

Manuscript Title

This manuscript ([permalink](#)) was automatically generated from [Akshay163/gp13 food flow@5fcea6d](#) on December 6, 2020.

Authors

- **John Doe**

 [XXXX-XXXX-XXXX-XXXX](#) ·  [johndoe](#) ·  [johndoe](#)

Department of Something, University of Whatever · Funded by Grant XXXXXXXX

- **Jane Roe**

 [XXXX-XXXX-XXXX-XXXX](#) ·  [janeroe](#)

Department of Something, University of Whatever; Department of Whatever, University of Something

Abstract

United States has been a hub of bilateral trade, both internationally and domestically. In our study, we looked specifically at the food flows, i.e. the bilateral trade of food produce, on a domestic scale among its 50 states. We implemented two prevalent machine learning algorithms (Neural Networks and Random Forest) to predict the food flows. We also compared our results with one of the most sought-after models for bilateral trade, the gravity model to validate the relevance of application of machine learning models in our analysis. We find out that the random forest model is the best model out of all the three selected models with a R^2 score of 0.86 as compared to 0.43 R^2 score of neural networks. The gravity model performed quite poorly because of its linear nature which emphasizes the pertinence of non-linear models in food flow predictions.

Introduction and Background

The focus of our study is the bilateral food trade between states in US. To model the food flow between the studies, we examined two key studies to understand the application and effectiveness of Machine learning models where they drew parallel between ML algorithms and the popular linear gravity model. In this section we will describe those studies briefly to set-up the background for our analysis.

The article Machine learning in gravity models: An application to agricultural trade written by Munisamy Gopinath et al. employed supervised and unsupervised machine learning (ML) method to decipher patterns of international agricultural trade. Gravity model is one of the most robust empirical models to illustrate the drivers in international trade: bilateral trade between two countries is proportional to size, mostly measured in GDP and inversely proportional to “distance” between them, which commonly fitted through Poisson Pseudo Maximum Likelihood (PPML) method. Munisamy et al. leveraged the decision trees (LightGBM, XGboost), random forests and extra tree regression supervised ML algorithm to predict the bilateral trade. The data used in this project is international bilateral trade information for seven most traded food commodities from 1962 to 2016. The gain and loss was measured by adjusted R-square. The values of maximum depth of the tree, learning rate, number of leaves and feature fraction were tuned to achieve a better prediction. ML methods show more accurate prediction results than gravity model with the adjusted R-square ranged between 45 and 83%. Specifically, LightGBM had the best performance for sugar; Random Forest provided the best fit for corn; and the extra tree regressor yielded highest R-square for beef and milk powder. The size of the two countries has the largest influence to the trade, the distance follows, which is consistent to the gravity model assumptions. Munisamy et al. also employed multilayer perceptron (MLP), one paradigm of unsupervised ML method for the same datasets. The loss was measured by stochastic gradient descent (SGD) method. Unsupervised ML techniques might be a better method for longer-term trade projections than supervised ML.

Food flows between counties in the United States: Xiaowen Lin et al. developed a novel methodology to estimate food transfers between counties in the United States. They exploited the Freight Analysis Framework (FAF) dataset in their analysis to downscale the commodity transfers from FAF zones scale to county scale. FAF data categorizes food commodities into SCTG groups with two-digit codes and in this study, they have studied focused on SCTG 01-07, which are related to food commodity. To achieve this, Lin et al developed a “Food Flow Model” which is a computational algorithm that integrates machine learning, linear programming, network constraints, and mass balance. This model incorporates a gamma mixture hurdle model that uses supervised learning to develop a functional form of regression models at the FAF zone scale which is then utilized for calculating potential food transfer between counties. A key assumption made here is that network properties remain consistent across scales (Konar et al). The gamma mixture hurdle model is a two-part model where a) Hurdle

model uses logistic regression to predict the presence or absence of a link, and b) Gamma mixture model estimates the mass of the estimated link in the previous part following the assumption that food flux distributions follow gamma distribution across scales (Konar et al). Finally, they use linear programming to minimize transportation distance of food flows between counties and maintain mass balance at FAF zone level to solve the flow system at the county scale. Lin et al. used Fourier amplitude sensitivity test (FAST) to determine the most influential variables. They found out that counties in California and the Great Lakes region have the highest outflow of commodities. Also, the network density of the county scale (0.016) is much less than the network density of FAF scale (0.675). To validate their results, they compared their results with a study performed by Smith et al. where they modeled county-scale corn flows. They used R2-squared values and simple matching coefficient (SMC) to compare the results. SMC values turned out to be 1 for potential links, addressing identical estimation of presence or absence of county flow links in both models. R2-squared values are highest for outflows with a value of 0.46. Through the global sensitivity and uncertainty analysis (GSUA), where they used FAST method, they found out the distance is the most influential variable.

Munisamy et al. clearly define the quantitative ML algorithm that we can follow to predict the trade and variables we can also consider for national flows. Comparison between the ML algorithm and PPML method shows that ML is a promising method for the topic we interested in. Also, Munisamy's work indicates a potential problem we might encounter: the zero values prevalent in the trade data might impair the ML model and deviate it to a wrong direction. We suppose that using the qualitative ML method like discriminant analysis and K-Nearest Neighbors to predict the existence of trade and build the quantitative ML model on the highly potential existing link might produce a better prediction. Also, because Munisamy work focus on the international trade and our targeted scope is the US, we need to find the alternative predictors for tariff, same language and participated international trade organizations.

The paper by Lin et al. is an extensively rigorous effort to understand the food supply chain dynamics at a much finer resolution. The lack of existing literature at such finer scale, with the only literature at county-scale for corn flows by Smith et al., was another hurdle which they surpassed using a data-driven framework by bringing supervised machine learning, mass balance, and linear programming together to predict food flows for all food commodities. In their article, they addressed the shortcomings of their models, which we will be addressing in our project. Firstly, the model estimate food flows only for 2012, which means that regression models are specific to each time period and cannot be compared across different periods. Second, the model does not capture the non-linearity between environmental variables and food flow that can be captured using deep learning. Also, they used the great-circle distance between counties, which is also used in the gravity model of international trade and is a simplification of transport pathways. This can also be addressed by using the roadway network for which we have data from FAF. In our project, we will be focusing on first two key shortcomings to further research in this field.

Methods

In our study, we used two different machine learning algorithms to predict the food flows between 50 states in the US. To achieve this, our group utilized separate Exploratory data analyses and data preprocessing techniques that are different from one another. In this section, those two different approaches are explained in detail. First let's have a look at the data sources.

Data sources

We used diverse data sources to obtain data for major years and those years are 1997, 2002, 2007, 2012 and 2017. Those data sources are described briefly below.

Bilateral Trade: We obtain data on agricultural and food commodity transfers between FAF zones in the United States. The Freight Analysis Framework Version 4 (FAF4) database provides empirical agricultural and food commodity transfers between FAF zones for the year 2012 (Oak Ridge National Laboratory 2015). Second, we obtain statistical information on economic production within each US county.

Production: We used United States Department of Agriculture (USDA) data to determine state level production values for 10 selected food and crop categories in accordance with SCTG categorization, namely corn, dairy, honey, milk, oats, rice, rye, wheat, aquaculture and sorghum. The value is measured in tons.

Distance: To calculate the distance between states, we used the latitude and longitude data and used the haversine formula.

Income and GDP: We obtained economic data, i.e. income and GDP data from Bureau of Economic Analysis portal. For income we used income per capita values. All the values are measured in million dollars.

Neural network

EDA

The EDA is divided into two sections. First is Data Cleaning where we go through various key datasets and tidy them up so that they all can work together. In the second section, we will be looking at their distributions and their possible correlation between each other. To begin with, we load the data for food flow. The data has coded features which has relevant meaning. They are described below. fr: Foreign. Trade across borders. We won't be considering this data for our analysis as we are focussed on food flows among US states.

orig: Origin

dms: Domestic

dest: Destination

st: State

inmode: mode of import of international flows

outmode: mode of export of international flows

sctg2: Food category. This is a standard code used by Bureau of Trade Statistics for classifying food [BTS Guide]

value: value of food in dollars. tons: food value in weight.

After introducing the flow data, we utilized the latitude and longitude of states to calculate the haversine distance between origin and destination of food flows. Next, we introduced the remaining data from their respective files. Eventually, these data-frames were merged with once origin states and destination states to obtain the final data-frame for statistical analysis and visualization. Results of statistical analysis and visualization are presented in the results and discussion section.

Data Preprocessing

In the exploratory data analysis, we manipulated the data to find appropriate features for the model. There are 39 features in total. In this analysis, we removed self-loops in the food flow network, i.e. the origin and destination of the domestic flows being the same. Before moving towards building the neural network, we formatted the data type of different features to meet the model requirements (i.e. categorical or numerical). To obtain categorical features, we create a function that takes the string value of the column and return the whole data-frame with a new one-hot encoded feature for the feature fed to the function. We also carefully removed some null values in many features to not lose any valuable data values. This was done by removing the features with missing percentage greater than 30%. To spot any collinearity and anomalies that might affect our results, we produced the correlation matrix and histogram of all variables (refer figures). Most of the crop production values are right tailed which means that we need to normalize the data before we feed it to our model. Also, we notice that the “value” feature is highly concentrated at one bin. The standard deviation for the dataset is very high and the quantile values are quite low, which suggests that we need to remove the major outliers from the dataset.

Model

In this analysis, we used neural network to predict the food flow between US states. We used two hidden layers with 30 and 12 neurons respectively with ReLu activation function. We used Adam optimizer function with a learning rate of 0.001, “mean squared error” loss function and RMSE as the metric. Batch size of 100 was considered with a validation split of 0.2 and number of epochs equals to 25. The results of the model run are discussed in the results and discussion section.

Results and Discussion

Conclusion

References
