

MOVIE RECOMMENDATION

Project proposal Using [MovieLens](#) dataset

Table of content:

1. Introduction
2. Literature Review
3. The Dataset
 - A. Users Dataset
 - B. Ratings Dataset
 - C. Movies Dataset
4. Features Selection and Exploratory Data Analysis (EDA)
5. Experiments
6. Methodology & Results
7. Application
8. Bibliography

1. Introduction:

Recommendation systems are very important applications of Machine learning that aim to improve the quality of the experience provided to the internet user. There are many available movies everywhere that users can watch. Recommendation systems are an important part of suggesting items especially in streaming services. Netflix, for example, uses these systems to help users find the movies that they may like. Similarly, this paper is suggesting different approaches for The MovieLens streaming service. Given a list of users, movies and ratings with their features, we want to build a recommendation system that takes the most informative features from the user's information and movies history to recommend for him/her other movies that he/she might enjoy.

1. Literature Review:

There are two well-known recommendations systems approaches which are the Collaborative Filtering and the Content Based Filtering. The Collaborative Filtering Approach is basically recommending new movies for the user based on how similar this user is to the rest of the users. The Content Based Filtering is basically focusing on the content of the movie being watched. So, if the user watched many Comedy movies, similar movies genres should be recommended. Many researchers contribute to this recommendation systems field by using different Machine and Deep learning algorithms to achieve higher accuracies. Lund and Ng (2018) use deep learning approaches in the Collaborative Filtering approaches and compared it to approaches like K-NN, and they got a lower RMSE. Vilakone & Park (2018) also used new approaches other than the K-NN and Cosine Similarity, k-clique methods and they called it the improved k-clique method. Similar projects and proposals suggested by Zhang, Yao et.al (2018) in the University of New South Wales. Most researchers are interested to beat the baseline and trivial approaches using the Deep learning and improved models approaches. In this paper, we are interested in the Collaborative Filtering Approach as well. We will use different approaches other than the ones suggested by the authors. Moreover, we will shift the problem into a classification problem, whether the user will dislike or like the movie.

2. Dataset Description:

Data Source: <http://www.grouplens.org/node/73>

This is the MovieLens Dataset. It has three tables: Users, Movies and Ratings. The Users data provides us with some information about the current MovieLens Users. The Movies data provides us with information about the current Movies in the MovieLens platform. The ratings data is the users' ratings for the movies they watched.

A. Users:

The users data consists of some information about the current users the MovieLens have.

Variables Description:

- User ID : Unique ID for the user
- Gender : 'F' for female , 'M' for male.
- Age : 6 different values representing 6 different age ranges.
 - 18 represent age in the range "18-24"
 - 25 represent age in the range "25-34"
 - 35 represent age in the range "35-44"
 - 45 represent age in the range "45-49"
 - 50 represent age in the range "50-55"
 - 56 represent age in the range "56+"
- Occupation: 20 Different fields cover all occupations in the Dataset
 - 0: "other" or not specified
 - 1: "academic/educator"
 - 2: "artist"
 - 3: "clerical/admin"
 - 4: "college/grad student"
 - 5: "customer service"
 - 6: "doctor/health care"
 - 7: "executive/managerial"
 - 8: "farmer"
 - 9: "homemaker"

- 10: "K-12 student"
- 11: "lawyer"
- 12: "programmer"
- 13: "retired"
- 14: "sales/marketing"
- 15: "scientist"
- 16: "self-employed"
- 17: "technician/engineer"
- 18: "tradesman/craftsman"
- 19: "unemployed"
- 20: "writer"
- Zip code

B. Ratings:

This dataset provides information about the users' ratings for the movies they watched.

- UserID: Unique
- MovieID:Unique
- Rating : on scale from 0 to 5
- Timestamp

C. Movies:

This dataset provides information about the movies' features.

- Movie ID: Unique
- Movie Name: has the format Movie Name (Year)
- Year of Production
- Genres: Comma separated values representing the different genres of the movie (The dataset has 20 different genres)
 - Action
 - Adventure
 - Animation
 - Children's
 - Comedy
 - Crime

- Documentary
- Drama
- Fantasy
- Film-Noir
- Horror
- Musical
- Mystery
- Romance
- Sci-Fi
- Thriller
- War
- Western

3. Features Selection and Exploratory Data Analysis (EDA):

The features that we are going to use are:

From User Data: [age, occupation, gender, Zip Code]

From Ratings Data: Rating

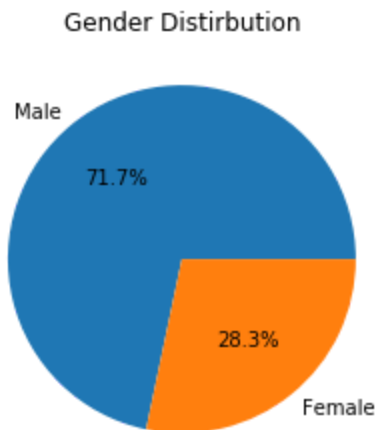
From Movies Data: [Genres, Year of Production]

In this section, we analyzed each variable separately, which helped us to understand our dataset more. We started by visualizing our data and notice the distribution of each variable.

A. Users Data

1.1. Gender

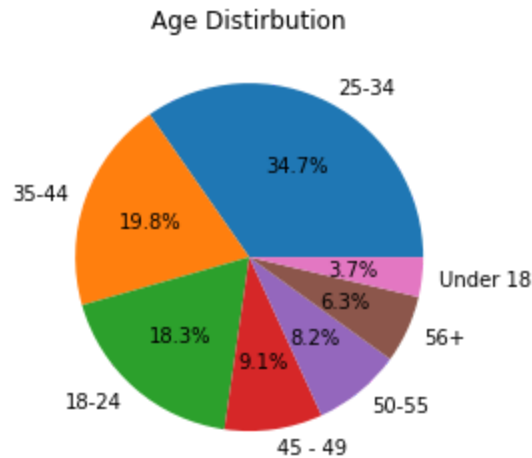
The below figure shows that the users data has more males (71.7%) than females (28.3%). We have unbalanced data. It might affect our Recommendation system. Various oversampling techniques will be performed in the next phases.



1.2. Age

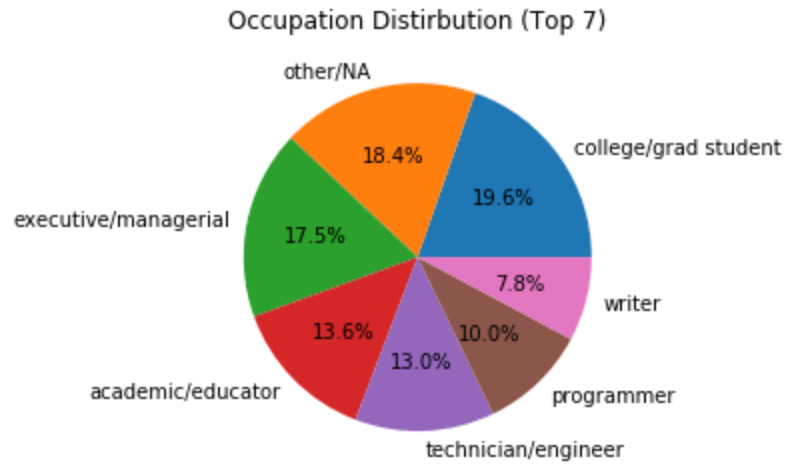
The below figure shows that the 25-34 years old category has the highest number of users (34.7%) among the other age categories. It was unexpected to see the

“Under 18” category has the lowest percentage of users. However, since MovieLens is an online platform. It might be the case that users below 18 years old do not really want to mention their true ages to avoid parental control issues or the fear of being prevented from creating their own accounts or being labeled as underage. As for the other categories, it was expected to see low percentages of users in the 45-49, 50-55, 56+ categories since they might not be interested in watching movies as the 18-24, 25-34 and 35-44 users do.

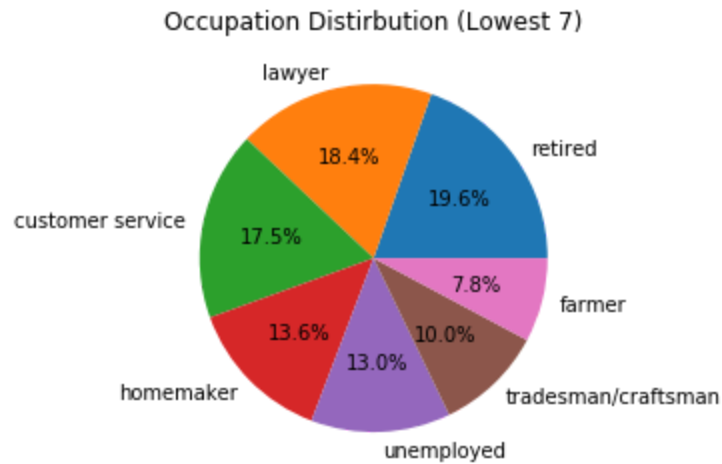


1.3. Occupation

The below figure shows the top 7 Users Occupation out of 20. Given the age distribution with the 25-34 is the highest percentage, it is reasonable to see that the college/grad student has the highest percentage (19.6%). Moreover, the college/ grad represents a significant part of the population. Other categories have lower percentage simply because they are not present in the population with the same percentage as students are. However, it is interesting to see 7.8% writers among the top 7 Occupations, writers are not present in the actual population as graduate students are. Once again we have an unbalanced data problem. Moreover, many users did not mention their Occupation (other/NA), this proves the point that it might be the case that the below 18 years users did not want to record their true ages.



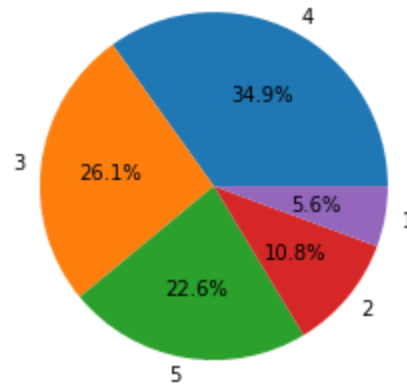
The below chart shows the lowest 7 Users Occupation. It was expected to see “unemployed” as one of the lowest 7 Occupation status (13%). People usually feel ashamed mentioning this type of information. It’s also not surprising to see such categories.



B. Rating Data

From the figure below, most users give the movies a rating of “4” (34.9%), and It is not common that the users give a rating of “1” (5.6%). The ratings data is negatively skewed.

Rating Distribution

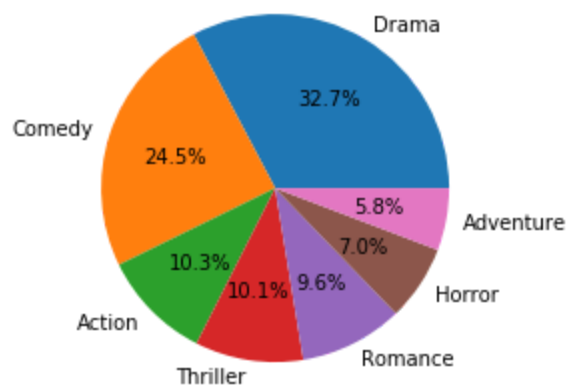


C. Movies Data

1.4. Genres Distribution

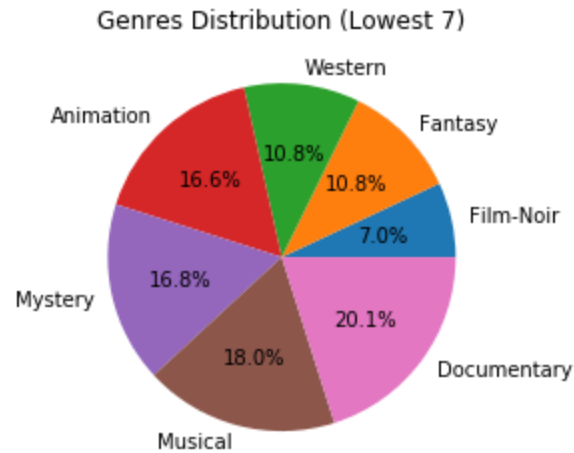
The figure below shows the Top 7 genres out of 18 Genres. As we can see, most of the movies have Drama (32.7%) and Comedy (24.5%) scenes. We can see the low percentage of Adventure (5.8%) and Horror (7.0%) because it is not common to have an Horror movie that can be labeled as “Comedy” or “Drama”.

Genres Distribution (Top 7)



It's expected to have lowest 7 since they are “Drama” and

these categories as the not as popular as “Comedy”.



D. Association measures

In this section, we will compute the Chi-Square Association measure, since all our features are categorical variables. We will declare that two features are independent if the evaluated p-value is less than 5%. (Reject the Null Hypothesis)

Action: 1, Adventure: 2, Animation: 3, Children's: 4, Comedy: 5, Crime: 6, Documentary: 7, Drama: 8, Fantasy: 9, Film - Noir: 10, Horror: 11, Musical: 12, Mystery: 13, Romantic: 14, Sci- Fi: 15, Thriller: 16, War: 17, Wester: 18

	1	2	3	4	5	6	7	8	9
1									
2	0								
3	0.007	0.02							
4	0.0002	0	0						
5	0	0	0.13	0.03					
6	0	0.11	0.02	0	0				

7	0	0	0.10	0	0	0.01			
8	0	0	0	0	0	0.73	0		
9	0.038	0	0.006	0.14	0.69	0.24	0.23	0	
10	0.019	0.32	0.77	0	0	0.002	0.4	0	0.76
11	0.001	0.94	0.007	0	0	0.02	0.5	0	0.02
12	0.001	0.1	0	0.08	0.28	0.003	0.1	0	0.72
13	0.72	0.2		0.07	0	0	0	0.02	0.3
14	0	0	0.01	0	0	0	0	0.36	0.78
15	0	0	0.99	0.4	0	0.02	0	0	0
16	0	0.4	0.009	0	0	0	0	0	0.01
17	0	0.7	0.47	0.02	0	0.006	0.12	0.004	0.5
18	0.8	0.8	0.3	0.34	0.35	0.08	0.23	0	0.5

	10	11	12	13	14	15	16	17
10								
11	0.2							
12	0.48	0.01						
13	0	0.32	0.13					
14	0.07	0	0.28	0.91				
15	0.711	0	0.04	0.69	0			
16	0	0.01	0	0	0	0		
17	0.37	0	0.72	0.075	0.57	0.9	0.01	
18	0.75	0.02	0.3	0.31	0.075	0.53	0	0

This table shows each feature and their highly associated features.

Note: If feature “Action” is associated with “Mystery”. It will not be stated that “Mystery” is associated with Action in the table.

Feature	Associated with
Action	Mystery, Western
Adventure	Crime, Film-Noir, Musical, Mystery, Romance, Thriller, Western
Animation	Comedy, Documentary, Film-Noir, Mystery, Sci-Fi, War, Western
Children’s	Film-Noir, Mystery, Sci-Fi, Western
Comedy	Fantasy, Musical, Western
Crime	Drama, Fantasy, Western
Documentary	Fantasy, Film-Noir, Musical, Mystery, War, Western
Drama	Romance
Fantasy	Film-Noir, Musical, Mystery, Romance, War, Western
Film-Noir	Horror, Musical, Sci-Fi, War, Western
Horror	Mystery
Musical	Mystery, Romance, War, Western
Mystery	Romance, Sci-Fi
Romance	War
Sci-Fi	War, Western
Thriller	-
War	Western

Initial Results

The Movie recommendation problem could be dealt with as a classification problem where the model classifies the users' preferences on a specific movie whether the user will like it or not or it can also predict the rating of the user to the specific movie in this section we will showcase different models and their performance on the Movielens dataset

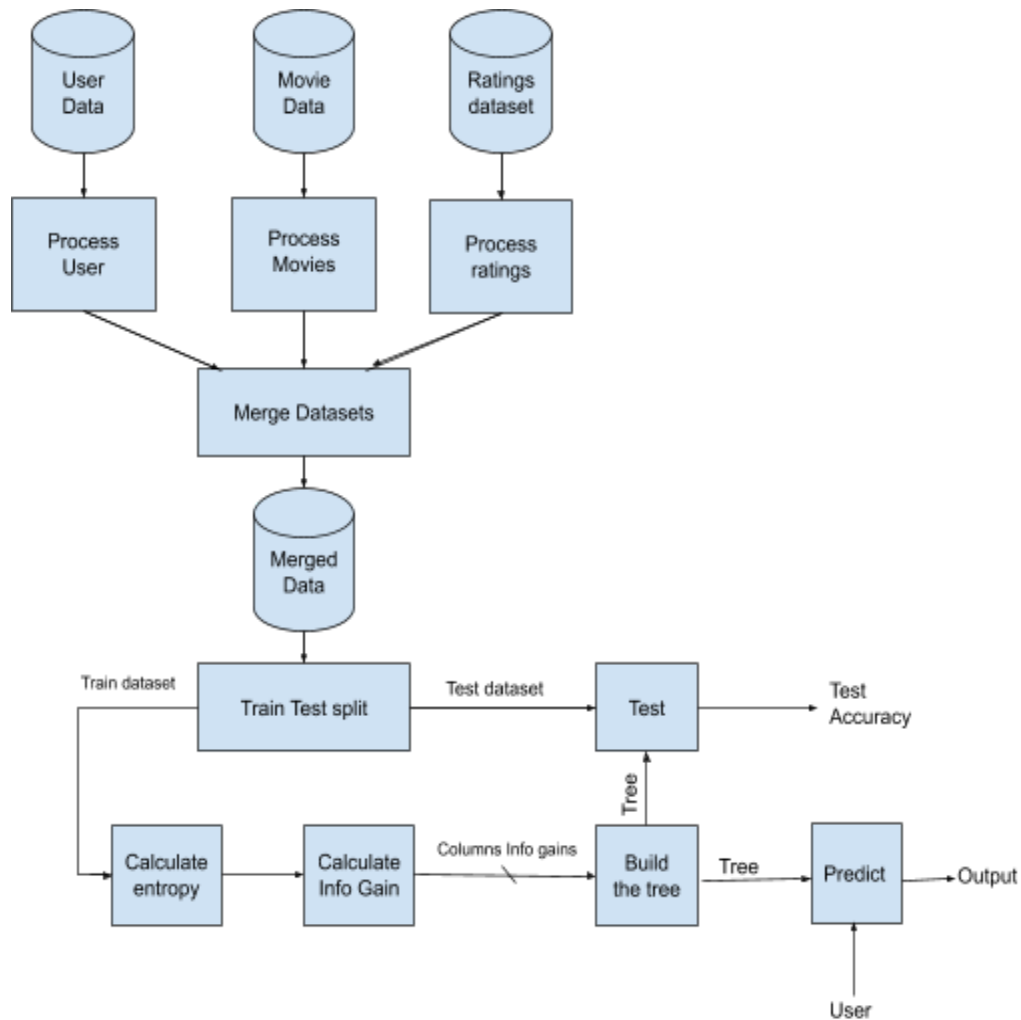
Model	Accuracy	Pros	Cons
Knn-Collaborative Filtering	35%	-It takes into consideration other users behaviour and provides a recommendation based on that	<ul style="list-style-type: none">- Very computationally expensive- Low Accuracy- Cold-Start: It doesn't work with cold-start user or items, since the dot product will be all 0s. It can't recommend anything.- Sparsity: Similarly, it doesn't work with sparse data, since the intersection between 2 users is 0, the dot product is also 0- Scalability: we need to calculate the user similarity or item similarity matrix. This is a large matrix that doesn't scale with a large number of users, movies
Logistic regression	40%	<ul style="list-style-type: none">-Not computationally expensive-Simple	<ul style="list-style-type: none">-Logistic regression requires the observations to be independent of each other-assumes linearity of independent variables and log odds which is very hard to guarantee in case of the movie lens
Neural Network	60%	<ul style="list-style-type: none">-High accuracy-Ability to capture complex functions	<ul style="list-style-type: none">- Complex- since the number of features is small there's a great tendency for overfitting the data
Decision tree	60%	Simple and computationally inexpensive (for predicting) once the tree is built	<ul style="list-style-type: none">-Doesn't perform well in unexpected events
Random forest	64%	High accuracy for the training data	The random forest algorithm tends to overfit and will have low performance for an external data

Other suggested models:

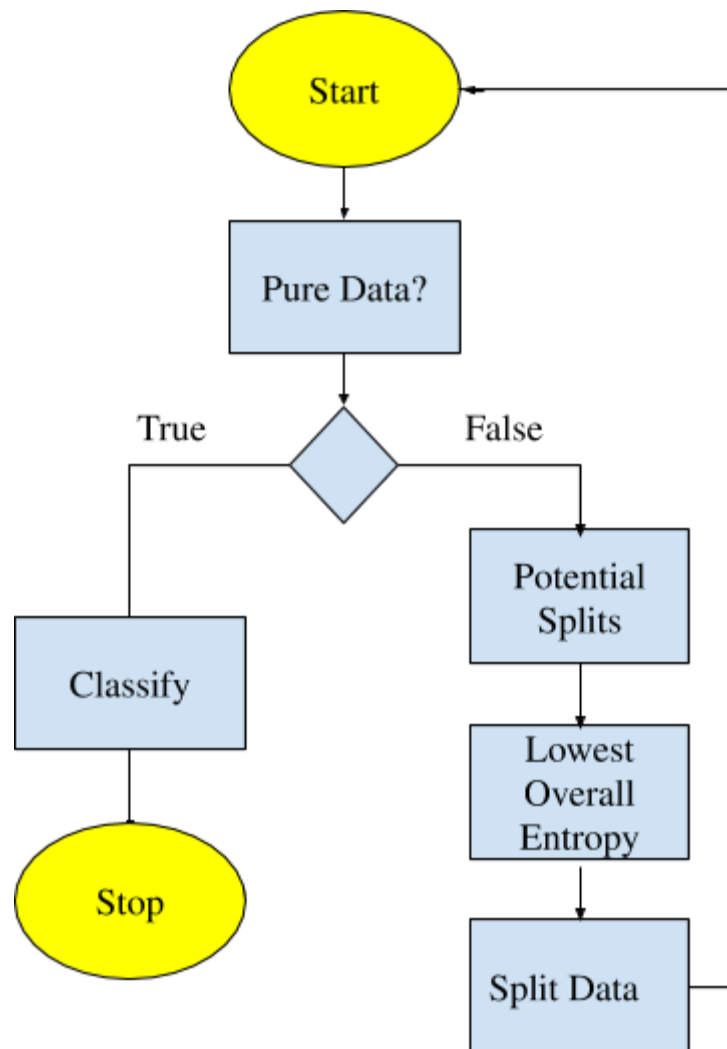
Model	Description	Accuracy (RMSE)	Pros	Cons
Mean approach	get the average rating of each movie across all users. In the testing data, if the movie exists in the training data, we get the average of its users' ratings. Otherwise, we return 3.0 (normal average rating).	0.98	Simple approach	- Affected by outliers.
Weighted Mean Approach	Compute an nxn cosine similarity matrix (n is the number of users in the training matrix). Then we use this matrix to get the weighted average of the users' ratings which will be compared to the existing movies in the testing data.	0.97	- Simple Approach -Not assuming that users are equally similar to each other.	- Not robust
Gender - Demographics	For each Movie ID we get the average rating per gender as we believe it's a significant feature from the users' info.	0.98	Using users' features not only their movie ratings.	-Not robust

Singular Value Decomposition (SVD)	<p>It's a linear algebra approach. What it does is that it decomposes a matrix into constituent arrays of feature vectors corresponding to each row and each column. It will take the training matrix and the parameter `k` which is the number of features into which each user and movie will be resolved into.</p>	0.92	<p>-It is a dimensionality reduction technique, and it avoids features dependency.</p> <p>-Robust approach.</p>	<p>-It may lead to some amount of data loss.</p>
---	---	------	---	--

Structure of the model:



Decision tree implementation:



Design Choices and tradeoffs:

1- Input Choices and tradeoffs:

We started with a feature set that has more than 60 features including:

- 1) The user data [Age, gender , occupation , state(zip code)] (one hot encoded)
- 2) The user interests scores which had 19 interest scores representing the user interests in all the different 19 movie genres
- 3) The movie data which had a 19 one hot encoded features representing the movie genre(s)
- 4) The movie production year

After research the the following improvements were made to Improve the previously mentioned data:

- 1- Removing columns with very low information gain. (less than a 0.0001)

Most of the movie interest scores had very low Information gain because the majority of the values were zeros because most of the users have only one rating and therefore it is very hard to determine their interests.

The following is a table with the features and their corresponding information gains

Feature	Information Gain
Age	0.00138
Gender	0.000234
Year of production	0.01313
Genre_1	0.0137
Genre_2	0.006
Occupation	0.001498
State	0.00221

2- Trying label encoding instead of one hot encoding for the most important features (such as the top two movie genres representing the movie)

This solves the sparsity problem of the one hot encoding method however it has another potential problem because it assumes an order or a rank but this problem is not a major problem because the decision tree algorithm doesn't depend on the numerical values of those labels

3- Converting numerical values such as age and movie production year to bins where every bin represents a range of values

Instead of having more than 60 possible values for the age we only have 4 which means less tree branches and better generalizability :

Child: representing the age from (1:12) years old

Teen: representing the age from (13:19) years old

Adult: representing the age from (20:45) years old

Old: representing the age from (45:69) years old

Instead of having 80 values for the movies production year we divided it into 10 bins representing the decade:

1920's: including years between [1910-1920]

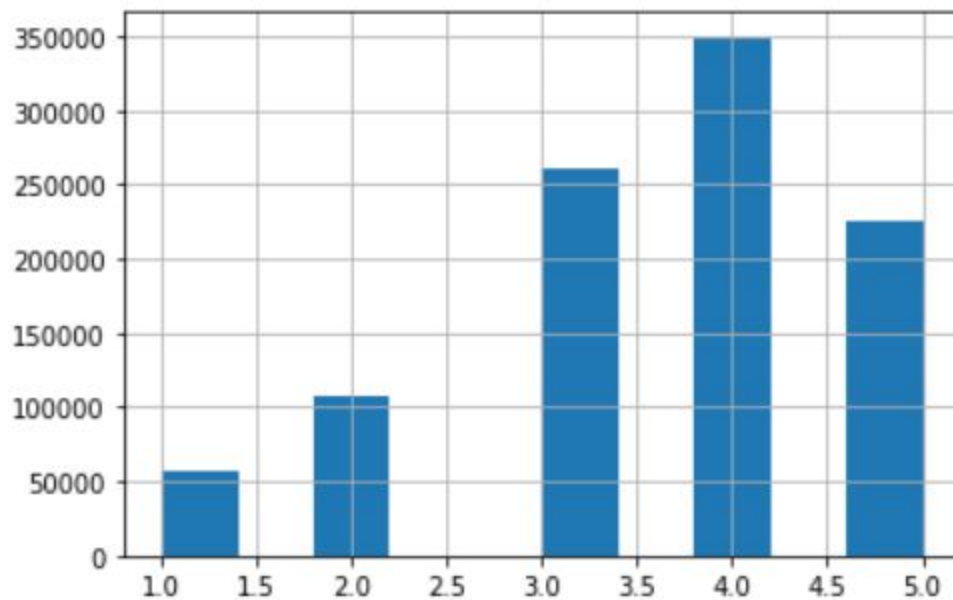
1930's: including years between [1921-1930]

And so on

Output Choices and tradeoffs:

For the output there are different design choices: the first design choice is to predict the movie rating and return a value of [1, 2, 3,4 or 5]

The problem with this design choice is the distribution of the dataset shown below.



This distribution shows that the percentage of rating= 1 is around 5%

And the percentage of instances where rating = 2 is around 10%

Whereas the percentage of instances where rating = 4 is around 35%

And the percentage of instances where rating = 3 is around 27%

And the percentage of instances where rating = 5 is around 22%

This unbalanced ratings problem could be solved by selecting only 5% randomly from each rating however this will decrease the number of ratings significantly and the effectiveness of the model

The other option was to split the ratings into two values (Likes and dislike) where like is represented by a 1 and dislike is represented by a zero.

The choice of the threshold is really tricky because the ratings are discrete values so there are no ratings with value 2.5 and therefore we have to split the ratings at either the value 2 or 3.

If we chose the likes to be [1,2] this will re-create the unbalanced data problem since the dislikes will be only 15% of the data. However, if we chose the dislikes to be at [1, 2, 3] and the likes to be at [4,5] this will create a perfectly balanced data but the accuracy and recall drop significantly as shown in the table

Likes range	Dislikes = [1,2]	Dislikes = [1,2, 3]
TP	15637	7208
FP	3062	5321
TN	1081	3161
FN	219	4309

Accuracy	83.5%	51%
Recall	83.6%	57%

So the output of the model is equal to zero (dislike) if the rating is equivalent to 1 or 2 and the output is one (Like) otherwise.

After implementing the previous adjustments we ended up with 7 features that represent most of the important 60 features and we ended with a training time equal to 1 hour and 20 minutes to build the ID3 decision tree based on 666564 (~ 600K) training instances.

Phase 5:

User Input

The user is expected to feed the application with the following inputs:

- 1- Gender
- 2- Age
- 3- Movie production year
- 4- Occupation (Selected)
- 5- ZipCode
- 6- First Movie Genre (Selected)
- 7- Second Movie Genre if exists (Selected)

** A better design choice is to let the users enter the movie title and retrieve the movie genres from the table however this choice given this dataset is not applicable because the movie titles are full of noise such as repeated characters, incorrect spelling and different non alpha characters that are very hard to determine if they belong to the movie name for example Toy Story 3 has the number '3' which distinguishes it from other toy story movie versions and other movies are written in this format (JOK3R) where the 3 represents the letter 'E' and cleaning the titles requires human analysis which is hard to implement given the time constraint and the size of the data. So the prediction is whether the user will like the movies based on the movie genres and the user data.

However this design choice allows the user to get the prediction for movies that don't even exist in the movie dataset.

User input snippet:

(1)

```
Please enter your Gender F for female and M for male : F
Please Enter your age : 1
Please enter the production year : 1990
Choose the movie genre(s) :
Choose 0 for Action
Choose 1 for Adventure
Choose 2 for Animation
Choose 3 for Children's
Choose 4 for Comedy
Choose 5 for Crime
Choose 6 for Documentary
Choose 7 for Drama
Choose 8 for Fantasy
Choose 9 for Film-Noir
Choose 10 for Horror
Choose 11 for Musical
Choose 12 for Mystery
Choose 13 for Romance
Choose 14 for Sci-Fi
Choose 15 for Thriller
Choose 16 for War
Choose 17 for Western
First Genre : 1
Second Genre : 3
Please Enter your ZipCode : 48067
```

(2)

```
Please select your occupation
0: other or not specified
1: academic/educator
2: artist
3:
4: clerical/admin
5: college/grad student
6: customer service
7: doctor/health care
8: executive/managerial
9: farmer
10: homemaker
11: K-12 student
12: lawyer
13: programmer
14: retired
15: sales/marketing
16: scientist
17: self-employed
18: technician/engineer
19: tradesman/craftsman
20: unemployed
21: writer
22:
23:
24:
25:
26:
27:
28:
29:
30:
31:
32:
33:
34:
35:
36:
37:
38:
39:
40:
41:
42:
43:
44:
45:
46:
47:
48:
49:
50:
51:
52:
53:
54:
55:
56:
57:
58:
59:
60:
61:
62:
63:
64:
65:
66:
67:
68:
69:
70:
71:
72:
73:
74:
75:
76:
77:
78:
79:
80:
81:
82:
83:
84:
85:
86:
87:
88:
89:
90:
91:
92:
93:
94:
95:
96:
97:
98:
99:
100:
101:
102:
103:
104:
105:
106:
107:
108:
109:
110:
111:
112:
113:
114:
115:
116:
117:
118:
119:
120:
121:
122:
123:
124:
125:
126:
127:
128:
129:
130:
131:
132:
133:
134:
135:
136:
137:
138:
139:
140:
141:
142:
143:
144:
145:
146:
147:
148:
149:
150:
151:
152:
153:
154:
155:
156:
157:
158:
159:
160:
161:
162:
163:
164:
165:
166:
167:
168:
169:
170:
171:
172:
173:
174:
175:
176:
177:
178:
179:
180:
181:
182:
183:
184:
185:
186:
187:
188:
189:
190:
191:
192:
193:
194:
195:
196:
197:
198:
199:
200:
201:
202:
203:
204:
205:
206:
207:
208:
209:
210:
211:
212:
213:
214:
215:
216:
217:
218:
219:
220:
221:
222:
223:
224:
225:
226:
227:
228:
229:
230:
231:
232:
233:
234:
235:
236:
237:
238:
239:
240:
241:
242:
243:
244:
245:
246:
247:
248:
249:
250:
251:
252:
253:
254:
255:
256:
257:
258:
259:
260:
261:
262:
263:
264:
265:
266:
267:
268:
269:
270:
271:
272:
273:
274:
275:
276:
277:
278:
279:
280:
281:
282:
283:
284:
285:
286:
287:
288:
289:
290:
291:
292:
293:
294:
295:
296:
297:
298:
299:
300:
301:
302:
303:
304:
305:
306:
307:
308:
309:
310:
311:
312:
313:
314:
315:
316:
317:
318:
319:
320:
321:
322:
323:
324:
325:
326:
327:
328:
329:
330:
331:
332:
333:
334:
335:
336:
337:
338:
339:
340:
341:
342:
343:
344:
345:
346:
347:
348:
349:
350:
351:
352:
353:
354:
355:
356:
357:
358:
359:
360:
361:
362:
363:
364:
365:
366:
367:
368:
369:
370:
371:
372:
373:
374:
375:
376:
377:
378:
379:
380:
381:
382:
383:
384:
385:
386:
387:
388:
389:
390:
391:
392:
393:
394:
395:
396:
397:
398:
399:
400:
401:
402:
403:
404:
405:
406:
407:
408:
409:
410:
411:
412:
413:
414:
415:
416:
417:
418:
419:
420:
421:
422:
423:
424:
425:
426:
427:
428:
429:
430:
431:
432:
433:
434:
435:
436:
437:
438:
439:
440:
441:
442:
443:
444:
445:
446:
447:
448:
449:
450:
451:
452:
453:
454:
455:
456:
457:
458:
459:
460:
461:
462:
463:
464:
465:
466:
467:
468:
469:
470:
471:
472:
473:
474:
475:
476:
477:
478:
479:
480:
481:
482:
483:
484:
485:
486:
487:
488:
489:
490:
491:
492:
493:
494:
495:
496:
497:
498:
499:
500:
501:
502:
503:
504:
505:
506:
507:
508:
509:
510:
511:
512:
513:
514:
515:
516:
517:
518:
519:
520:
521:
522:
523:
524:
525:
526:
527:
528:
529:
530:
531:
532:
533:
534:
535:
536:
537:
538:
539:
540:
541:
542:
543:
544:
545:
546:
547:
548:
549:
550:
551:
552:
553:
554:
555:
556:
557:
558:
559:
560:
561:
562:
563:
564:
565:
566:
567:
568:
569:
570:
571:
572:
573:
574:
575:
576:
577:
578:
579:
580:
581:
582:
583:
584:
585:
586:
587:
588:
589:
590:
591:
592:
593:
594:
595:
596:
597:
598:
599:
600:
601:
602:
603:
604:
605:
606:
607:
608:
609:
610:
611:
612:
613:
614:
615:
616:
617:
618:
619:
620:
621:
622:
623:
624:
625:
626:
627:
628:
629:
630:
631:
632:
633:
634:
635:
636:
637:
638:
639:
640:
641:
642:
643:
644:
645:
646:
647:
648:
649:
650:
651:
652:
653:
654:
655:
656:
657:
658:
659:
660:
661:
662:
663:
664:
665:
666:
667:
668:
669:
670:
671:
672:
673:
674:
675:
676:
677:
678:
679:
680:
681:
682:
683:
684:
685:
686:
687:
688:
689:
690:
691:
692:
693:
694:
695:
696:
697:
698:
699:
700:
701:
702:
703:
704:
705:
706:
707:
708:
709:
710:
711:
712:
713:
714:
715:
716:
717:
718:
719:
720:
721:
722:
723:
724:
725:
726:
727:
728:
729:
730:
731:
732:
733:
734:
735:
736:
737:
738:
739:
740:
741:
742:
743:
744:
745:
746:
747:
748:
749:
750:
751:
752:
753:
754:
755:
756:
757:
758:
759:
760:
761:
762:
763:
764:
765:
766:
767:
768:
769:
770:
771:
772:
773:
774:
775:
776:
777:
778:
779:
780:
781:
782:
783:
784:
785:
786:
787:
788:
789:
790:
791:
792:
793:
794:
795:
796:
797:
798:
799:
800:
801:
802:
803:
804:
805:
806:
807:
808:
809:
810:
811:
812:
813:
814:
815:
816:
817:
818:
819:
820:
821:
822:
823:
824:
825:
826:
827:
828:
829:
830:
831:
832:
833:
834:
835:
836:
837:
838:
839:
840:
841:
842:
843:
844:
845:
846:
847:
848:
849:
850:
851:
852:
853:
854:
855:
856:
857:
858:
859:
860:
861:
862:
863:
864:
865:
866:
867:
868:
869:
870:
871:
872:
873:
874:
875:
876:
877:
878:
879:
880:
881:
882:
883:
884:
885:
886:
887:
888:
889:
890:
891:
892:
893:
894:
895:
896:
897:
898:
899:
900:
901:
902:
903:
904:
905:
906:
907:
908:
909:
910:
911:
912:
913:
914:
915:
916:
917:
918:
919:
920:
921:
922:
923:
924:
925:
926:
927:
928:
929:
930:
931:
932:
933:
934:
935:
936:
937:
938:
939:
940:
941:
942:
943:
944:
945:
946:
947:
948:
949:
950:
951:
952:
953:
954:
955:
956:
957:
958:
959:
960:
961:
962:
963:
964:
965:
966:
967:
968:
969:
970:
971:
972:
973:
974:
975:
976:
977:
978:
979:
980:
981:
982:
983:
984:
985:
986:
987:
988:
989:
990:
991:
992:
993:
994:
995:
996:
997:
998:
999:
1000:
1001:
1002:
1003:
1004:
1005:
1006:
1007:
1008:
1009:
1010:
1011:
1012:
1013:
1014:
1015:
1016:
1017:
1018:
1019:
1020:
1021:
1022:
1023:
1024:
1025:
1026:
1027:
1028:
1029:
1030:
1031:
1032:
1033:
1034:
1035:
1036:
1037:
1038:
1039:
1040:
1041:
1042:
1043:
1044:
1045:
1046:
1047:
1048:
1049:
1050:
1051:
1052:
1053:
1054:
1055:
1056:
1057:
1058:
1059:
1060:
1061:
1062:
1063:
1064:
1065:
1066:
1067:
1068:
1069:
1070:
1071:
1072:
1073:
1074:
1075:
1076:
1077:
1078:
1079:
1080:
1081:
1082:
1083:
1084:
1085:
1086:
1087:
1088:
1089:
1090:
1091:
1092:
1093:
1094:
1095:
1096:
1097:
1098:
1099:
1100:
1101:
1102:
1103:
1104:
1105:
1106:
1107:
1108:
1109:
1110:
1111:
1112:
1113:
1114:
1115:
1116:
1117:
1118:
1119:
1120:
1121:
1122:
1123:
1124:
1125:
1126:
1127:
1128:
1129:
1130:
1131:
1132:
1133:
1134:
1135:
1136:
1137:
1138:
1139:
1140:
1141:
1142:
1143:
1144:
1145:
1146:
1147:
1148:
1149:
1150:
1151:
1152:
1153:
1154:
1155:
1156:
1157:
1158:
1159:
1160:
1161:
1162:
1163:
1164:
1165:
1166:
1167:
1168:
1169:
1170:
1171:
1172:
1173:
1174:
1175:
1176:
1177:
1178:
1179:
1180:
1181:
1182:
1183:
1184:
1185:
1186:
1187:
1188:
1189:
1190:
1191:
1192:
1193:
1194:
1195:
1196:
1197:
1198:
1199:
1200:
1201:
1202:
1203:
1204:
1205:
1206:
1207:
1208:
1209:
1210:
1211:
1212:
1213:
1214:
1215:
1216:
1217:
1218:
1219:
1220:
1221:
1222:
1223:
1224:
1225:
1226:
1227:
1228:
1229:
1230:
1231:
1232:
1233:
1234:
1235:
1236:
1237:
1238:
1239:
1240:
1241:
1242:
1243:
1244:
1245:
1246:
1247:
1248:
1249:
1250:
1251:
1252:
1253:
1254:
1255:
1256:
1257:
1258:
1259:
1260:
1261:
1262:
1263:
1264:
1265:
1266:
1267:
1268:
1269:
1270:
1271:
1272:
1273:
1274:
1275:
1276:
1277:
1278:
1279:
1280:
1281:
1282:
1283:
1284:
1285:
1286:
1287:
1288:
1289:
1290:
1291:
1292:
1293:
1294:
1295:
1296:
1297:
1298:
1299:
1300:
1301:
1302:
1303:
1304:
1305:
1306:
1307:
1308:
1309:
1310:
1311:
1312:
1313:
1314:
1315:
1316:
1317:
1318:
1319:
1320:
1321:
1322:
1323:
1324:
1325:
1326:
1327:
1328:
1329:
1330:
1331:
1332:
1333:
1334:
1335:
1336:
1337:
1338:
1339:
1340:
1341:
1342:
1343:
1344:
1345:
1346:
1347:
1348:
1349:
1350:
1351:
1352:
1353:
1354:
1355:
1356:
1357:
1358:
1359:
1360:
1361:
1362:
1363:
1364:
1365:
1366:
1367:
1368:
1369:
1370:
1371:
1372:
1373:
1374:
1375:
1376:
1377:
1378:
1379:
1380:
1381:
1382:
1383:
1384:
1385:
1386:
1387:
1388:
1389:
1390:
1391:
1392:
1393:
1394:
1395:
1396:
1397:
1398:
1399:
1400:
1401:
1402:
1403:
1404:
1405:
1406:
1407:
1408:
1409:
1410:
1411:
1412:
1413:
1414:
1415:
1416:
1417:
1418:
1419:
1420:
1421:
1422:
1423:
1424:
1425:
1426:
1427:
1428:
1429:
1430:
1431:
1432:
1433:
1434:
1435:
1436:
1437:
1438:
1439:
1440:
1441:
1442:
1443:
1444:
1445:
1446:
1447:
1448:
1449:
1450:
1451:
1452:
1453:
1454:
1455:
1456:
1457:
1458:
1459:
1460:
1461:
1462:
1463:
1464:
1465:
1466:
1467:
1468:
1469:
1470:
1471:
1472:
1473:
1474:
1475:
1476:
1477:
1478:
1479:
1480:
1481:
1482:
1483:
1484:
1485:
1486:
1487:
1488:
1489:
1490:
1491:
1492:
1493:
1494:
1495:
1496:
1497:
1498:
1499:
1500:
1501:
1502:
1503:
1504:
1505:
1506:
1507:
1508:
1509:
1510:
1511:
1512:
1513:
1514:
1515:
1516:
1517:
1518:
1519:
1520:
1521:
1522:
1523:
1524:
1525:
1526:
1527:
1528:
1529:
1530:
1531:
1532:
1533:
1534:
1535:
1536:
1537:
1538:
1539:
1540:
1541:
1542:
1543:
1544:
1545:
1546:
1547:
1548:
1549:
1550:
1551:
1552:
1553:
1554:
1555:
1556:
1557:
1558:
1559:
1560:
1561:
1562:
1563:
1564:
1565:
1566:
1567:
1568:
1569:
1570:
1571:
1572:
1573:
1574:
1575:
1576:
1577:
1578:
1579:
1580:
1581:
1582:
1583:
1584:
1585:
1586:
1587:
1588:
1589:
1590:
1591:
1592:
1593:
1594:
1595:
1596:
1597:
1598:
1599:
1600:
1601:
1602:
1603:
1604:
1605:
1606:
1607:
1608:
1609:
1610:
1611:
1612:
1613:
1614:
1615:
1616:
1617:
1618:
1619:
1620:
1621:
1622:
1623:
1624:
1625:
1626:
1627:
1628:
1629:
1630:
1631:
1632:
1633:
1634:
1635:
1636:
1637:
1638:
1639:
1640:
1641:
1642:
1643:
1644:
1645:
1646:
1647:
1648:
1649:
1650:
1651:
1652:
1653:
1654:
1655:
1656:
1657:
1658:
1659:
1660:
1661:
1662:
1663:
1664:
1665:
1666:
1667:
1668:
1669:
1670:
1671:
1672:
1673:
1674:
1675:
1676:
1677:
1678:
1679:
1680:
1681:
1682:
1683:
1684:
1685:
1686:
1687:
1688:
1689:
1690:
1691:
1692:
1693:
1694:
1695:
1696:
1697:
1698:
1699:
1700:
1701:
1702:
1703:
1704:
1705:
1706:
1707:
1708:
1709:
1710:
1711:
1712:
1713:
1714:
1715:
1716:
1717:
1718:
1719:
1720:
1721:
1722:
1723:
1724:
1725:
1726:
1727:
1728:
1729:
1730:
1731:
1732:
1733:
1734:
1735:
1736:
1737:
1738:
1739:
1740:
1741:
1742:
1743:
1744:
1745:
1746:
1747:
1748:
1749:
1750:
1751:
1752:
1753:
1754:
1755:
1756:
1757:
1758:
1759:
1760:
1761:
1762:
1763:
1764:
1765:
1766:
1767:
1768:
1769:
1770:
1771:
1772:
1773:
1774:
1775:
1776:
1777:
1778:
1779:
1780:
1781:
1782:
1783:
1784:
1785:
1786:
1787:
1788:
1789:
1790:
1791:
1792:
1793:
1794:
1795:
1796:
1797:
1798:
1799:
1800:
1801:
1802:
1803:
1804:
1805:
1806:
1807:
1808:
1809:
1810:
1811:
1812:
1813:
1814:
1815:
1816:
1817:
1818:
1819:
1820:
1821:
1822:
1823:
1824:
1825:
1826:
1827:
1828:
1829:
1830:
1831:
1832:
1833:
1834:
1835:
1836:
1837:
1838:
1839:
1840:
1841:
1842:
1843:
1844:
1845:
1846:
1847:
1848:
1849:
1850:
1851:
1852:
1853:
1854:
1855:
1856:
1857:
1858:
1859:
1860:
1861:
1862:
1863:
1864:
1865:
1866:
1867:
1868:
1869:
1870:
1871:
1872:
1873:
1874:
1875:
1876:
1877:
1878:
1879:
1880:
1881:
1882:
1883:
1884:
1885:
1886:
1887:
1888:
1889:
1890:
1891:
1892:
1893:
1894:
1895:
1896:
1897:
1898:
1899:
1900:
1901:
1902:
1903:
1904:
1905:
1906:
1907:
1908:
1909:
1910:
1911:
1912:
1913:
1914:
1915:
1916:
1917:
1918:
1919:
1920:
1921:
1922:
1923:
1924:
1925:
1926:
1927:
1928:
1929:
1930:
1931:
1932:
1933:
1934:
1935:
1936:
1937:
1938:
1939:
1940:
1941:
1942:
1943:
1944:
1945:
1946:
1947:
1948:
1949:
1950:
1951:
1952:
1953:
1954:
1955:
1956:
1957:
1958:
1959:
1960:
1961:
1962:
1963:
1964:
1965:
1966:
1967:
1968:
1969:
1970:
1971:
1972:
1973:
1974:
1975:
1976:
1977:
1978:
1979:
1980:
1981:
1982:
1983:
1984:
1985:
1986:
1987:
1988:
1989:
1990:
1991:
1992:
1993:
1994:
1995:
1996:
1997:
1998:
1999:
2000:
2001:
2002:
2003:
2004:
2005:
2006:
2007:
2008:
2009:
2010:
2011:
2012:
2013:
2014:
2015:
2016:
2017:
2018:
2019:
2020:
2021:
2022:
2023:
2024:
2025:
2026:
2027:
2028:
2029:
2030:
2031:
2032:
2033:
2034:
2035:
2036:
2037:
2038:
2039:
2040:
2041:
2042:
2043:
2044:
2045:
2046:
2047:
2048:
2049:
2050:
2051:
2052:
2053:
2054:
2055:
2056:
2057:
2058:
2059:
2060:
2061:
2062:
2063:
2064:
2065:
2066:
2067:
2068:
2069:
2070:
2071:
2072:
2073:
2074:
2075:
2076:
2077:
2078:
2079:
2080:
2081:
2082:
2083:
2084:
2085:
2086:
2087:
2088:
2089:
2090:
2091:
2092:
2093:
2094:
2095:
2096:
2097:
2098:
2099:
2100:
2101:
2
```

The output:

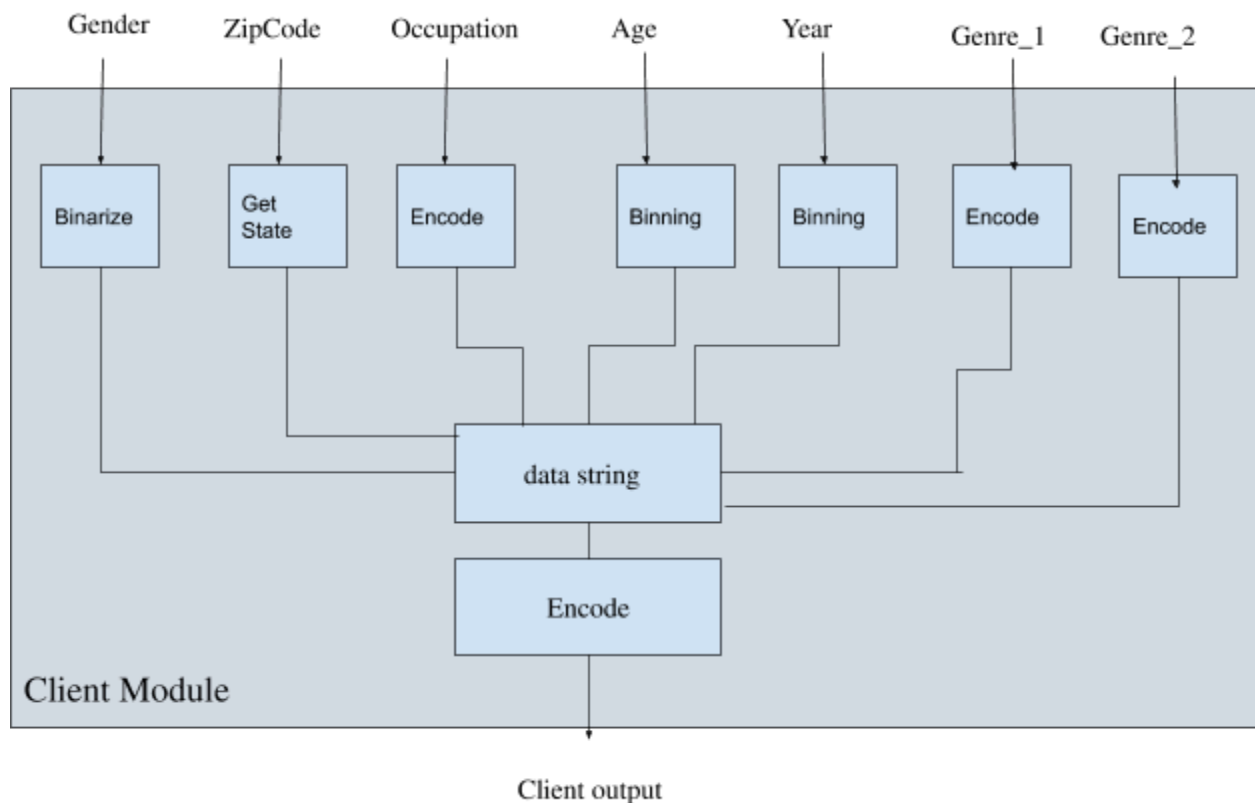
The output either suggests that the user watches the given movie or suggest that the user tries another one based on the value of the prediction (0 or 1)

Go Ahead you will like this movie!

Client :

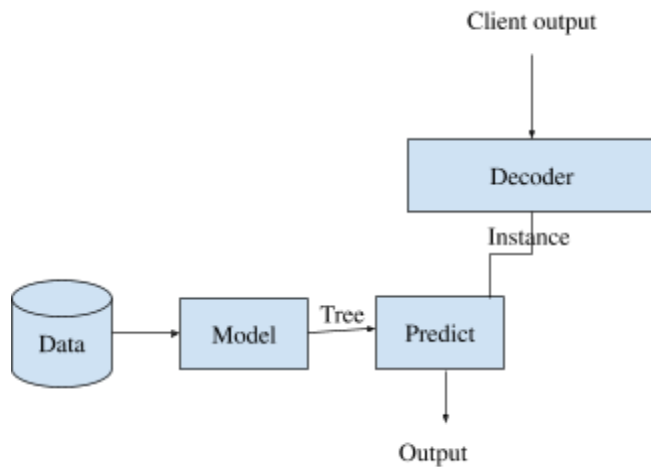
Client Design:

The client takes the input and processes it and provides in to the server in the form of an encoded string and the server takes the string , decodes it and feed it to the model so and provides a prediction for the user to either watch the movie or not.



- 1- Binarize: A function responsible for taking a categorical variable and returning 0 or 1 based on a threshold fed to the function (0 if value < threshold) and (1 if value > threshold)
- 2- Get state: function that uses the [Uszipcode](#) package to return the state.
- 3- Encode: Takes a value and returns a code to that value based on label encoding
- 4- Binning: takes a value and returns the bin that this value belongs to based on bin ranges and the value

Server:
Server Design:



The Model is described above and the output is either to watch or not to watch the movie.

Code:

Client Server Application:

<https://github.com/mjaafar97/Client-Server-implementation-of-ID3>

Model Design Code:

<https://github.com/mjaafar97/MoviesRecommendationSystem>

Demo:

[MovieRecommenderDemo](#)

Sources: