

# Recipe Recommendations

Easton Potokar

March 2020

## Abstract

We seek to explore making a personalized recipe recommendation system based on ingredients, tags, and past user reviews. We analyze both the accuracy of the recommendations as well as the temporal complexity of making them. We find that we can make very accurate recommendations to satisfy what cravings people may have and those they didn't know they had.

## 1 Problem Statement and Motivation

Machine Learning algorithms that make recommendations to users are becoming more and more common and are seen all over the internet for movies, friends, web searches, etc. A less common application is for recipe suggestions. Is it possible to find a perfect recipe for someone given their past preferences of foods? Can we do it in fast enough so we don't lose their interest?

This has been attempted using matrix factorization [1], which essentially utilizes both components of an SVD decomposition. There has also been work done to use a hybrid approach that utilizes both the reviews of users and the attributes of the recipes, whereas more traditional methods use one or the other [2]. We seek to continue along this path of more traditional methods, while also analyzing the speed of our algorithms to see how viable they would be in production.

## 2 Data

Our data was gathered from [3], whose dataset can be found [here](#). The data consists of more than 180 thousand recipes (made up of their ingredients, calorie level, and various tags describing them), and 700 thousand user reviews of the 180 thousand recipes all scraped from Food.com and ranging from 1 to 5 stars. The original authors of the data used various natural language processing techniques to parse the recipes into about 8 thousand unique ingredients, around 500 unique tags (both characterized with a 1 when they are present, 0 when they aren't), and a calorie level that is a 0, 1, or 2 denoting low-calorie to high-calorie recipes. They did this to create a personalized recipe generation system, similar to our task here. To visualize

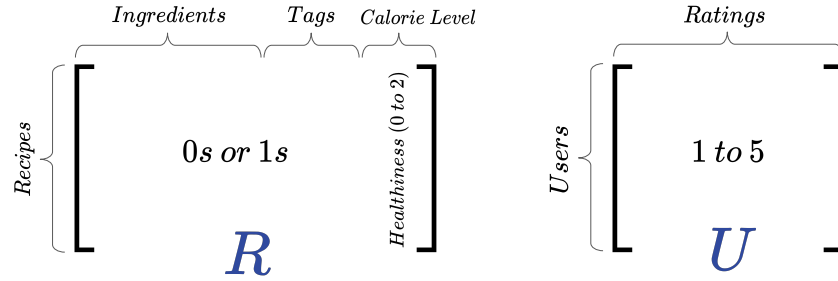


Figure 1: Data Overview

what we’re saying, see Figure 1. We denote the recipe dataset as  $R$  and the user dataset as  $U$  throughout the rest of the paper.

The data appears to be very reliable, having been scraped directly from Food.com with the scraper being publicly available on github [4]. Furthermore, it was used in a published research article, so one would hope everything was done ethically and correctly. The dataset is plenty large enough to give use meaningful answers to our questions and should be able to make a fairly robust recommendation system.

### 3 Ethical Ramifications

As with all recommendation systems, we must be careful of what sort of bias we program into our algorithms. It would be very easy to heavily weight the algorithms to favor low-calorie recipes (healthier options) in an effort to push the public towards eating more healthy. While this would likely be well intentioned, in a way it is removing one’s freedom of options, and is manipulating users’ behavior. Many similar biases must also be carefully considered, whether it be toward a certain ingredient/brand (company sponsorship), a certain ethnic preference (potentially racist), etc.

## 4 Methods

### 4.1 Data Preparation

Our data is on the format as seen in Figure 1. This is a far from perfect solution. Note that the importance of a recipe having salt will be just as important as a recipe having jasmine or a jalapeno pepper. Commonly used in document classification, we use "TF-IDF" (term frequency–inverse document frequency) to help with this problem. In our case, we change each entry using the conversion (where item means either ingredient/tag):

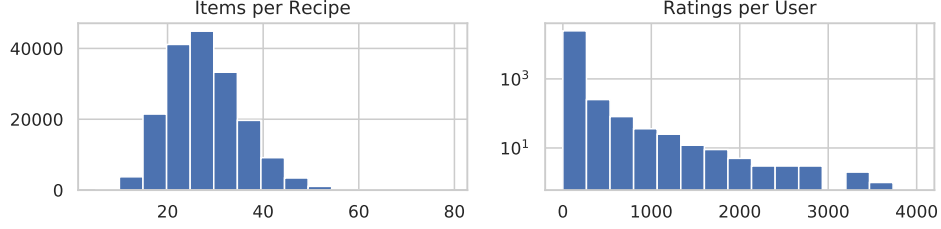


Figure 2: Data Statistics

$$\frac{\text{Quantity of } i\text{th item}}{\text{Total \# of items in Recipe}} * \log \left( \frac{\text{Total \# of Recipes}}{\text{\# of Recipes with } i\text{th item}} \right)$$

We perform the same transformation on the user dataset, replacing item with rating, and recipes with users. Throughout this document we analyze our data in both forms in order to visualize which data format helps the algorithms perform the best. We denoted these datasets as *Rhat* and *Uhat* as compared to *R* and *U* for the original datasets.

To give ourselves a little intuition, we also visualize the distributions of how many items a recipe has and how many recipes a user has rated. This can be seen in Figure 2.

## 4.2 Recommendation Process

Most recommendation systems or "filtering" systems are considered either content-based or collaborative-based filtering. Content-based systems make recommendations on attributes of the items; in our case, this would be ingredients and tags of our recipes. Collaborative-based filtering uses the history of other users who are similar to us to recommend items. We'll explore using both. Our workflows will be shown in Figure 3.

### 4.2.1 Collaborative-Based Filtering

It's here that we'll use our new algorithm, density-based spatial clustering of applications with noise (DBSCAN). It has two hyperparameters,  $\epsilon$  and  $\text{min\_samples}$ . It begins with a random point and checks if there is  $\text{min\_samples}$  within  $\epsilon$  of it. If so, all these samples are added to it's cluster. It then checks the same thing for those samples, and iteratively continues checking all points in the cluster. Once there is no more samples to check in the cluster, randomly choose another point and perform the same to create another cluster. This allows DBSCAN to detect non-convex clusters, unlike KMeans and GMM, but still be relatively faster than MinCut and t-SNE.

We'll go over what was attempted in each step of collaborative-based filtering.

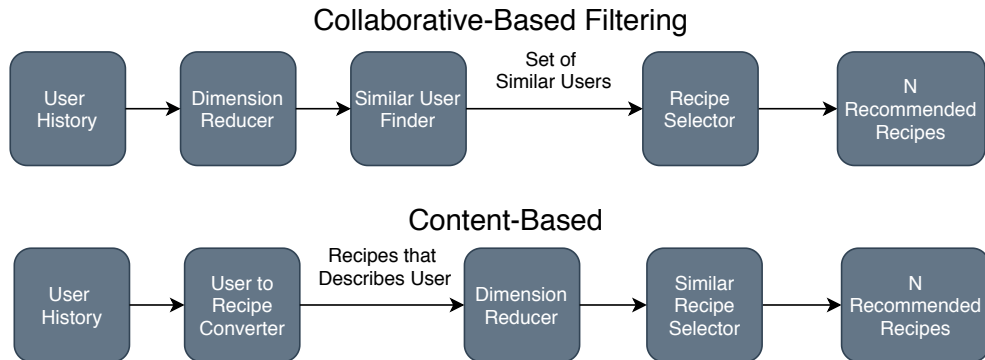


Figure 3: Recommendation Systems Flow

1. **Dimension Reducer.** PCA, LDA, NMF and Kernelized PCA were all used. Probabilistic PCA and Sparse PCA were also attempted, but don't play well with sparse input. Kernelized PCA was also too computationally expensive and often abandoned.
2. **User Finder.** We tried various algorithms for this:
  - k Nearest Neighbor (kNN).
  - Various Clustering Algorithms. Return all other users found in same cluster. Attempted using KMeans, GMM, MinCut, and DBSCAN.
  - Combination of clustering and kNN. Return k nearest neighbors found in cluster. Could potentially be more computationally effective.

We also take care to make sure our user finder doesn't return the input user.

3. **Recipe Selector.** Taking the users passed to it, count number of times recipes are rated by user set along with their rating. Choose from among recipes by:
  - Sorting by sum of ratings and return largest N sums.
  - Sorting by frequency and then by average rating. Return N largest.
  - Sorting by average rating and then frequency. Return N largest.

Note if our algorithm returns too few recommendations, (ie a cluster ends up being a singleton or something of the sort), we simply fill in the rest with random recipes.

#### 4.2.2 Content-Based Filtering

We'll go over what was attempted in each step of content-based filtering.

1. **User to Recipe Converter.** We attempt this in two ways:
  - Take all recipes with rating from user above a threshold. This method keeps each recipe separate (SR).
  - Take all ingredients from recipes with rating from user above a threshold and combine into one huge recipe. (TR)
2. **Dimension Reducer.** Dimensions were reduced in the same way as in the Collaborative-Based Approach.
3. **Similar Recipe Selector.** In the case where the set of user describing recipes is a singleton, simply take the N NNs. In the case it isn't, take k NNs of each recipe and return the N most frequent.

### 4.2.3 Scoring

To evaluate our models, we removed around 100 thousand reviews from random users. We will then take in the user as input and try to recommend the removed recipe that we know they rated or a recipe similar to it.

To quantitatively define how good a prediction is, we use two different loss functions. Let  $R$  be the recommended recipe,  $G$  the goal recipe and  $I_R, I_G$  be the set of items that each has. We define our loss functions as:

$$Int(R, G) = |I_R \cap I_G| \quad Com(R, G) = -|(I_R \cap I_G)^C| \quad (1)$$

We desire to maximize both.  $Int(R, G)$  rewards us for correctly recommending items;  $Com(R, G)$  punishes us for recommending incorrect items. For example, sometimes we may recommend something with cheddar cheese instead of American cheese and  $Com(R, G)$  punishes that, even though it's likely close enough. Similarly,  $Int(R, G)$  rewards recipes that have many ingredients. Using both gives us more insight on how our algorithms perform.

To get our total loss, we simply take the mean of all N (we have arbitrarily chosen N=5 for our tests) recommendations and then the mean of all the data points together.

## 5 Results

We ran a grid search over each combination of dimension reducer and similarity finder to find the best hyperparameters for each combination of the two. We also timed both the fitting as well as the prediction steps to be able to see which methods are actually feasible in practice. We did this separately for both the intersection and the complement metric to see how each could be maximized. The results can be seen in Figures 4, 5 and 6.

After a few test were done, we found that the best recipe selector for collaborative-based filtering was sorting by sum of ratings, almost uniformly. To help narrow down our search, we used it exclusively.

Certain combinations had to be abandoned as well due to high processing requirements that just didn't make them feasible in practice. This includes kernelized PCA with content-based filtering, as well as using MinCut along with NN.

As a baseline, we tested the performance of randomly predicting N recipes and recommending the N most popular recipes. Random recipes got 7.66 and -39.12 and popular recipes got 9.58 and -42.515, in the intersection and complement metrics respectively.

We also include the best results in each category in Figure 6. We calculated best overall by averaging their rank in each metric.

## 6 Analysis

### 6.1 Collaborative-Based Approach

First note that nearly none of these algorithms were particularly good in the complement metric. This means that on average, no particular algorithm was extremely

	Mean Intersection Metric		Mean Complement Metric		Mean Fit Time		Mean Prediction Time	
kNN	8.58	8.44	-40.31	-40.42	2121.33	1409.83	103.43	132.49
NNBall	9.39	9.58	-40.00	-42.45	331.94	605.91	5337.02	8701.91
KMeans	9.68	9.77	-41.94	-41.44	1113.96	625.85	5041.26	6460.48
GMM	9.68	9.80	-41.85	-41.18	1224.96	550.77	3320.89	3771.90
MinCut	9.72	9.62	-41.78	-42.40	1372.13	977.84	6676.89	7379.82
DBSCAN	9.65	9.59	-41.41	-42.50	2262.34	1855.48	6370.40	9905.02
KMeansNN	8.42	8.50	-40.27	-40.01	343.03	1259.54	2093.82	2616.44
GMMNN	8.48	8.47	-40.10	-39.83	2267.96	505.38	1762.95	1289.28
DBSCANN	8.52	8.42	-40.13	-40.46	736.97	847.87	2104.43	3326.59
	U	Uhat	U	Uhat	U	Uhat	U	Uhat
	Mean Intersection Metric		Mean Complement Metric		Mean Fit Time		Mean Prediction Time	
PCA	9	9.1	-40.37	-41.36	132.66	152.57	3189.71	4951.51
KPCA	9.2	9.1	-40.61	-40.94	216.45	323.12	2649.86	5444.50
NMF	9.1	9.1	-41.53	-41.37	184.75	187.44	4986.22	5338.93
LDA	9.2	9.2	-40.95	-41.08	4699.30	3176.19	3756.91	3635.70
	U	Uhat	U	Uhat	U	Uhat	U	Uhat

Figure 4: Comparison of Collaboration-Based Recommendation Algorithms

adept at the complement metric under all circumstances, but doesn't mean that a certain combination won't perform extremely well.

Also, the clustering algorithms all did almost uniformly well, with very little separating them. This surprises me since I expected the clusters to be relatively non-convex, but based on these results, I would actually assume that most of the resulting clusters are convex. They performed a little better than recommending popular recipes, and I expect would bring better variety to users as well, making them particularly worth having around. It also appears the shortcomings of kNN also fell upon the combo cluster and kNN algorithms as they performed better than random, but still worse than just recommending popular algorithms. I would also guess that NNBall essentially acted as a "local" cluster by just selecting everything within a radius of itself, thereby circumventing the shortcomings of kNN (which may reach very far away to select a neighbor).

Also, the choice of data, original or TF-IDF, made very minimal difference in the results, which is extremely surprising. This could be due to most users having rated so few recipes out of the total amount of recipes that it ends up looking like they're close/not close regardless of the data.

Note that the choice of dimension reducer had little to no effect on average either. Again, this doesn't mean a certain combination of dimension reducer and similar user selector might be dynamite, but as a whole, it didn't matter too much.

	Mean Intersection Metric		Mean Complement Metric		Mean Fit Time		Mean Prediction Time	
TR-3	11.54	9.20	-48.69	-38.52	503.79	585.98	155.12	150.23
TR-4	11.41	8.96	-48.26	-38.74	457.17	343.03	165.62	161.36
TR-5	11.13	9.04	-47.89	-38.53	472.26	353.83	155.19	162.20
SR-3	8.57	8.58	-37.92	-38.99	335.75	391.46	2683.47	2593.36
SR-4	8.56	8.64	-37.88	-38.95	314.13	304.49	2892.06	3135.03
SR-5	8.51	8.63	-37.95	-38.91	329.46	357.37	2681.70	3201.14
	R	Rhat	R	Rhat	R	Rhat	R	Rhat
	Mean Intersection Metric		Mean Complement Metric		Mean Fit Time		Mean Prediction Time	
PCA	10.26	9.59	-44.68	-39.55	8.33	9.00	1431.81	1441.70
NMF	9.96	8.44	-42.31	-38.11	192.29	98.85	1455.37	1607.03
LDA	9.64	8.50	-42.30	-38.66	1005.66	1060.24	1479.39	1652.92
	R	Rhat	R	Rhat	R	Rhat	R	Rhat

Figure 5: Comparison of Content-Based Recommendation Algorithms

## 6.2 Content-Based Approach

Right off the bat, it can be seen that mashing all previously liked ingredients of a user into a recipe and finding similar ones (TR) likely just sought out recipes that have many, many ingredients, which means it should probably be avoided. On the other hand, separating out the recipes had the opposite effect, performing extremely well in the complement metric and poorly in the intersection metric.

Once again, the choice of dimension reducers appears to mostly be a wash.

Finally, the choice of data in this scenario did play a huge role in the results. Notice that the use of TF-IDF really improved performance in the complement metric, possibly playing a pretty active role in transforming the vector space.

## 7 Conclusion

Recommendations systems are common all over the web from Netflix to eBay. One space where the usefulness of these hasn't been researched extensively is the food recipe industry. In our analysis, we found no clear-cut perfect algorithm for all scenarios. Depending on what you want your users to experience - either extremely similar recipes, or a little more freedom with some overlap in previous history - the choice of machine learning algorithm is going to play a large role.

Best Intersection Metric		Best Complement Metric		Best Overall	
(R) PCA → TR-3	12.16	(U) PCA → NNBall	-34.58	(Uhat) KPCA → KMeans	9.92 -39.15
(R) PCA → TR-4	12.08	(Rhat) NMF → TR-3	-36.68	(Rhat) PCA → TR-3	10.28 -39.50
(R) NMF → TR-3	11.69	(Rhat) NMF → TR-4	-37.17	(Rhat) PCA → TR-5	10.06 -39.68
(R) PCA → TR-5	11.65	(Rhat) NMF → TR-5	-37.19	(Rhat) PCA → TR-4	10.22 -39.71
(R) NMF → TR-4	11.61	(R) NMF → SR-5	-37.21	(U) PCA → NNBall	8.78 -34.58
(R) NMF → TR-5	11.21	(R) NMF → SR-4	-37.25	(U) LDA → KMeans	9.81 -40.26
(R) LDA → TR-3	10.77	(R) NMF → SR-3	-37.27	(U) PCA → DBSCAN	9.73 -40.22
				Int	Com
Fastest Fit		Fastest Predict			
(Rhat) PCA → SR-3	2.94	(U) PCA → kNN	71.00		
(Rhat) PCA → SR-4	3.06	(Uhat) PCA → kNN	78.99		
(Rhat) PCA → SR-5	3.30	(U) KPCA → kNN	86.32		
(Rhat) PCA → TR-5	7.33	(Uhat) NMF → kNN	86.77		
(Rhat) PCA → TR-3	8.48	(U) NMF → kNN	91.40		
(R) PCA → TR-5	8.56	(R) NMF → TR-5	105.23		
(R) PCA → TR-4	8.92	(R) NMF → TR-3	110.74		

Figure 6: Best Results



## References

- [1] Sai Bharath Goda. Recommender system for recipes. Master's thesis, Kansas State University, 5 2014.
- [2] Anirudh Jagithyala. Recommending recipes based on ingredients and user reviews. Master's thesis, Kansas State University, 5 2014.
- [3] Bodhisattwa Prasad Majumder, Shuyang Li, Jianmo Ni, and Julian McAuley. Generating personalized recipes from historical user preferences, 2019.
- [4] Bodhisattwa Prasad Majumder, Shuyang Li, Jianmo Ni, and Julian McAuley. Recipe personalization. <https://github.com/majumderb/recipe-personalization>, 2019.