

# Classification of Review Sentiments

**Vivek Kumar**

Civil and Environmental Engineering  
Princeton University  
vivekk@princeton.edu

**Victor Charpentier**

Civil and Environmental Engineering  
Princeton University  
@princeton.edu

## Abstract

Applying machine learning techniques for social good is one of the bright outcomes of the technological revolution of the past years. In homework 2 we have performed a prediction and analysis of 6 key outcomes at age 15 of the fragile family challenge. The methods implemented depend on the type of outcome. For each type of outcome over 6 prediction methods are tested. The performance of regressions and classifiers is assessed with xx and receiver operating curves. Overall, the random forest regression methods has the best results for continuous variables and XX has the best result for binary outcomes. Finally, the confidence intervals of the predictors are computed to quantify the quality of the prediction.

## 1 Introduction

The fragile family challenge (FFC) is a Princeton University-led initiative opening the data collected in the long term Fragile Family and Child Wellbeing Study (FFCWS) to a variety of data scientists. The study follows an ensemble of nearly 4700 families with children born in 20 cities across the U.S. between 1998 and 2000. The study has an over-representation of non-marital birth compared to national average. The goal of the study is to provide information on the family factors that influence the life of the child [3]. The FFC was officially closed in May 2017, and its winner announced in the fall of 2017. With an increased understanding of the data gained from the challenge, we aim to apply regressions and classifications to the dataset of the FFC to predict the 6 key outcomes. Regressions and classification constitute the backbone of supervised learning. However beyond the purely mathematical aspect of the challenge, feature selection and dataset preparation are some of the most important aspects of a real machine learning project. So important that some describe the refinement of feature or feature engineering as a "black art" [2]. With respect to this mystical reference we have worked to prepare the training data for prediction. The continuous outcomes (GPA, grit and material hardship) have been predicted using Random Forest Regressor (RF), Ada Boost Regressor (AB), lasso regression (LRCV), ridge regression (RRCV), elastic net (EN), extra trees regressor (ET) and multi-layer perceptron regressor (MLP). The binary outcomes (eviction, job training and layoff) were predicted with quadratic discriminant analysis (QLA), logistic regression (LR), random forest classifier (RFC), ada boost classifier (AB), multi-layer perceptron classifier (MLP) and extra trees classifier (EF).

## 2 Description of the Data

In the domains that generate a lot of data, social sciences are maybe one of the least explored by data scientist. Social sciences excel in causal inference, they have handled data from observations for a long time. Data science on the other end has experience with large dataset but is not focused so much on causal inference [? ]. The hope of applying machine learning techniques to the FFCWS dataset is to combine both sciences for the understanding of fragile families outcomes. Improving the methods to treat large datasets is vital for social sciences and the potential outcomes could be transformative for society.

Inequalities in the U.S. are threatening the model of the American Dream. Some voices from the highest economic international agencies are calling for a new American Dream, "*based on equality and sustainable growth*" [? ]. In this context, understanding what factors positively influence children growing up in some of the most difficult family environments may be key to creating policies more able to grant equal opportunities for all. The study follows an ensemble of nearly 4700 families with children born in 20 cities across the U.S. between 1998 and 2000. The study has an over-representation of non-marital birth compared to national average. The goal of the study is to provide information on the family factors that influence the life of the child [3].

### 3 Overview of the methods and fit to reference data

#### 3.1 Data wrangling

The input data for the homework consists of two files. The true values of the outcomes for 2121 families and the data set of over 15 000 features for all the 4242 families. As identified in the introduction, feature selection and engineering can be seen as an "art". Therefore we spent some time on making sure the quality of data that was going to be used for prediction was satisfying. The preparation of the true outcomes is a two step process. In the first step, the families that have outcomes values equal to 6 NAs are dropped. The goal of that step was to reduce the skew of the training data toward the average of each outcome. 655 families were removed from the training data as a result. This is a big number but we assume that the missing information of the outcomes carries some meaning. The training outcomes still contain some NA values. However for these cases some of the other outcomes are available. In the second step the dataset is cleaned. First the constant variables provided in the file `constantVariables.txt` are removed. Some survey weights of the type "q5natwt\_XXXX", "q5natwtX\_XXXX" and "q5citywt\_XXXX" were placed in the list of constant variables. We decided to consider the information embedded in missing data. All values of missing data less than -5 are considered empty from meaning and set to 0. The missing data from -1 to -4 is maintained. Finally the answers of the type "if not in the list specify:" (e.g. d3c5a\_14ot) are coded to 1 if there is a string in the cell and to 0 if the cell is empty.

#### 3.2 Prediction

### 4 Detailed presentation of the regression forest

The random forest predictor used for regression, also called regression forest (RF), is based on the stand alone model of decision trees. Decision trees are prone to over fitting and sensitive to missing data. Therefore, a forest stems (pun) from the aggregate of many weak models to produce a better overall prediction. As opposed to a classification forest, a regression forest provides prediction of continuous variables. The input data seen on Figure 1 is continuous. Therefore the leaf nodes predict real values (as opposed to classes). The data-set is split based on homogeneity of data (with the standard deviation). This leads to subsets of the data contains similar values of the data to predict (cf. Figure 1(a)). The similarity is quantified by entropy, a measure of predictability. The form of the predictor can vary once the tree is fully grown. Several types of predictors can be used [1]. They can be constant, linear, polynomial or probabilistic-linear among others 2.

In the case of the constant predictor, the value of the predictor is given by the minimization of the sum of squared error (SSE) as given by Equation 1 for the subdomain  $D_k$ .

$$\hat{y}_k = \arg \min_y \sum_{i \in D_k} (y - y_i)^2 \quad (1)$$

Two parameters appears as main control nobs of the method: the number of regression trees in the ensemble and the depth of those trees. Each of those parameter has a distinct effect on the prediction. The depth of the trees controls the closeness of the fit. A ensemble of trees of depth 1 will correspond to a linear regression since each tree corresponds to its root node. In comparison, an ensemble of very deep trees will risk being overfitted. The number of trees in the forest influences the smoothness of the prediction, the more trees the smoother the direction of the prediction. Finally

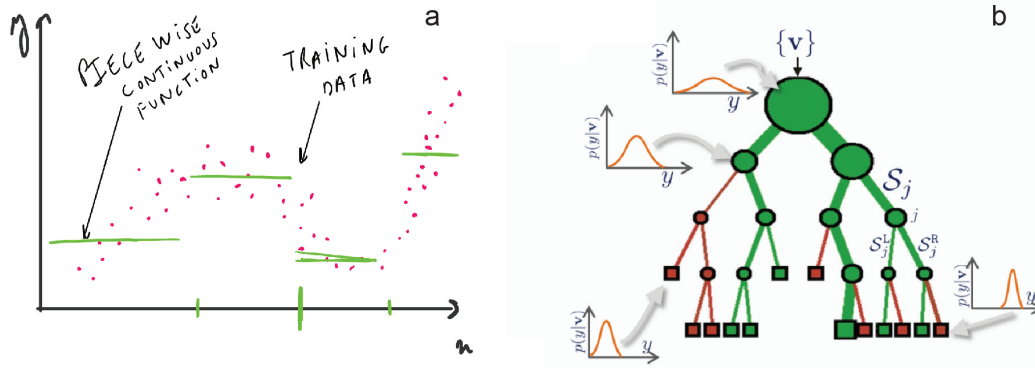


Figure 1: Detail of a regression forest - (a) the continuous data is approximated by a piecewise continuous function. The subdivision of the  $x$  interval stops when the entropy or SSE has reached a threshold value - (b) an example of regression tree from [1], the leaves of the tree determine values of continuous outputs for subdomains of  $x$

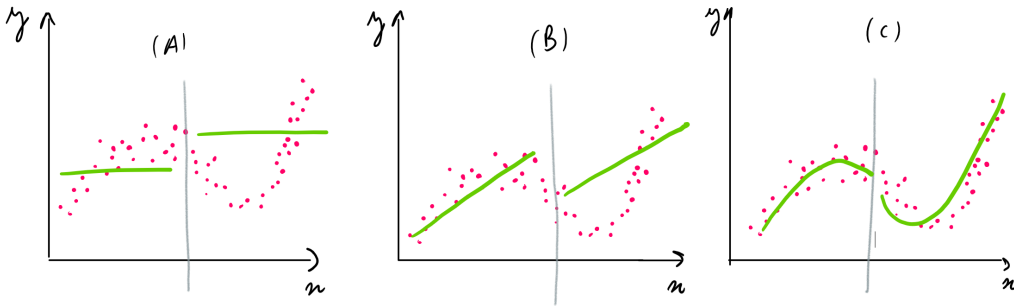


Figure 2: Example of three predictor models (a) constant - (b) linear and (c) polynomial

regression forest is a preferred method due to its speed, efficiency, ability to cope with missing data and flexibility of use.

## 5 Presentation of results

1. Build a model for predicting 1-type of variable well.
2. GPA, Grit, Material Hardship
  - (a) Hypothesis : Positive Environment + Lack of Negative leads to positive GPA + ..
  - (b) Hypothesis : Positive Environment + Plus some negativity + ..
  - (c) Hypothesis : Negative Environment + Lack of parent + ..
3. Drop NA values for training: only in case of NA in all.

## 6 Discussion of the results

## 7 Conclusions

## References

- [1] A Criminisi, J Shotton, and E Konukoglu. Decision forests for classification, regression, density estimation, manifold learning and semi-supervised learning. *Microsoft Research*, 2011.
- [2] Pedro Domingos. A few useful things to know about machine learning. *Communication of the ACM*, 55(10):78–87, 2012.

162 [3] Nancy E Reichman, Julien O Teitler, Irwin Garfinkel, and Sara S McLanahan. Fragile families:  
163 Sample and design. *Children and Youth Services Review*, 23(4-5):303–326, 2001.  
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