

Gaussian Process Regression for Movie Recommendations

Group 7

National University of Singapore

Abstract

In recent years, online streaming and video on-demand services have grown exponentially. The online nature of such services necessarily multiplies the number of products that can be shown to the consumer. Consequently, there is a need for such services to narrow down a selection of products which should be recommended. In order to produce accurate predictions, some factors that can be considered are: users' preferences, item profiles and correlations between items and users. In this paper, we attempt to apply Gaussian Process (GP) regression on movie ratings to generate recommendations.

1 Introduction

Today, it is easy to obtain online information pertaining to a movie and its critics through a search on Google, IMDB, Rotten Tomatoes, etc. Hence, modern day movie watchers would often research online about a movie before buying tickets in their local cinema or paying for online streams on movie streaming platforms such as Netflix, Amazon Video, or HBO.

With the rise of on-demand video services, an accurate movie recommendation system will aid the user in navigating the wide variety of entertainment options. As movie-watching shifts into the online world, consumers are faced with more choices as compared to the traditional way of browsing limited movies available in physical dvd rental stores or the theatre. Within the online ecosystem, services like Netflix are faced with the double-edged sword of being able to provide virtually any movie in existence to the consumer. Letting the consumer conduct an exhaustive search on this massive search space is certainly a computational burden that we would like to avoid.

2 Motivating Application

Moreover, motivating reasons such as word of the mouth, reviews and advertising are insufficient in the context of digitized movies. To cater to the consumer who faces the paradox of choice, there is a need for quality recommender systems which can enable movie watchers to discover new movies to their liking.

Our application firstly exploits the predictive mean and uncertainty provided by the Gaussian process model. In particular, it is important for our application to obtain a ranking of movies that we should recommend. The Gaussian process model is particularly well-suited in this context as compared to traditional regression models. This is because our selection can be made on the basis of the predictive uncertainty that the Gaussian process model usefully provides. For example, we might obtain similarly high ratings for 2 movies, A and B. If the predictive uncertainty given by the Gaussian process model is such that $\sigma_A^2 > \sigma_B^2$, then we are more likely to recommend movie B, or rank it higher on our recommendation list. Additionally, it is possible to use multiple models and weight our predictions inversely proportional to the variance. We can also exploit this predictive uncertainty in the following manner: Suppose a new user joins our service and gives ratings for a small set of movies. With this information, we can begin to learn that there are certain movies for which our knowledge is limited and the uncertainty given by the Gaussian process model is high. Asking the user to rate movies they have watched in this new group of movies, we can continually update our model.

Moreover, movie preference is highly personal and complex to model. The director, cast, visual effects and genres are just of some of the most prominent features of a movie. To give a trivial example, an individual might not like mainstream A-listers action movies such as Fast & Furious, Bourne or James Bond, but he may like Marvel action movies because he likes superhero comics. Even then, he might not like the newer Spiderman movies because of the main lead. He might love DC Comic Batman trilogy directed by Christopher Nolan but still be disappointed by Batman v Superman: Dawn of Justice. The usage of Gaussian process models in this context is particularly helpful in enabling us to exploit the correlations between users with similar preferences. Since users with similar tastes are likely to give similar ratings to a specific movie, the Gaussian process model enables us to exploit these correlations to recommend movies.

It is clear that modelling human behaviour and preference is complex. This implies that using a parametric model could unintentionally constrain our model. Since the Gaussian process model is a distribution over an infinite number of functions, we are better able to model the full spectrum

of highly variable human preferences by exploiting the non-parametric characteristic of Gaussian processes.

Finally, let us suppose that our Gaussian process model can generate relatively accurate recommendations. With increased user satisfaction and hence higher utilisation of the system, we are likely to obtain even more data from new ratings. This will allow us to further exploit the correlations among user preferences, as well as the Bayesian nature of Gaussian process models by updating the model with new data.

In conclusion, to tackle the problem of movie recommendation, we propose the application of Gaussian processes for rating prediction. Given a particular movie, we can compute the predicted ratings for each user. Further, given a particular user, we can compute the predicted ratings for all movies in the dataset. Of course, in the context of a recommender system, not all ratings must be used. Instead, we can simply select the top-N number of recommendations which can then be displayed to the user.

3 Technical Approach

A Gaussian process is a collection of random variables, any finite subset of which have a multivariate Gaussian distribution. It is completely specified by a mean function $\mu(\mathbf{x})$ and the covariance function $k(\mathbf{x}, \mathbf{x}')$. For a real process $f(\mathbf{x})$:

$$\begin{aligned}\mu(\mathbf{x}) &= E[f(\mathbf{x})] \\ k(\mathbf{x}, \mathbf{x}') &= E[(f(\mathbf{x}) - \mu(\mathbf{x}))(f(\mathbf{x}') - \mu(\mathbf{x}'))]\end{aligned}$$

The GP can then be denoted as:

$$f(\mathbf{x}) \sim \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

We assume that our data $\mathcal{D} = \{(x_1, y_1), \dots, (x_i, y_i)\}$ are such that y_i are noisy observations originating from a GP-distributed random function $f(\mathbf{x}_i)$ such that:

$$\begin{aligned}y_i &= f(\mathbf{x}_i) + \epsilon_i \\ \epsilon_i &\sim \mathcal{N}(0, \sigma_n^2)\end{aligned}$$

Given $\mathbf{y} = [y_1 y_2 \dots y_i]^\top$, suppose we have a new input \mathbf{x}_* for which we would like to obtain a prediction for. In other words, we would like to infer $p(f_* | \mathbf{y})$. Then, the predictive mean and variance from the GP can be given by:

$$\begin{aligned}E[f_* | \mathbf{y}] &= \mu(\mathbf{x}_*) + \mathbf{k}_*^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} (\mathbf{y} - \boldsymbol{\mu}) \\ V[f_* | \mathbf{y}] &= k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}_*\end{aligned}$$

Problem Definition

Traditionally, there are 2 approaches to implementing recommender systems:

1. Content-based systems make recommendations based on item attributes. For example, a user who tends to watch many movies of the horror genre will be recommended a movie from that genre.

2. Collaborative filtering systems analyse the similarities between users and/or items to make recommendations. Users will be recommended items that are preferred by users with similar tastes.

The task is thus rating prediction, given item and user attributes. In particular, suppose we are given a set of movies, \mathcal{M} , and a set of users, \mathcal{U} . We formulate an input matrix X based on 2 types of models, per-user and per-movie, to predict user ratings $y = [y_1, y_2, \dots, y_n]^\top$, where X is given by

$$X = \begin{bmatrix} x_1(1) & x_2(1) & \dots & x_d(1) \\ \vdots & & \ddots & \vdots \\ x_1(N) & x_2(N) & \dots & x_d(N) \end{bmatrix}$$

Model Definition

In this paper, we assume that user ratings follow a Gaussian distribution. We formulate two different types of models, each exploiting a different set of features:

1. **Per-user:** The per-user model can be viewed as follows – imagine the entire universe of movies, \mathcal{M} – the known ratings by this particular user consist of a subset of movies, \mathcal{M}_i . By utilising the item profiles of these $|\mathcal{M}_i|$ movies, the Gaussian process predicts the ratings for the unknown subset of movies, \mathcal{M}_u , where $\mathcal{M}_i \cup \mathcal{M}_u = \mathcal{M}$. The item profile consists of movie attributes such as year of release and average critic's rating.
2. **Per-movie:** The per-movie can be viewed using a similar analogy, but flipping the roles of users and items. Given the user profiles of all known user ratings for a specific movie, we predict the ratings that the remaining users are likely to give, based on the correlations between the user profiles. The user profile consists of user attributes such as age and gender.

Here is an example of how a prediction is made from these two models. Suppose Cara is a user in the system and she has provided ratings for a few movies. From these ratings provided by Cara, we train a per-user model for her. In addition we provide these ratings to their respective per-movie model so as to further train the models. One movie that Cara has not rated is Frozen and we are interested in predicting her rating for that movie. To do so, we first pass features from the movie Frozen (2013, Animation, Adventure and Comedy) into Cara's per-user GP model, which gives us a prediction for her rating. Likewise, we supply Frozen's per-movie model with Cara's features (age, gender) and obtain another prediction for the rating. Finally, we take a weighted average of the prediction to obtain the final prediction. This approach for taking the weighted average will be further explained under experimental evaluation.

The reason for building two models is so that we can exploit both personal and demographic preferences. The per-user model tells us a user's personal preference for certain movie features such as genre or year of release whereas the per-movie model tells us demographic preferences for a particular movie. Another advantage of utilising two models to

compute predictions is that we can still make reasonable predictions even when data is sparse. A new user would not provide many ratings on the system but we can give more weights to the per-movie model's prediction. Likewise for a new movie, we give more weights to the per-user model's prediction.

Choosing Model Parameters

The GP function is mainly characterized by its covariance function after normalizing the data to attain a mean of 0. The covariance function produces a covariance matrix which is utilised by the Gaussian process model for inference. Our choice of kernel for both experimental models, per user and per movie, was the radial basis function(RBF) kernel

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left(-\frac{1}{2} \sum_{i=1}^d \frac{(x_i - x'_i)^2}{\ell_i^2}\right)$$

where the hyperparameters σ^2 is the variance, ℓ is the length-scale, d is the dimension of the input vector x_i .

The choice of the rbf kernel was due to its ability to produce similar predictive ratings for two inputs with similar features based on their distance. If two users have contrasting preferences, it is highly unlikely that both will assign the same rating to their opposing's favorite movie.

Learning the hyper-parameters

Lemma: The product of two Gaussians function gives another (un-normalized) Gaussian function.

$$p(\mathbf{y}|X, \theta) = \int p(\mathbf{y}|\mathbf{f}, X, \theta)p(\mathbf{f}|X, \theta)d\mathbf{f}$$

where the prior $\mathbf{f}|X, \theta \sim N(0, K)$. The log marginal likelihood of the gaussian process can be computed by marginalizing the likelihood with respect to function \mathbf{f} and by using the lemma we obtain

$$\log p(\mathbf{y}|X, \theta) = -\frac{1}{2}\mathbf{y}^T(K + \sigma_n^2 I)^{-1}\mathbf{y} - \frac{1}{2}\log |K + \sigma_n^2 I| - \frac{n}{2}\log 2\pi$$

where the $\theta = \{\ell, \sigma_n^2, \dots\}$ are the hyper-parameters which can be tuned by Markov Chain Monte Carlo methods or through Bayesian Optimisation (Neal 1997; Sundararajan and Keerthi 2001). Learning the hyperparameters is possible as the log likelihood partial can be derived. Letting Q to be $K + \sigma_n^2 I$,

$$\frac{\partial}{\partial \theta} \log p(\mathbf{y}|X, \theta) = -\frac{1}{2}\text{tr}(Q^{-1} \frac{\partial Q}{\partial \theta_i}) + \frac{1}{2}\mathbf{y}^T Q^{-1} \frac{\partial Q}{\partial \theta_i} Q^{-1} \mathbf{y}$$

We used a python library GPy to optimize our per-user and per-movie models. GPy implements the bayesian optimization framework for maximizing the likelihood function. Through the batch gradient descent method, the maximum a posteriori(MAP) estimation gives us the maximum likelihood(Rasmussen, 1997).

This feedback to us critical information on the which are the important features in the per-user model through the Automatic Relevance Determination(ARD) represented by the lengthscale ℓ . Suppose ℓ is short for the per user model input vector dimension "genre", then this suggests that "genre" is an important feature for predicting the rating.

4 Experimental Setup

Dataset

Our primary dataset is the MovieLens¹ 1M Dataset which is a stable benchmark dataset widely used for evaluating recommender systems. The Dataset consists of three files, namely "ratings.dat", "movie.dat" and "user.dat" which contains around 1,000,000 ratings of approximately 3,900 movies made by 6,040 MovieLens users who joined MovieLens in 2000. Each rating is accompanied by a movie_id and a user_id so that we can map the rating to the respective user and movie.

Besides ratings, the dataset also provides movie attributes. The following attributes were extracted from "movie.dat":

movie_id movie identification number

title movie title

genres a list of genres that the movie belongs to

We also utilised user demographic information such as age, gender and occupation available from "user.dat":

user_id user identification number

gender user's gender

age user's age

occupation user's occupation

To supplement this dataset, we combined it with the MovieLens+IMDb/Rotten Tomatoes Dataset released by HetRec 2011. Our motivation for doing so was to obtain movie critic ratings from Rotten Tomatoes. As the project progressed, we also extracted the relevant directors and actors for feature engineering purposes. The notable features extracted from this dataset include:

AllCriticRating average critics' rating for the movie

AudienceRating average audience rating for the movie

CriticsNumRatings number of critics' ratings for the movie

AudienceNumRatings number of audience ratings for the movie

By combining all these datasets, we obtained a processed dataset consisting of user ratings augmented with user attributes and movie attributes.

¹MovieLens is a movie recommender system created by GroupLens Research.

Procedure

The pre-processed dataset contains both numerical and categorical features. We experimented on simple models based solely on numerical features as well as models that utilizes categorical features. We also experimented with different kernels such as the linear kernel and cosine kernel in addition to the RBF kernel. In order to test the performance of our models, the data was split using stratified sampling into train and test sets².

The per-movie model uses of the following numeric features: user age, user gender³.

The basic per-user model uses of the following numeric features: year of movie release, average critics and audience rating, number of critics and audience rating.

In the next section, we shall explain the different approaches taken to integrate non-numeric features into the models.

Feature Engineering

From our dataset, we were only able to obtain 3 demographic features: age, gender and occupation. However, the dataset contains a much richer set of features, including title and genres. We were also able to obtain the directors and actors by merging the original dataset with the supplementary dataset from HetRec. To utilise these non-numeric features so that we could input them into the Gaussian process, we have to first transform them into numerical features.

One approach was to use one-hot encoding over the categorical features. We attempted this on the genre feature where a column is created for each of the genres. For each movie, the value is 1 if the movie belongs to that particular genre, 0 otherwise. However, this technique does not scale very well as the input dimension increases with the number of categorical classes. Since there are 20 movie genres in total, we sought to reduce the number of dimensions for the input matrix.

To transform the remaining text data into useful numerical features and reduce the dimensions of the input, we employed Word2vec, which enables us to create word embeddings, or vector representations of words with a given corpus. Similar words are constructed such that they are close to each other within the vector space. We thus extracted the textual data from our datasets and processed them into a corpus suitable for training a Word2vec model on. The Word2vec model is then queried for vectors, $x = [x_1, \dots, x_k]^T$ representing each film. Additionally, we were able to obtain vectors representing each distinct genre present in the dataset. In training the model, it was found that a value of $k = 8$ produced reasonable results from a test of similarity between selected films and genres. The following table gives some of the cosine similarities between genre vectors most similar to each other:

Genre	Genre most similar to	Similarity
Drama	Crime	0.9132
Thriller	Mystery	0.9834
Animation	Children	0.9621
Sci-fi	Fantasy	0.8729

Another approach to reduce the number of dimensions is to select several genres as "buckets" and calculate the probability of a movie's genres given a bucket. This probability gives the intuition of how similar two genres are. For example, if a movie has the following 3 genres, adventure, animation, family, and the given bucket is comedy, we calculate the numerical feature as follows:

$$P(\text{adventure, animation, family} \mid \text{comedy}) \\ = \max \left\{ \begin{array}{l} P(\text{adventure} \mid \text{comedy}) \\ P(\text{animation} \mid \text{comedy}) \\ P(\text{family} \mid \text{comedy}) \end{array} \right\}$$

Since the conditional independence assumption does not necessarily hold, we instead take the max of the individual probabilities. The individual probabilities are calculated based on their various counts in the training set. In addition, the probabilities are smoothed using Witten-Bell smoothing to ensure non-zero probabilities.

$$P(\text{animation} \mid \text{comedy}) = \frac{C(\text{animation, comedy})}{C(\text{comedy})}$$

Choosing from the top popular genres (based on their counts), we identify 9 buckets: drama, comedy, crime, action, thriller, horror, fantasy, family and animation. This probabilistic approach allows us to compare the distance between genres and buckets. It also reduces the dimension of the input vector significantly.

5 Experimental Evaluation

Our approach involves combining the predictions from the per-user and per-movie models. We thus train a GP model for each of the 6,040 users and each of the 3,090 movies using the processed data. The number of ratings per user ranges from 18 to 2264 with a median of 94. The number of ratings per movie ranges from 1 to 3428 with a median of 135. Given that the maximum training samples per GP model is only in the thousands, there was no need to use any techniques for handling large datasets in GP model.

Evaluating the accuracy of the models is slightly more involved as we have to feed the right entry to the right model. Suppose we want to evaluate the accuracy of the per-user models using the test set. For each entry on the test set, we pass the entry to the per-user model that matches the entry's user_id and obtain the prediction for that entry. We do this for every entry on the test set and that gives us the full prediction for the test set. Likewise for per-movie models.

As mentioned, we experimented with various combinations of features and kernels. Below is a subset of the experimental results⁴ for the different combinations: The Mean

²A 70-30 split was used, while stratification was done over the user id.

³Encoded as Female: 0 Male: 1

⁴Full result is provided in the appendix

Absolute Error (MAE) and R-Squared (R^2) for the models trained are given as follows:

Type	Features	Kernel	MAE	R^2
Per-Movie	Numeric	RBF	0.7823	0.2308
		Cosine	0.7823	0.2308
		Linear	0.7823	0.2307
			0.8127	0.1663
Per-User	Numeric	RBF	0.8273	0.1424
	OneHotEncoding		0.8210	0.1544
	Word2vec Genres		0.8278	0.1424
	Word2vec Movies		0.8204	0.1534
	Probabilistic			

We observed that when predicting with just per-movie or per-user models by themselves, per-movie models outperforms per-user models. This could be due to per-movie models having more training points on average.

Combining the models

The predictions from the per-user and per-movie models can be combined to give an estimate with a lower variance. While we can take the average of both predictions as the final prediction, we can do better than that as every prediction comes with a variance which allows to calculate the optimal weights that minimises variance of the final prediction. Below we shall show how to obtain the optimal weights given the prediction variances. Suppose the per-user model returns $(\hat{y}_{user}, \sigma_{user}^2)$ and the per-movie model returns $(\hat{y}_{movie}, \sigma_{movie}^2)$. Given that $\alpha + \beta = 1$:

$$\begin{aligned}
\hat{y}_{final} &= \alpha \hat{y}_{user} + \beta \hat{y}_{movie} \\
V[\hat{y}_{final}] &= V[\alpha \hat{y}_{user} + \beta \hat{y}_{movie}] \\
&= \alpha^2 V[\hat{y}_{per_user}] + \beta^2 V[\hat{y}_{per_movie}] \\
&= \alpha^2 \sigma_{user}^2 + \beta^2 \sigma_{movie}^2 \\
&= \alpha^2 \sigma_{user}^2 + (1 - \alpha)^2 \sigma_{movie}^2 \\
&= (\sigma_{user}^2 + \sigma_{movie}^2) \alpha^2 - (2\sigma_{movie}^2) \alpha + \sigma_{movie}^2
\end{aligned}$$

To find the minimum, we differentiate $V[\hat{y}_{final}]$:

$$\frac{\partial}{\partial \alpha} V[\hat{y}_{final}] = 2(\sigma_{user}^2 + \sigma_{movie}^2) \alpha - (2\sigma_{movie}^2)$$

and equating it to zero give us the optimal value of α :

$$\begin{aligned}
0 &= 2(\sigma_{user}^2 + \sigma_{movie}^2) \alpha - (2\sigma_{movie}^2) \\
\alpha &= \sigma_{movie}^2 / (\sigma_{user}^2 + \sigma_{movie}^2)
\end{aligned}$$

Substituting the result into the equation of \hat{y}_{final} , we get the optimal method to combine the two predictions:

$$\begin{aligned}
\hat{y}_{final} &= \frac{\sigma_{movie}^2 \hat{y}_{user} + \sigma_{user}^2 \hat{y}_{movie}}{\sigma_{user}^2 + \sigma_{movie}^2} \\
V[\hat{y}_{final}] &= \frac{\sigma_{movie}^2 \sigma_{user}^2}{\sigma_{user}^2 + \sigma_{movie}^2}
\end{aligned}$$

Here are some of the validation results for the combined predictions using the RBF kernel:

Combination	Per-User	Per-Movie	MAE	R^2
Variance Avg	Numeric	Numeric	0.6648	0.2752
	OneHotEnc		0.6641	0.2691
	Word2Vec Genres		0.6635	0.2752
	Probabilistic		0.6640	0.2724
Simple Avg	Numeric	Numeric	0.6653	0.2742
	OneHotEnc		0.6647	0.2665
	Word2Vec Genres		0.6643	0.2733
	Probabilistic		0.6647	0.2706

The combined predictions are much stronger than just taking one set of predictions. Although the effect is not as dramatic as we hoped, there was a small improvement in accuracy when we used the variance weighted average over the simple average counterpart.

6 Improvements

The rbf kernel exploits the similarity through the euclidean distance. However, Wang(2008) mentions that the cosine distance could also give us important information to predict the ratings in the per-user and per-movie model.

Therefore, to improve the prediction ratings using Gaussian Process Regression, we could create a new kernel function as proposed by Aftab(2014) used in neural networks.

$$k(\mathbf{x}, \mathbf{x}') = \alpha \cos(x, x') + \beta \sigma^2 \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$$

where the hyperparameters $\theta = \{\alpha, \beta, \sigma^2, \gamma, \dots\}$ can be learn from cross validation.

Since the cosine kernel and the rbf kernel are both valid kernels, then the new kernel can be shown to be positive semi-definite and symmetric which satisfies the criteria for a valid kernel.

7 References

Aftab, W., Moinuddin, M., & Shaikh, M. S. (2014). A Novel Kernel for RBF Based Neural Networks. *Abstract and Applied Analysis*, 2014, 1-10. doi:10.1155/2014/176253

Neal, R. M. (1997, January 1). Monte Carlo implementation of Gaussian process models for Bayesian regression an. Retrieved October 13, 2016, from http://adsabs.harvard.edu/cgi-bin/bib_query?arXiv:physics/9701026

Rasmussen, & Edward, C. (1997, January 1). Evaluation of gaussian processes and other methods for non-linear regression. Retrieved October 13, 2016, from <http://dl.acm.org/citation.cfm?id=927743>

Sundararajan, S., & Keerthi, S. S. (2001). Predictive Approaches for Choosing Hyperparameters in Gaussian Processes. *Neural Computation*, 13(5), 1103-1118. doi:10.1162/08997660151134343

Wang, J., Vries, A. P., & Reinders, M. J. (2008). Unified relevance models for rating prediction in collaborative filtering. *ACM Transactions on Information Systems*, 26(3), 1-42. doi:10.1145/1361684.1361689