

Two-Factor Recommendation Models

Fall 2014

6.867 Final Project

Geoffrey Gunow and Neha Patki

1 Introduction

In 2014, Yelp published an academic dataset containing over 42,000 businesses that span 5 metropolitan areas. Types of data included the full star rating, and reviews for each business, as well as aggregated statistics for each user's reviews. In this project, we use the dataset to formulate, model, and solve the problem of estimating a user's overall review for a business, given the user's history. The algorithm we present estimates a single rating, but can be applied to multiple business to determine the best one for a given user. Thus, our problem falls in the category of recommendation systems.

1.1 Motivation

Yelp does not currently give personal recommendations to its users, but the data it collects about reviews and users make it possible to learn a user's preferences. Typical approaches fall in two categories. Collaborative filtering analyzes similarities between different users, and the items they rate highly. Content-based filtering focuses on a single user's personal history to determine individual likes and dislikes.

For Yelp's data set, we find that a collaborative filter would not take into account an individual's own preferences for cuisines and restaurant sub-categories. However, we cannot use pure content-based filtering because a single user may not have reviewed a significant number of restaurants, making the data points sparse. This motivates us to use a hybrid approach in estimation, where we combine a user's history with a cluster-based approach meant to account for the sparsity of data points. We introduce a hidden factor that models an overall group of similar users, and use this factor to separately learn preferences for each group.

An additional challenge and motivation for the problem is analyzing the given data to determine relevant features. Yelp offers a vast array of data ranging from the ambiance in each restaurant, to the number of times a user was voted funny by their peers. Furthermore, the 5 metropolitan areas are diverse in their offerings of food, and users may show preferences towards particular types of restaurants. We spend a significant time engineering feature vectors to reflect the variety of information. After the learning process, this enables us to comment on significant features to reveal trends in user preferences.

1.2 Problem Formulation

The overall goal is to accept, as inputs, a user and business, and to output a floating point value in $[1, 5]$ that predicts what the user will rate the business. We do this using two steps.

1. Categorize the users into groups who have roughly similar interests. This is the hidden factor.

2. Separately train each category to learn the average business rating for only the users in that group.

For example, assume a particular group we learn in step 1 contains users who enjoy restaurants with a hipster ambiance, which means restaurants matching this profile will achieve higher ratings by them. In step 2, the expected output is the average rating given only by its members. This means that another group may be trained to rate the same business differently, which shows its users have different broad preferences. Figure 1 shows an overall schematic of the system.

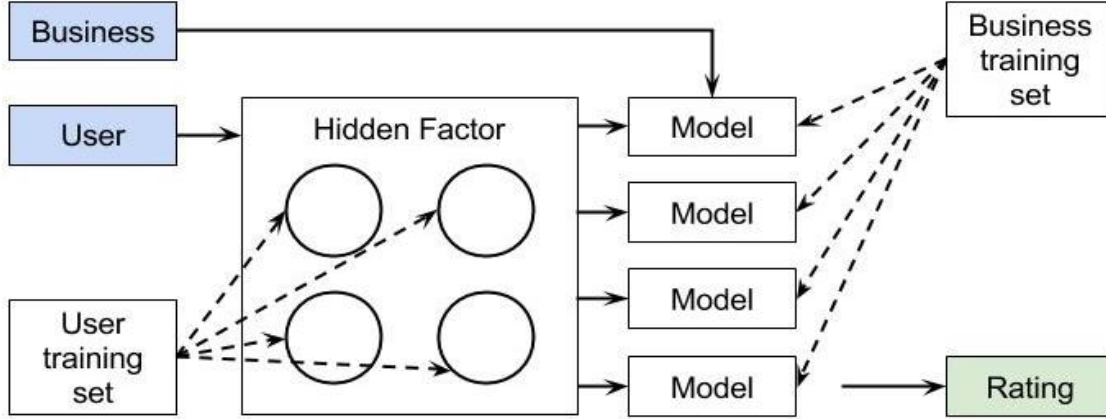


Figure 1: The architecture for the overall system. Inputs are colored blue while the output is shown in green. Dotted lines represents inputs during training, and solid lines represent the flow of information at testing time. The hidden factor is used to select a model, or a mixture of models, that the business data is tested on. The output is the rating.

Modeling the hidden factor in step 1 is an unsupervised learning problem, because we are constructing groups for which we have no prior labels. However, after learning the hidden factor, the step 2 is fully supervised because we are presented with a full list of ratings for each business. Considering the full system, we see that we can perform an end-to-end test given a user, a business, and an expected rating. This allows us to indirectly supervise the hidden factor to determine the appropriate hyperparameters for step 1. The overall complexity of our system grows, and we perform a train-validate-test strategy recursively for step 2 within the context of the overall system.

In the next sections, we discuss constructing the feature vectors that create the user and business data shown in Figure 1. We describe multiple approaches to implement both factors, and describe our train-validate-testing strategy. Finally, we compare the different methodologies in terms of accuracy and computational time.

2 Features

2.1 User Data

2.2 Business Data

3 Methodology

3.1 Hidden Factor

3.1.1 K-Means Clustering

3.1.2 K-Neighbors

3.1.3 Mixture of Gaussians

3.2 Second Factor

3.2.1 MLE

3.2.2 Lasso

3.2.3 Bayesian Ridge

3.2.4 Random Forests

4 Experiments

4.1 Setup

4.2 Results

5 References