numbers to as high as 70%, same as textual.

5 TEXTUAL FEATURES

We will use a representation for words - Word2Vec

model. As the total number of words is small, we even don't need to forward propagate our sample.

```
from gensim import models
# mod-
el2 = models.Word2Vec.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=True)
mod-
el2 = models.KeyedVectors.load_word2vec_format('/content/drive/MyDrive/AI project/GoogleNews-vectors-negative300.bin.gz', binary=True)
```

Even that has been done for us, and the result is stored in the form of a dictionary. simply just look up word in the dictionary and get the Word2Vec features for the word. For eg.

```
print(model2['king'].shape)
print(model2['dog'].shape)
```

output

```
(300,)
(300,)
```

Words can be represented in overview using word2vec model. After that, we can use that as our X representations. So, instead of count of words, we are using a representation which is based on the semantic representation of the word. Mathematically, each word went from 3 dimensional length to 300 dimensions!

For the same set of movies above, let's predict the genres from deep representation of overviews!

Textual content desires some pre-processing earlier than we are able to teach the version. The best preprocessing we do here is - we delete commonly taking place words which we know aren't informative about the style. think of it as the litter in a few experience. those words are frequently removed and are called "stop phrases". can look them up on line. these encompass simple phrases like "a", "and", "but", "how", "or" and so on. They may be effortlessly eliminated using the python bundle NLTK.

```
!pip install stop-words
```

output

Collecting stop-words

Downloading stop-words-2018.7.23.tar.gz (31 kB)

Building wheels for collected packages: stop-words

Building wheel for stop-words (setup.py) ... done

Created wheel for stop-words: filename=stop_words-2018.7.23-py3-none-any.whl size=32912

sha256=b3dcc4320580e7328700862016eb77cd142cd5a9037720713cdcb680e16c373

Stored in directory:

/root/.cache/pip/wheels/fb/86/b2/277b10b1ce9f73ce15059bf6975d4547cc4ec3feeb651978e9

Successfully built stop-words

Installing collected packages: stop-words

Successfully installed stop-words-2018.7.23

```
from nltk.tokenize import RegexpTokenizer
```



```

from stop_words import get_stop_words
tokenizer = RegexpTokenizer(r'\w+')

# create English stop words list
en_stop = get_stop_words('en')

movie_mean_wordvec=np.zeros((len(final_movies_set),300))

movie_mean_wordvec.shape
print(final_movies_set)

output
[{'adult': False, 'backdrop_path':
'/6MQmtWk4cFwSDyNvIgoJRBiHUT3.jpg', 'genre_ids': [14, 28, 12], 'id': 559, 'original_language': 'en',
'original_title': 'Spider-Man 3', 'overview': 'The seemingly invincible S.....'}]]

```

From the above dataset, films with overviews which comprise handiest stop words, or movies with overviews containing no words with word2vec representation are unnoticed. Others are used to construct our mean word2vec representation. placed for every film overview

i)Take film overview
 ii)Throw stop phrases
 iii)For nonstop words:
 tv)If in word2vec - take it is word2vec representation that is three hundred dimensional If no longer - throw word

v)For each film, calculate the arithmetic imply of the 300-dimensional vector representations for all phrases in the assessment which weren't thrown out

300-dimensional representation for movies. All these are stored in a NumPy array.

X matrix becomes (1553,300). And Y (1553,19) i.e., binarized 19 genres

[Appendix 1.9]

7 CONCLUSION

In this project we have tried to solve the multimodal movie genre classification as a multi-label classification problem. To perform the classification, we used different sources of information, namely movie overview and poster of the movie. We scrapped data information from TMDb and IMDb APIs and build the data set from it. from the DENSENET 169 modal we achieved high F1 score. With the help of F1 score calculate we differentiated the movie genre class. We learned the technique clustering and bi-clustering with this regroup genre classes of the movie. We used the SVM and multinomial naïve bayes modal. We used deep learning to extract visual features from poster of the movies (VGCNET). Deep learning helps us understand the textual feature from overview of the movie. The average precision is 0.53 and recall score is 0.54 for our sample from SVM and Multinomial Naive Bayes

modal. From word2vec average precision 0.539 and recall score 0.580. These models were trained on CPU's, and a simple 1-layer model was used to show that there is a lot of information in this data that the models can extract. With a larger dataset, and more training we able to bring these numbers to as high as 70%.

Our predictions for the movies are -

Predicted: ['Adventure', 'Comedy', 'Fantasy']
 Actual: ['Adventure', 'Fantasy', 'Action']

Predicted: ['Animation', 'Adventure', 'Action']
 Actual: ['Animation', 'Action', 'Science Fiction', 'Family']

Predicted: ['Animation', 'Action', 'Adventure']
 Actual: ['Adventure', 'Fantasy', 'Action']

Predicted: ['Drama', 'Action', 'Thriller'] Actual: ['Action', 'Western', 'Science Fiction']

Predicted: ['Mystery', 'Music', 'Documentary']
 Actual: ['Documentary']

Predicted: ['Comedy', 'Crime', 'Music'] Actual: ['Documentary', 'Music']

Predicted: ['Family', 'Science Fiction', 'Comedy'] Actual: ['Adventure', 'Comedy', 'Science Fiction', 'Family']

REFERENCES

1. <https://wikipedia.org/>
2. Howard Zhou, Tucker Hermans, Asmita V Karandikar, and James M Rehg. Movie genre classification via scene categorization. In Proceedings of the 18th ACM international conference on Multimedia, 2010.
3. Giuseppe Portolese and Valéria Delisandra Feltrim. On the use of synopsis-based features for film genre classification. In Anais do XV Encontro Nacional de Inteligência Artificial e Computacional, SBC, 2018.
4. Gabriel S Simões, Jônatas Wehrmann, Rodrigo C Barros, and Duncan D Ruiz. Movie genre classification with convolutional neural networks. In 2016 International Joint Conference on Neural Networks (IJCNN) IEEE, 2016.
5. Rafael C Gonzalez, Richard E Woods, et al. Digital image processing. Prentice hall Upper Saddle River, 2002.
6. Beth Logan et al. Mel frequency cepstral coefficients for music modeling. In ISMIR, 2000.
7. Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. Neural Computation, 1997.
8. Klaus Greff, Rupesh K Srivastava, Jan Koutník, Bas R Steunebrink, and Jürgen Schmidhuber. Lstm: A search space odyssey. IEEE Transactions on Neural Networks and Learning Systems, 2017.

GITHUB LINK

*HTTPS://GITHUB.COM/TRUTHSEARCHERS/ML_PR
OBJECT*

\