

TELL ME YOUR DEMOGRAPHIC, I’LL GUESS YOUR LEVEL OF EDUCATION

Using Census Data To Predict Levels Of Education

†Richard Muniu
‡Department of Computer Science, Swarthmore College



Introduction

Demographic data is often used in Machine Learning to predict people’s levels of income. The target label represents their level of income, which is discretized into two classes: above 50,000 per year, or below 50,000 per year.

Instead of using different machine learning techniques to perform and optimize this classical classification problem, we chose to look into how the very same demographic data, including one’s income level, can be used to predict a preson’s level of education.

Data

Our dataset was the 1994 Census Data, obtained from the UCI Machine Learning Repository. It contained about 48,842 data points, with a spread of continuous and discrete features. This included 32,561 training examples and 16,281 testing ones. We randomly shuffled the data and split it into train and test sets, using the labels *{HS, Some HS, HS grad, Some College, College Grad, Masters, Doctorate}* for multi-class classification and *{No College vs College}* for binary classification.

The data attributes included *Age, Work Class, Income, Marital Status, Occupation, Relationship, Sex, Hours Per Week, Native Country*, and *Education*, which we used as our label.

Methods

We used three different classifiers in this task:

- SUPPORT VECTOR MACHINES
 - To pre-process the data for this classifier, we used a one-hot encoding for discrete features, leaving continuous features as they were.
 - Unknown values were replaced with a ‘0’ in our encoding scheme.
 - For multi-class classification, we used a one-vs-rest classifier.
 - This method involved hyper-parameter tuning for our LinearSVC.
- DECISION TREES
 - To pre-process our data, we converted our continuous features to discrete ones, and binarized discrete features. The continuous features were split using arbitrary bin sizes.
 - We used an off-the-shelf implementation of Decision Trees from **scikit-learn**.
- NAIVE BAYES
 - We created an ”unkonwn” value for each unknown feature in the original dataset.
 - Next, we discretized all the continuous features.
 - We then used an off-the-shelf implementation of Decision Trees from **scikit-learn**.

We also used a Most Frequent Class (MFC) classifier as a performance baseline for the aforementioned classifiers. The MFC classifier uses the most frequent label in the data as its classification output (prediction label).

Results

Our MFC classifier had an accuracy score of 36.5% over the 7 possible outputs.

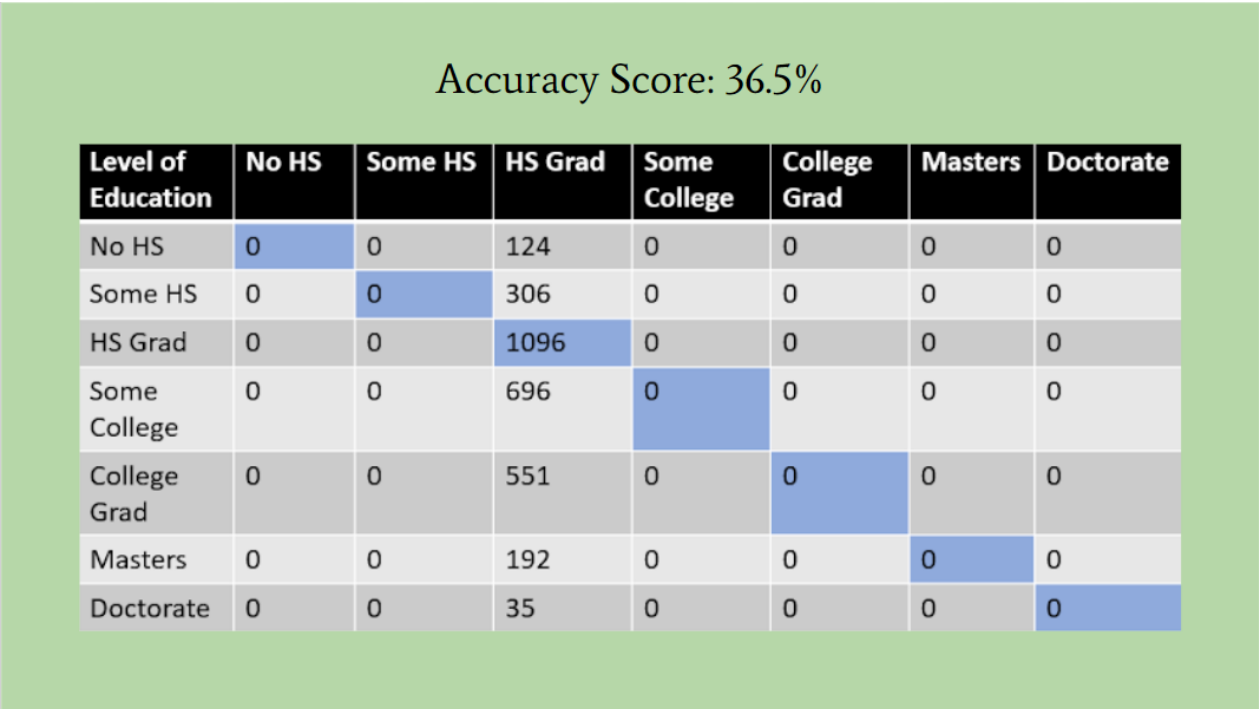


Fig. 1: Most Frequent Classifier Accuracy

Our SVM classifier achieved the highest accuracy score, at 44.0%. This was closely followed by the Decision Tree classifier at 41.4% and lastly Naive Bayes, which achieved 35.6% accuracy.

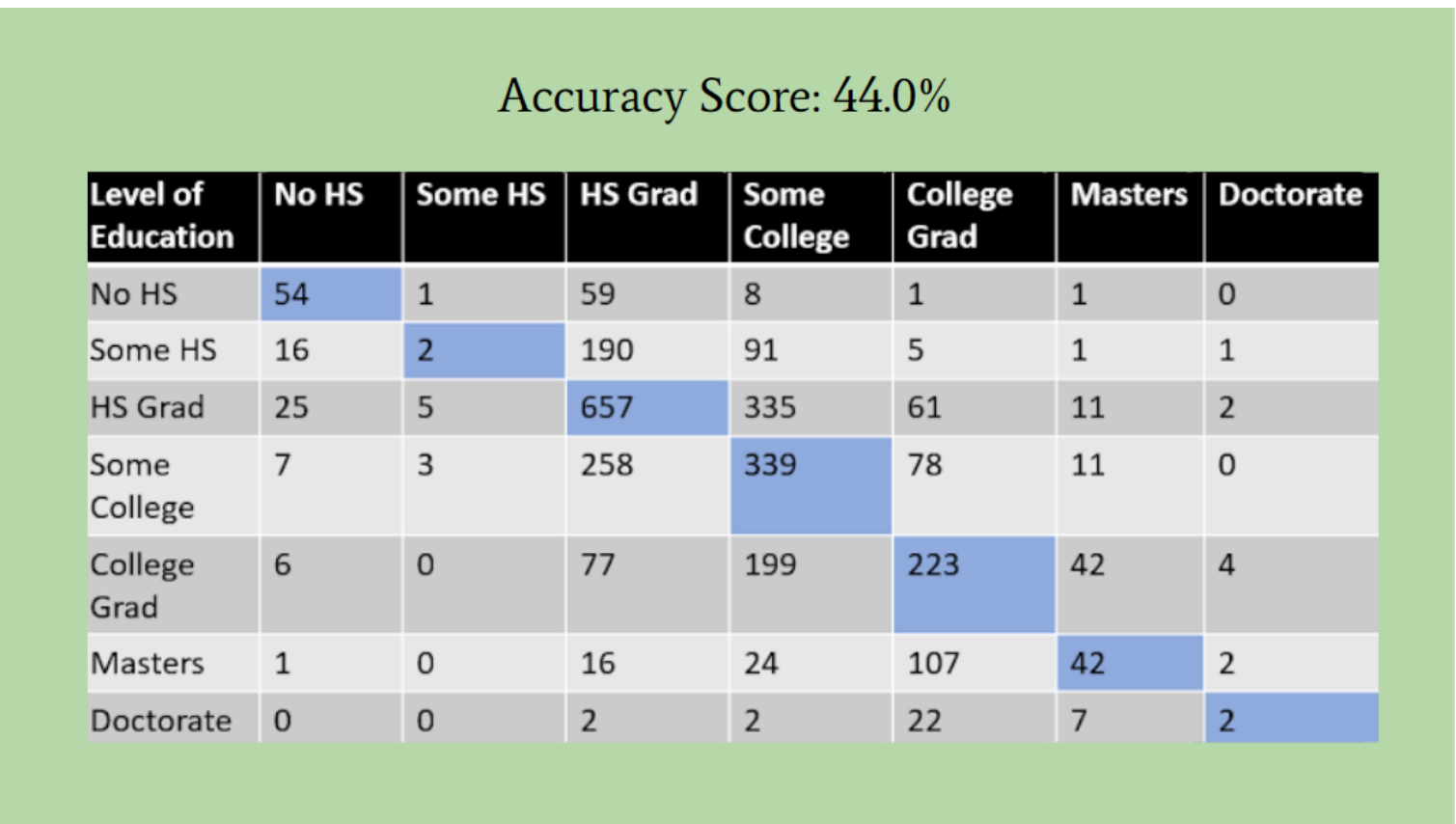


Fig. 2: SVM accuracy

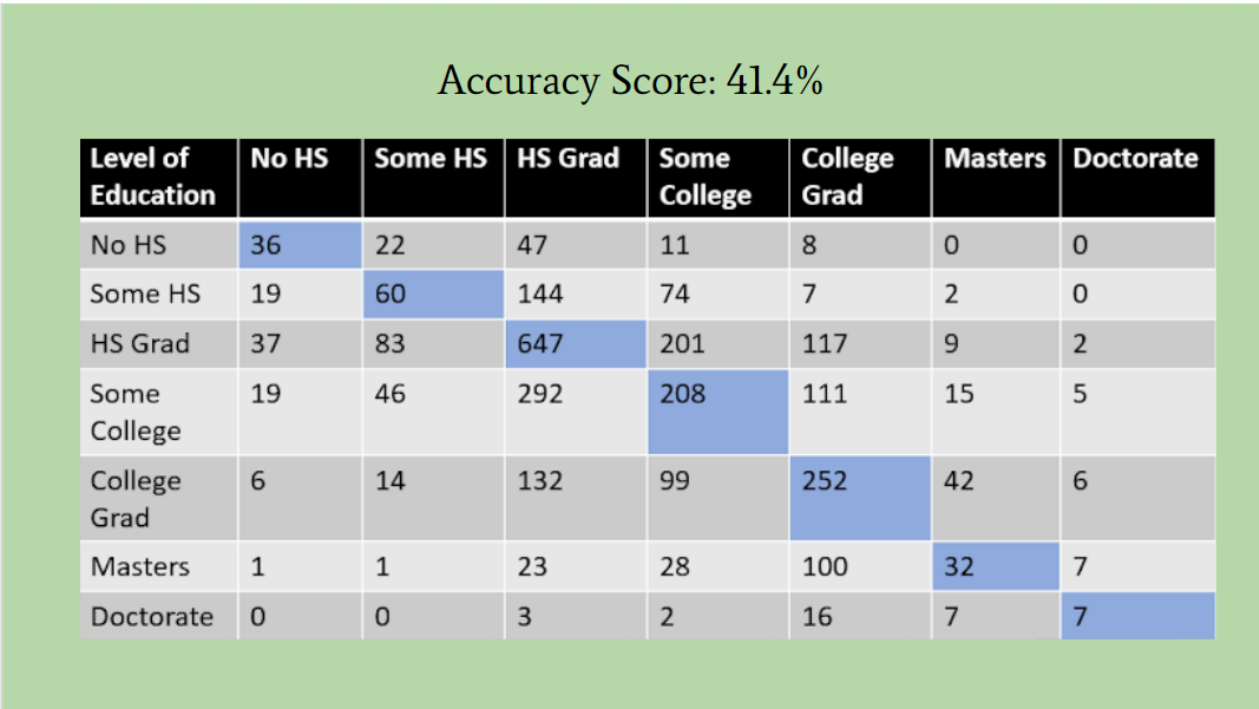


Fig. 3: Decision Tree Accuracy

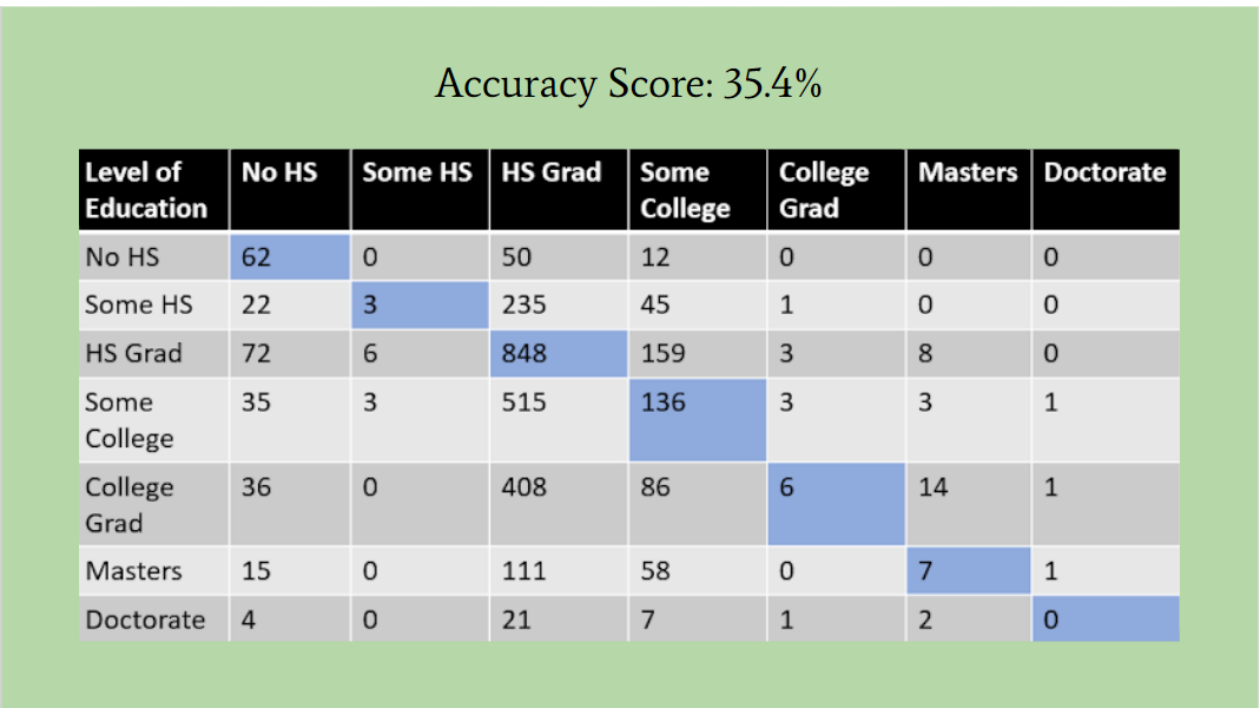


Fig. 4: Naive Bayes Accuracy

Comparison

Our overall comparison of all the classifiers for both binary and multi-class classification is as follows:

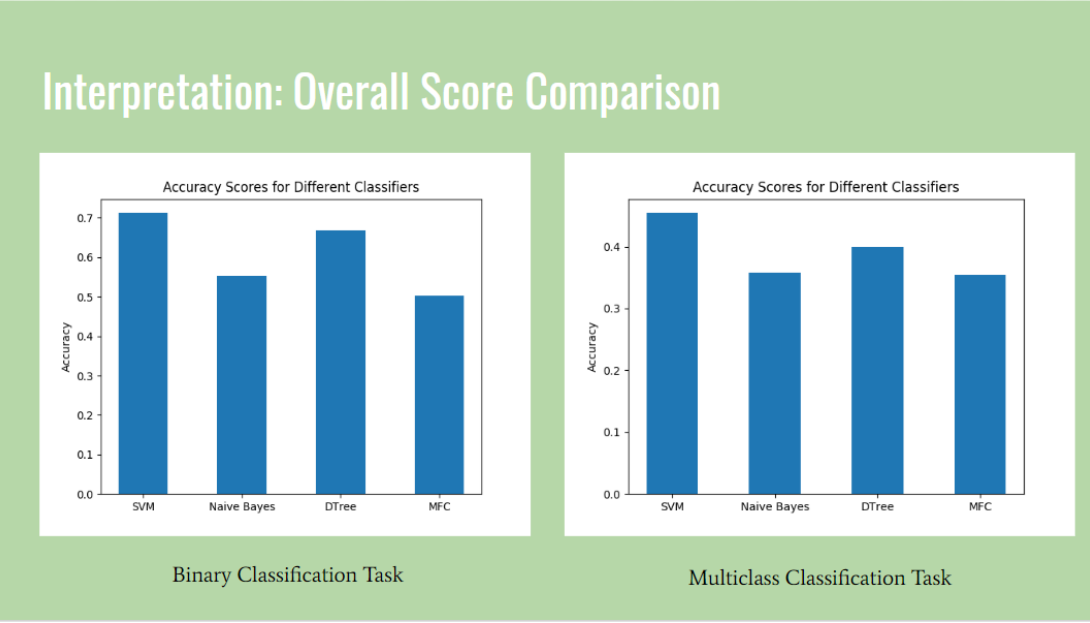


Fig. 5: Overall Score Comparison

Feature Analysis

We used our LinearSVC to compare the relative importance of features, by analyzing the nxn matrix of coefficients. The five most important features determining one’s level are shown below:

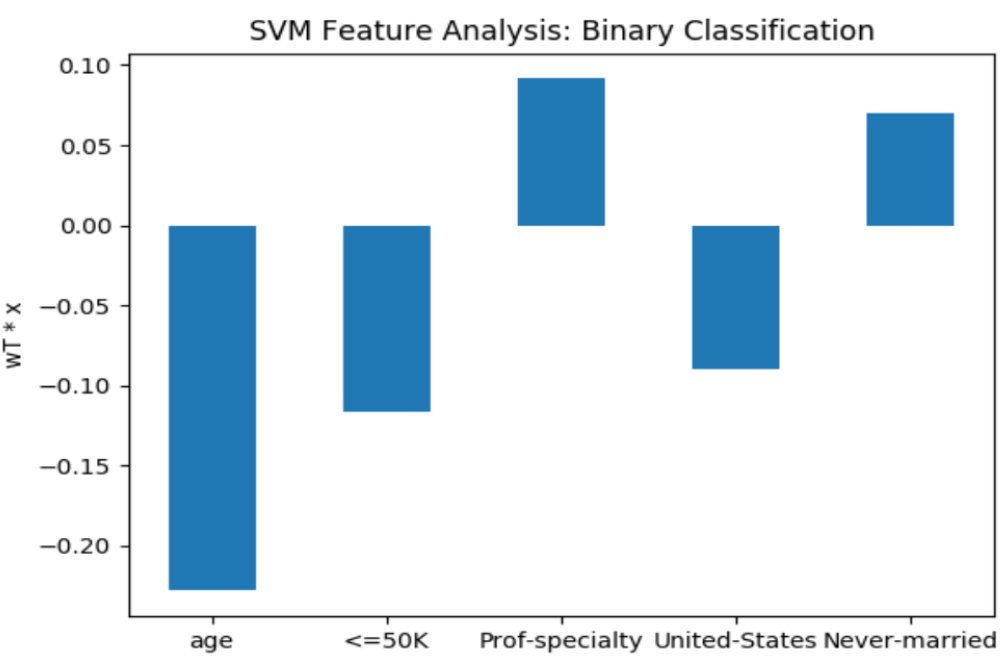


Fig. 6: Overall Score Comparison

Conclusion

Our SVM model produced best results overall, and income, occupation, and age were seen to be the most important features (according to our Tree and SVM analysis).

Future Considerations

We could account for colinear independent variables and view problem as regression rather than classification.

Acknowledgements

I would like to thank Sara Matheson for her constant encouragement and support throughout the project.I’d also like to thank my teammates Kenny Gwon and Raymond Liu for their collaborative efforts. Special thanks to the UCI Repository for making the dataset available..