

ANALYSIS AND PREDICTION OF POVERTY

T. H. Johnny Yiu

July 2019

1. EXECUTIVE SUMMARY

This report presents an analysis on data concerning the Poverty Probability Index (PPI) of individuals. By exploring the relationships between characteristics of the individuals and their PPI, we extract useful information that would allow us to make prediction of an individual's probability of poverty based on its socioeconomic indicators.

The dataset was retrieved from datasciencecapstone.org and contains the Poverty Probability Index (PPI) along with 58 features of 12,600 individuals across 7 different countries. Original sources of the dataset include the PPI website and Financial Inclusion Insights household surveys conducted by InterMedia. The PPI is a measure to calculate an individual's probability of living below the poverty line at the \$2.50/day threshold using 10 questions about a household's characteristics and asset ownership.

The report splits into two main parts. The first part being an exploratory data analysis, of which we calculate summary statistics and create data visualizations of the data that identify potential relationships between individuals' characteristics and their PPI. Based on the information extracted in the first part, we then create a predictive machine learning model in the second part that predicts the PPI of individuals.

It is concluded that the following features are the most important indicators of an individual's PPI (in no particular order):

- **avg_shock_strength_last_year** - Average strength of shocks experienced the past year
- **num_financial_activities_last_year** - Number of financial activities conducted the past year
- **country** - Unique identifier for country of residence (masked)
- **can_calc_percents** - Ability to calculate percentages
- **religion** - Unique identifier for religion (masked)
- **is_urban** - Residence in the urban area (vs. rural)
- **age_group** - Age group
- **female** - Sex (True = female, False = male)
- **married** - Marital status
- **relationship_to_hh_head** - Role in the family

Due to high variability and noise around the data, our machine learning model could only make a prediction on an individual's PPI with a low effect size. The adjusted r-squared for our model is 0.443.

2. EXPLORATORY DATA ANALYSIS

2.1. DATA PRE-PROCESSING

Before calculating summary statistics and creating visualizations for our data, the dataset was prepared into suitable formats for analysis. We have:

- removed duplicates
- dropped uninformative features with largely missing values, they include:
 - *bank_interest_rate*
 - *nm_interest_rate*
 - *mfi_interest_rate*
 - *other_fsp_interest_rate*
- substituted missing values for *education_level* and *share_hh_income_provided* with a category of their own

2.2. EXPLORING PPI

The summary statistics for PPI is shown in the table below:

	Count	Mean	Std	Min	25%	Median	75%	Max
PPI	12600	0.611272	0.291476	0	0.394	0.633	0.879	1

Table 1. Summary Statistics of PPI

It is observed that the median is greater than the mean and that the standard deviation is about half of the mean. The distribution should be left-skewed and has some variance, which is verified by the following histogram:

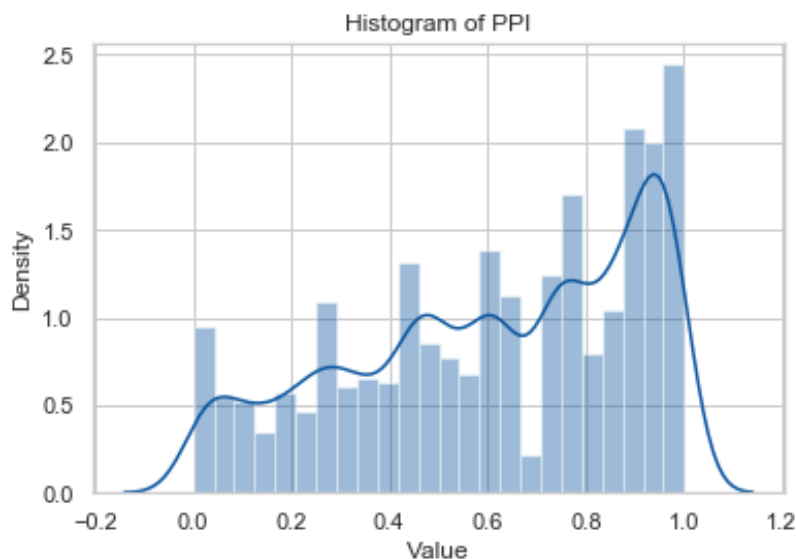


Figure 1. Histogram of PPI

Our PPI histogram shows that there are more people having a higher probability of poverty than a lower one. Besides, the multimodal characteristic entails peaks in 6 different PPIs. As an attempt

to make the distribution more symmetric, transformations of the values are performed. They include natural log (the natural log of 0 was converted into 0), cube root, and square transformations. However, these transformations did not make the distribution more normal, nor did they improve the performance of our machine learning model.

2.3. EXPLORING NUMERIC FEATURES

There are 3 numeric features in the dataset, namely *age*, *avg_shock_strength_last_year* and *num_financial_activities_last_year*. The respective summary statistics and histograms are shown below:

	age	avg_shock_strength_last_year	num_financial_activities_last_year
Count	12600	12600	12600
Mean	36.28071	2.112765	1.559683
Std	15.14594	2.019239	2.043831
Min	15	0	0
25%	25	0	0
Median	33	2	1
75%	45	4	3
Max	115	5	10

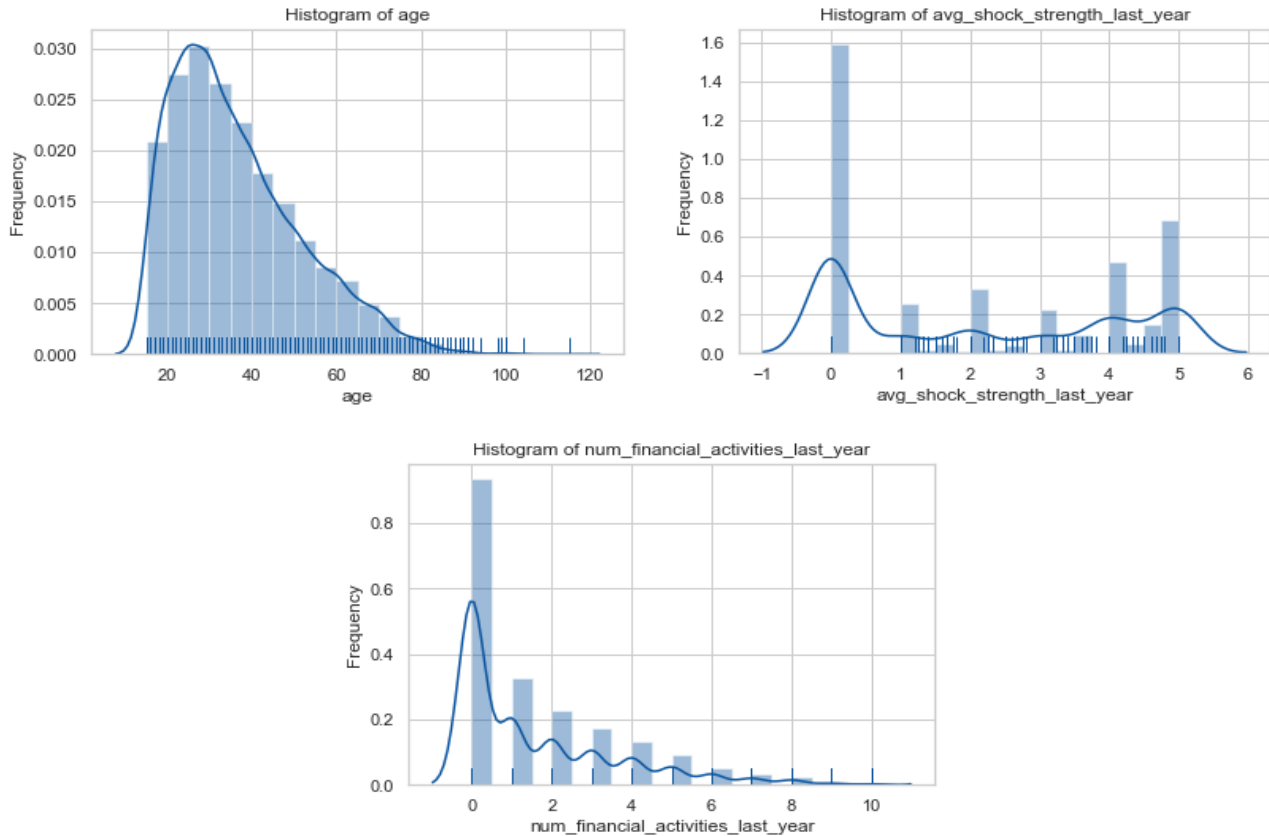


Figure 2. Histograms of Numeric Features

They indicate the following:

- Most respondents are in the age range of 20-40
- The majority of respondents experienced a shock strength of 0 in the past year, and a significant number of people experienced a 4 or above
- Most respondents were not involved in any financial activities last year
- Having 6 financial activities or above is uncommon

A pair-plot and a correlation matrix of the numeric features are created:

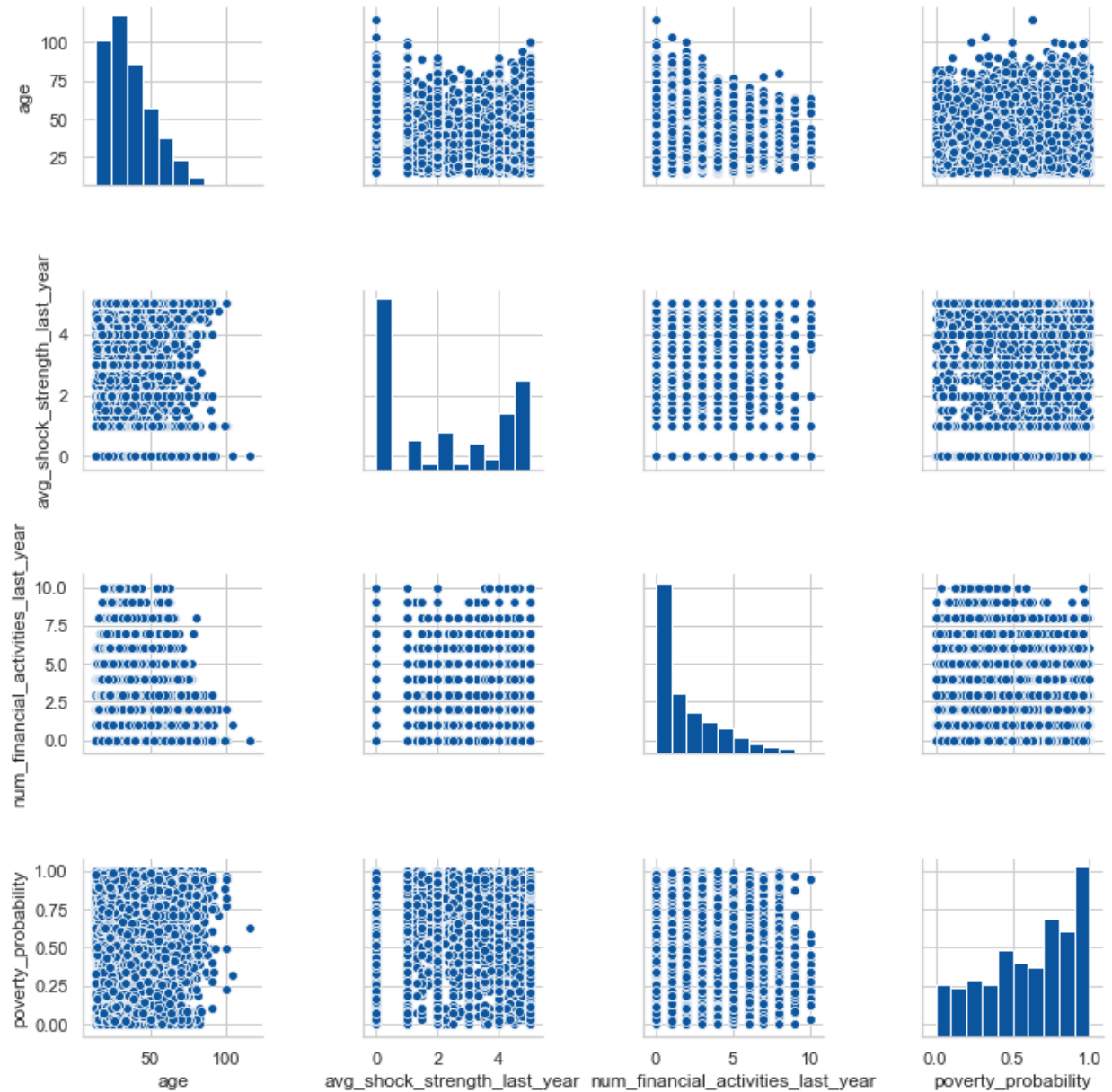


Figure 3. Pair-plot of Numeric Features

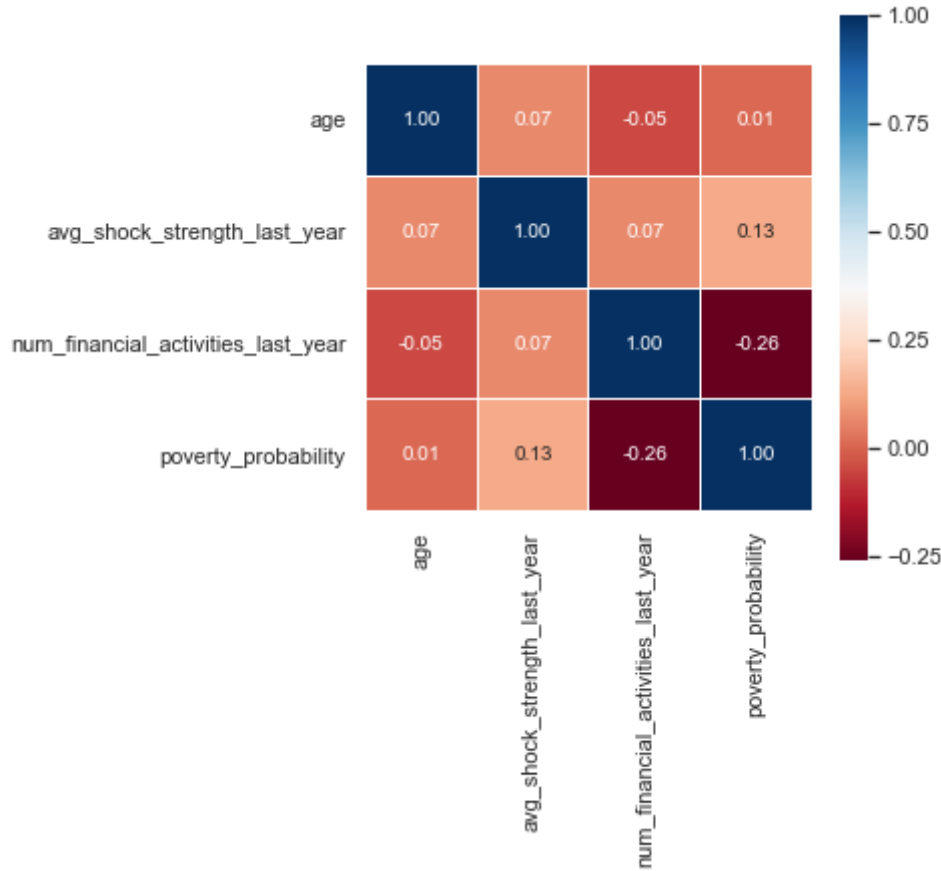


Figure 4. Correlation Matrix of Numeric Features

Little apparent relationships are found within these features and between the label and the features, other than that:

- *On average, the amount of financial activities seems to decrease with age*

2.4. EXPLORING CATEGORICAL FEATURES

After assigning some discrete variables as categorical features, we ended up with 52 categorical features. Bar charts were created to study the frequency of these features, and indicate the following:

- 86% of the respondents are either in the religion Q or X
- Over 80% of the respondents can either add or divide
- 75% of the respondents are either the head of the family or the spouse
- 75% of the respondents do not have the ability to use internet on their phone
- 72% of the respondents do not have a bank account under their own name
- 70% of the respondents do not have any form of investment
- Only 13% of the respondents have insurance
- Only 6% of the respondents received income from the government last year
- Only 5% of the respondents experienced 4 shocks or above last year

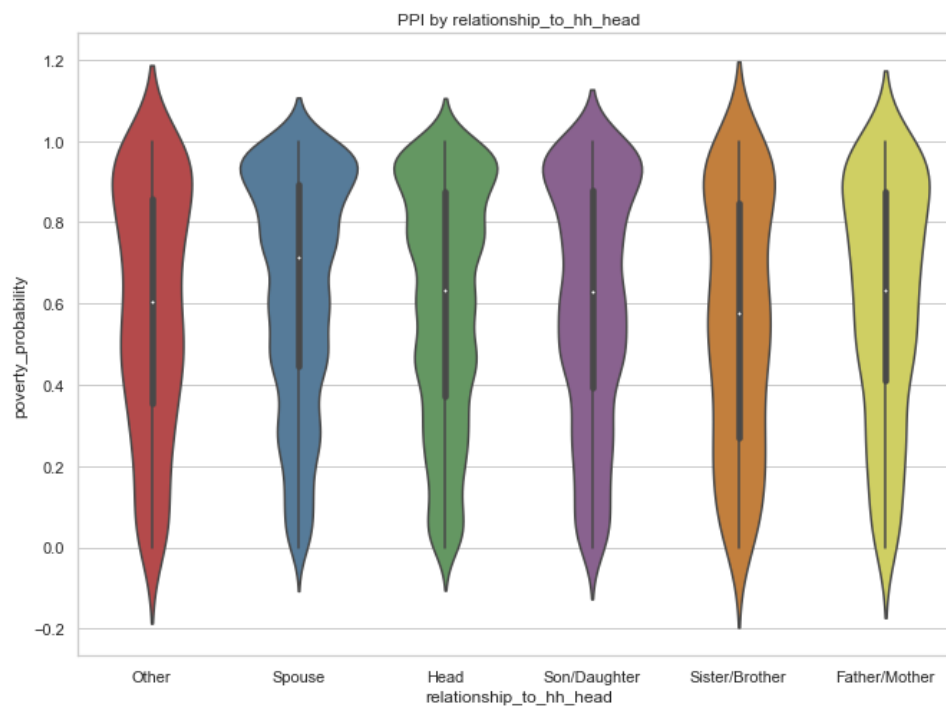
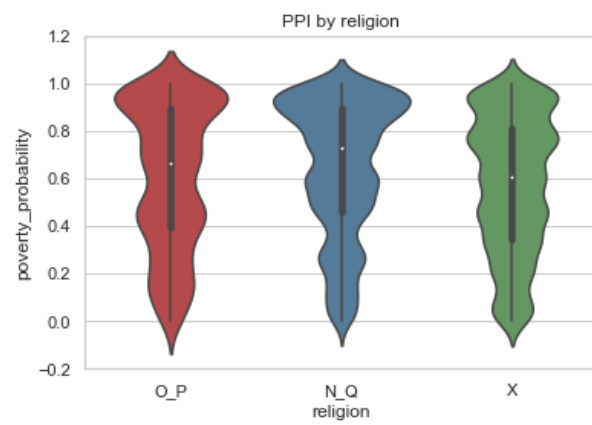
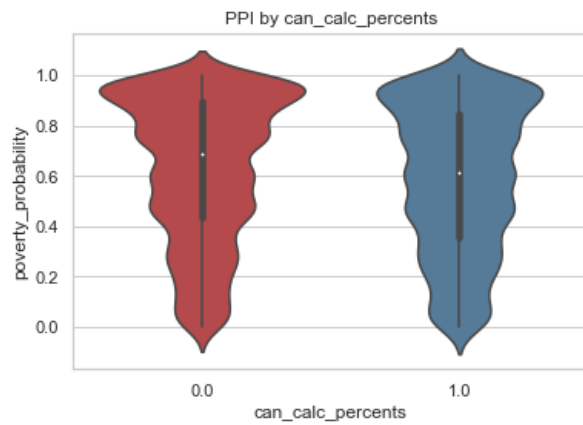
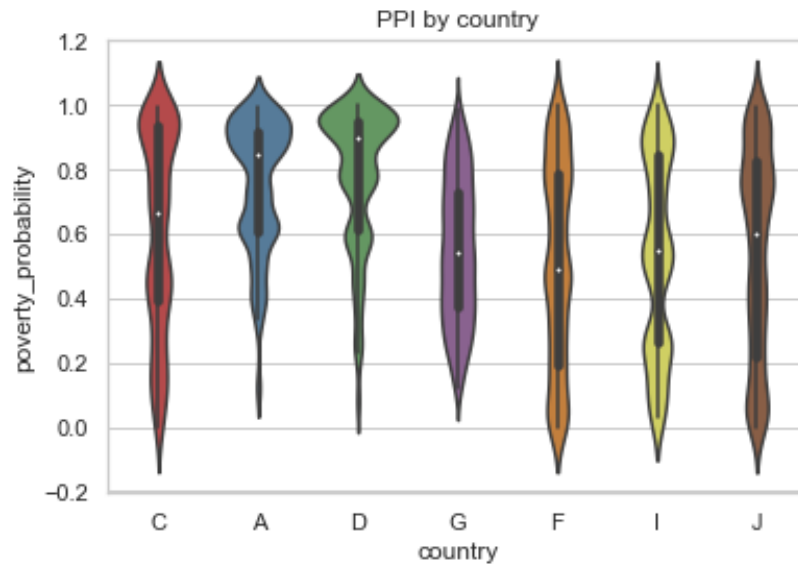
- Only 4% of the respondents received income from the public sector last year
- Less than 3% of the respondents used services of 3 or more formal financial institutions last year
- Only 2% of the respondents used services of 2 or more informal financial institutions last year
- Respondents in the religion *O* or *N* account for less than 1% of the sample

2.4.1. FEATURE ENGINEERING

The aggregation of categorical features was performed to reduce the number of categories. For discrete variables, rare values are combined to form a range of values. For categorical variables, rare values are combined with common values that share a more similar distribution in **PPI**:

- Religion *O* is combined with religion *P* to form religion *O_P*, and religion *N* is combined with religion *Q* to form religion '*N_Q*'
- The number of shocks of 4 or above are combined to form '*4_5*'
- The number of formal financial institutions used of 3 or above are combined to form '*3_4_5_6*'
- The number of informal financial institutions used of 2 or above are combined to form '*2_3_4*'
- The number of shocks of 4 or above are combined to form '*4_5*'
- The category '*Unknown*' for relationship to the head of the household is combined with '*Other*'

Violins plots of the most important features are created to study their relationships with **PPI**:



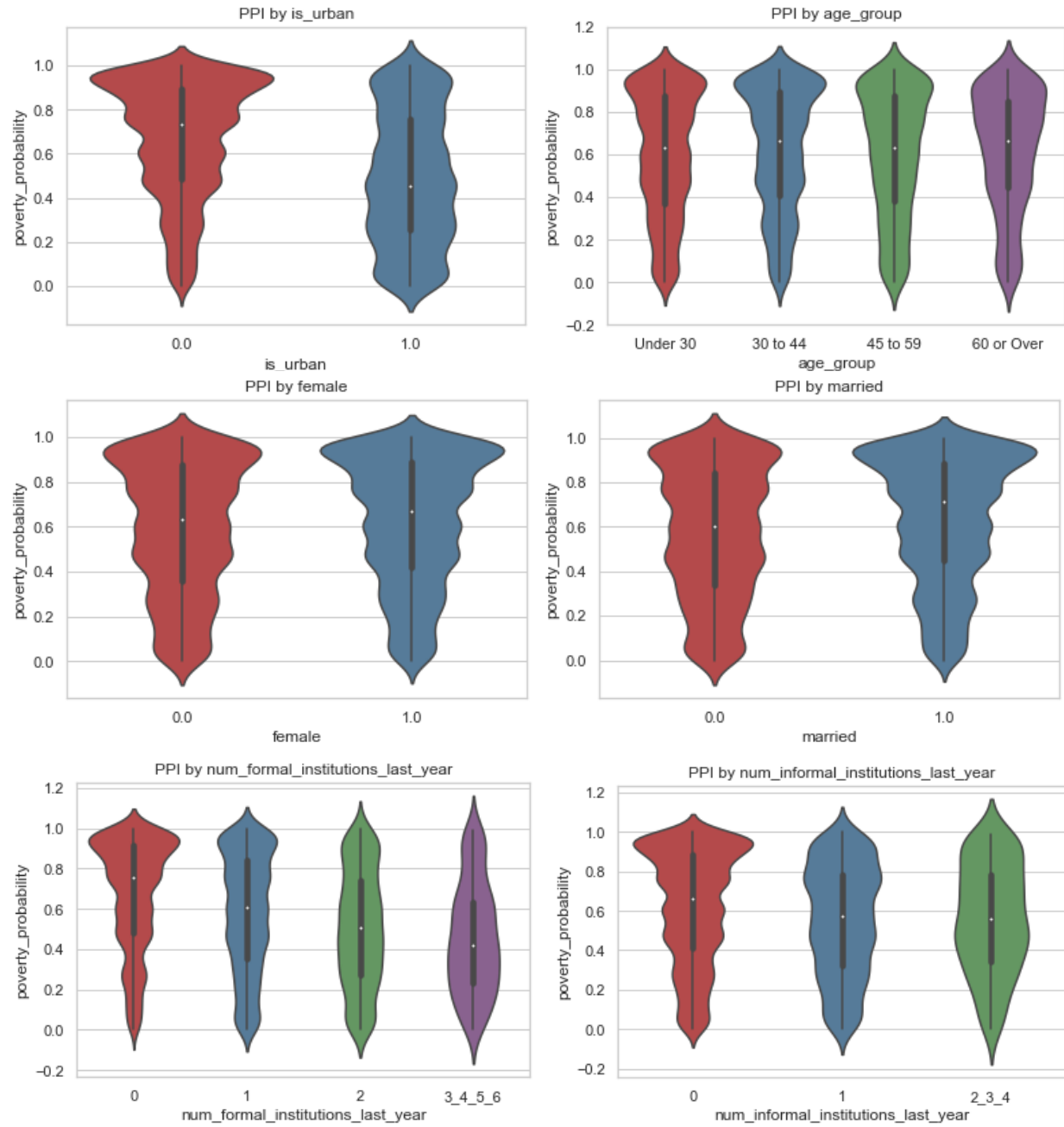


Figure 5. Violin Plots of Categorical Features

They indicate the following:

- On average, country D has the highest probability of poverty
- On average, respondents that cannot calculate percentages have a higher probability of poverty
- On average, respondents in the religion N or Q have a higher probability of poverty
- On average, spouses have a higher probability of poverty
- On average, rural residents have a much higher probability of poverty
- On average, females have a higher probability of poverty
- On average, married persons have a higher probability of poverty
- On average, respondents that do not use services of a financial institution have a higher probability of poverty

After aggregating the categorical features, we eliminated features with low variance. For the variance threshold

$$Var(x) = p(1-p)$$

, we have used a $p = 0.95$ and eliminated any features under the threshold. We further selected informative features with Recursive Feature Elimination using cross validation, and have listed those with the highest feature importance scores in the Executive Summary.

3. POVERTY PREDICTION

Based on the analysis of the probability of poverty, a regression model to predict PPI was created.

The model was created based on a boosted decision trees algorithm and trained with 80% of the data. Testing the model with the remaining 20% of the data yielded the following results:

Mean Square Error	0.0469
Root Mean Square Error	0.217
Mean Absolute Error	0.174
Median Absolute Error	0.147
Adjusted R-squared	0.443

Table 2. Prediction Model Test Results

The result of our test shows that the model only explains 44.3% of the variability of the response data around its mean. Due to high variability and noise around the data, our model could only make a prediction on an individual's PPI with a low effect size. Plot on residuals are shown below:

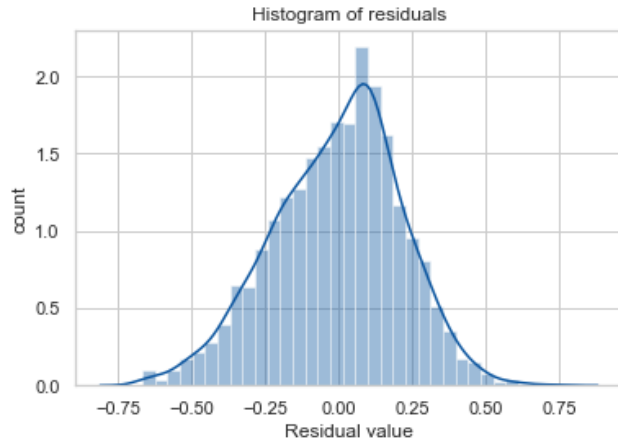


Figure 6. Histogram of Residuals

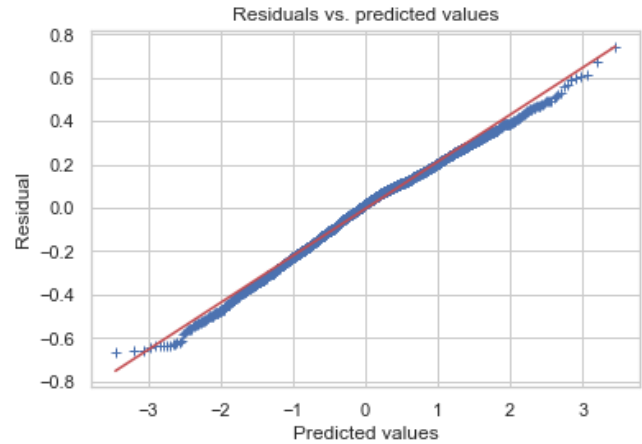


Figure 7. Quantile-Quantile (Q-Q) Plot of Residuals vs. Predicted Values

The above plots show a wide range of residuals across predicted values. However, the distribution of residuals is close to normal, which indicates that the amount of predictive ability the model has is consistent across the full range of the label (PPI).

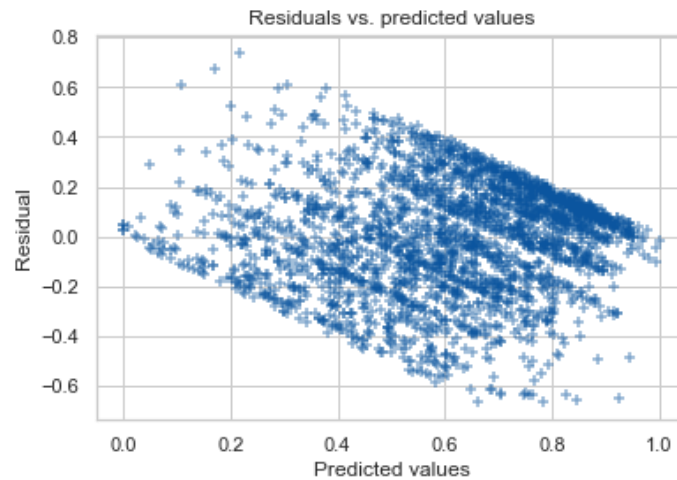


Figure 8. Scatter Plot of Residuals Vs. Predicted Values

The scatter plot of residuals vs. predicted values above shows that the model performs well when predicting individuals as having a probability of poverty of 0 or 1 but becomes progressively worse when predicting individuals as having a medium probability of poverty.

4. CONCLUSION

This report has extracted important demographic information of the respondents and concluded that an individual's probability of poverty can only be roughly predicted from its socioeconomic indicators, especially when they have a probability of poverty within the medium range in reality. Nonetheless, we have identified more important indicators of the Poverty Probability Index. They include the average strength of shocks experienced the past year, number of financial activities conducted the past year, country of residence, the ability to calculate percentages, religion, residence in the urban area, age group, sex, marital status, and the role in a family.