

Predicting Years of Education and Human Development Index Values Using 1994 Census Data Grouped by Controllable and Uncontrollable Attributes

Final Project Report
Matrix Methods & Data Science
MAS 4106 - CRN 14440

Instructor: Dr. Alberto Condori

Benton Stacy
(bmstacy7127@eagle.fgcu.edu)
Katarya Johnson-Williams
(kajohnsonwilliam3168@eagle.fgcu.edu)

Table of Contents

Abstract.....	4
Introduction.....	4
Methodology.....	5
Data, Figures, and Tables.....	7
Discovery Models.....	7
Discovery Models: Initial Least Squares.....	8
Least Squares Model 1: Education Prediction Using Age and Hours Worked Per Week.....	10
Least Squares Model 2: Education Prediction Using Work Class and Occupation.....	10
Least Squares Model 3: Education Prediction Using Marital Status and Native Country.....	11
Least Squares Model 4: Education Prediction Using Race and Sex.....	11
Least Squares Model 5: Education Prediction Using Gross Domestic Product and Income.....	12
Least Squares Model 6: Education Predicting Using Human Development Index and Income.....	12
Discovery Models: Cross-Validation.....	13
Cross-Validation Model 1: Education Prediction Using Age and Hours Worked Per Week.....	15
Cross-Validation Model 2: Education Prediction Using Work Class and Occupation.....	15
Cross-Validation Model 3: Education Prediction Using Marital Status and Native Country.....	16
Cross-Validation Model 4: Education Prediction Using Race and Sex.....	16
Cross-Validation Model 5: Education Prediction Using Gross Domestic Product and Income.....	17
Cross-Validation Model 6: Education Predicting Using Human Development Index and Income....	17
Predicting Years of Education with Controllable and Uncontrollable Attributes.....	18
Predicting Years of Education: Initial Least Squares.....	18
X-Hat Values Model 7: Education Prediction Using Uncontrollable Attributes.....	18
X-Hat Values Model 8: Education Prediction Using Controllable Attributes.....	19
Least Squares Model 7: Education Prediction Using Uncontrollable Attributes.....	20
Least Squares Model 8: Education Prediction Using Controllable Attributes.....	20
Predicting Years of Education: Cross-Validation.....	21
Cross-Validation Model 7: Education Prediction Using Uncontrollable Attributes.....	22
Cross-Validation Model 8: Education Prediction Using Controllable Attributes.....	22
Predicting Human Development Index with Controllable and Uncontrollable Attributes.....	23
Predicting Human Development Index: Initial Least Squares.....	23
X-Hat Values Model 9: Human Development Index Prediction Using Uncontrollable Attributes....	23
X-Hat Values Model 10: Human Development Prediction Using Controllable Attributes.....	24
Least Squares Model 9: Human Development Index Prediction Using Uncontrollable Attributes....	26

Least Squares Model 10: Human Development Index Prediction Using Controllable Attributes.....	26
Predicting Human Development Index: Cross-Validation	27
Cross-Validation Model 9: Human Development Index Prediction Using Uncontrollable Attributes	28
Cross-Validation Model 10: Human Development Index Prediction Using Controllable Attributes.	28
Results and Future Work	29
References.....	30

Abstract

The reason for writing this report is to explore the possible ways to predict years of education using census data. The processes and methods developed in this report could be applied to future census data to create an effective method of predicting significant census attributes. The problem this report seeks to solve is identifying patterns in the data provided by the census that can be used to predict years of education an individual has completed. A secondary problem the methodology in this report can be used to solve is to fill in missing data points in the census data. This report creates ten models that analyze the 15 different attributes included with the data in addition to adding two additional attributes (gross domestic product and Human Development Index for native countries). The models use a least squares solution which is validated by cross-validation. The results of this report did not find high variance in the error of each model but did make interesting discoveries regarding the relationship between the attributes indicating education and the native country of an individual. The overall conclusion made is that the United States population is significantly diverse and therefore it is difficult to accurately predict values such as years of education. The results of this report leave many avenues for future work including applying the models created to modern census data or adding additional census attributes (i.e. number of children or parental education levels) to see how these attributes affect the model.

Introduction

Every 10 years, the United States Census Bureau conducts a survey count of every resident within the United States. This information is critical to tracking the progress and demographic changes of the country throughout time. This data is used to make important decisions about the composition of the electoral members of the United States government and to delegate budgets for various organizations (United States Census Bureau). The Adult dataset from the UC Irvine Machine Learning Repository contains 1994 census data extracted by Barry Becker from the Census Bureau database. The data set contains 48,842 instances and 15 attributes: age, work class, final weight, education, education number, marital status, occupation, relationship, race, sex, capital gain, capital loss, hours worked per week, native country, and income level (above or below \$50,000).

A more detailed description for each attribute is as follows:

1. Age – age of the individual. This data is numerical, in whole numbers of years.
2. Work Class – sector the individual works in (government, private, self-employed, not working). The mode of this attribute was private, with nearly 75% of the entries.
3. Final Weight – a (continuous) number estimating the number of people in the United States population which matches the individual's demographics. This value considers the potential for sampling bias in the census and assigns a "weight" to the person based on their demographics. Final weight may be calculated at the federal, state, or local level and thus values in this category may not be consistent.
4. Education – label describing the highest level of education completed.

5. Education Number – arbitrary number assigned according to education label.
6. Marital Status – current marital status of the individual (single, married, previously married).
7. Occupation – a variety of labels of the general job title of the individual.
8. Relationship – a label of the relationship the individual has (wife, husband, unmarried, child, not in family).
9. Race – the race of the individual (white, Asian/Pacific Islander, Native American (“American Indian/Eskimo” in the data), black, other).
10. Sex – the sex of the individual (male, female).
11. Capital Gain – a continuous category which measures an individual’s capital gain. Note that this is different from income and might include income earned via stocks and investments.
12. Capital Loss – a continuous category which measures an individual’s capital loss.
13. Hours Worked per Week – a numerical category which details the number of hours per week an individual works.
14. Native Country – a label describing an individual’s native country. This attribute is highly skewed, with the mode of United States comprising of over 91% of the data.
15. Income Level – a Boolean (true/false) attribute which describes if the individual made a yearly income of greater than \$50,000.

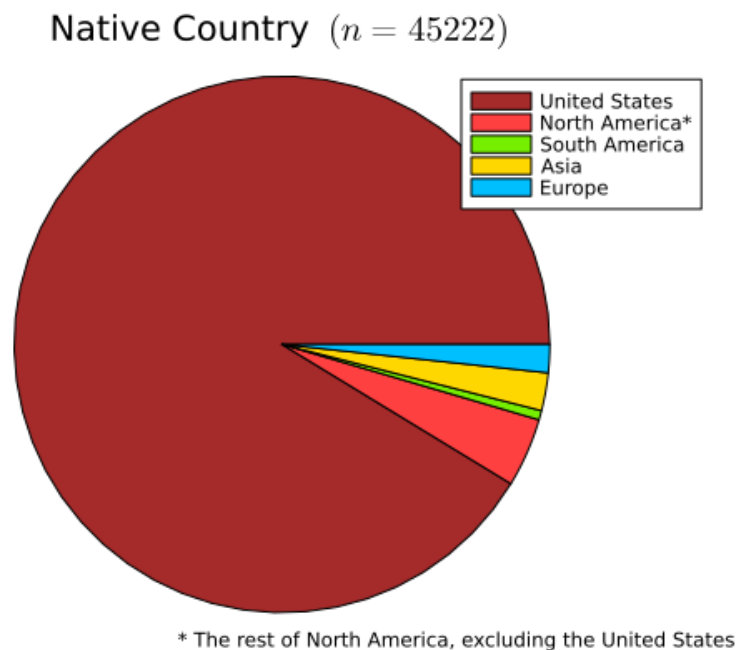
The purpose of this report is to apply knowledge of matrix methods and data science to predict the years of education of different individuals listed in the census data. The years of education an individual experiences can greatly vary and are influenced by a variety of factors. These factors can influence the experience and overall quality of life an individual endures. The ability to predict the years of education an individual in the census has undergone can point out key areas of improvement or success in the United States’ education system. Education is an important area of improvement for the United States because education is linked to economic growth. The growth of the United States economy ultimately improves the quality of life for United States’ citizens. Research suggests that better education contributes to an increase in community and long-term economic growth (Roser, Max, and Esteban Ortiz-Ospina).

While the dataset used for this report is nearly 30 years old, the methodologies described in this report could be used to analyze census data in the current decade. Being able to use the data of today to predict educational trends in the future could help the United States focus economic resources on the correct groups that need support to increase the years of education they receive on average. The intended audience for this report is individuals interested in data science and population analysis (i.e. demographers) for improvement of the human condition.

Methodology

Several columns required feature engineering before model generation could begin. Detailed documentation for all the feature-engineered columns is provided in the comments in the Jupyter Notebook code. All 3,620 rows that contained missing values were assigned the missing data type in Julia and then removed from the matrix. The resulting final matrix had 45,222 entries.

The attributes that were separated into more than one column or translated into numerical values include: work class, education, marital status, occupation, race, sex, and native country. The work class column was separated into three separate columns to represent private, self-employed, and government employees with an implicit fourth column for non-working individuals. The years of education column was created by assigning a value (in years) to each label in the education column. Then these values were consolidated into one numerical column. The marital status column was separated into two separate columns to represent single and married individuals with a third implicit column to represent previously married individuals. The occupation column was separated into four separate columns to represent engineering, business, technical, and non-degree workers with a fifth implicit column to represent government occupations. The race column was separated into four separate columns for white, Asian Pacific-Islander, Native American (“American Indian/Eskimo” in the data), and black with a fifth implicit column for individuals using the label “other”. The sex column was assigned a numerical Boolean value for female or male. The native country column was separated into three columns based on the continent where the country was located: North America, South America, or Asia. An implicit fourth column was added for Europe (other continents were not present in the data). An important note about this column is that it is extremely skewed toward North America because most entries have a native country of the United States. This may reduce the accuracy of predictions.



The attributes of final weight, capital gain, capital loss, education number, and relationship were eliminated from the data set for the purpose of this report. The decision for removal was based on a paper utilizing the same data set (Lemon, Chet, et al). Final weight was found to be inconsistently calculated between censuses, states, and counties so these values were determined to be unreliable for analysis (Weighting, United States Census Bureau). Capital gain and capital loss did not serve a purpose for prediction in this report as the columns did not appear

to contain significant values, and the distribution of the values was heavily skewed towards zero. The education number column appeared to follow some sort of pattern, but ultimately this was deemed unclear based on the documentation and did not correlate with years of education. Relationship appeared to also have an unclear meaning in the dataset documentation, so it was decided not to use that column since there is also a marital status column. The dataset will be treated as a sample of the United States population in 1994 since the final weight value is not being used.

Two additional attributes were added to the data including the Gross Domestic Product (GDP) (countryeconomy.com) and Human Development Index (HDI) (United Nations) for each native country. The GDP for each native country acts as a secondary economic indicator to supplement the lack of an exact income for each entry in the census data (each income is either above or below \$50,000). The HDI for each native country acts as a secondary measure of education for each entry, but also factors the country's GDP. Most entries have the native country as the United States, but observing patterns created by the HDI of native countries other than the United States may reveal insights into meaning behind the models produced.

The attributes provided in the data will be grouped into two categories: controllable and uncontrollable. Controllable attributes are assumed to include work class, occupation, marital status, income, and hours per week because each of these attributes are directly influenced by decisions made by an individual. Uncontrollable attributes are assumed to include native country, GDP, race, sex, and age because each of these attributes are not directly influenced by decisions made by an individual.

To start exploratory analysis, for each model the data was split into 20% testing/80% training data and a single least-squares solution predicting the target attribute was calculated, then the model's root mean square (RMS) error was calculated using testing data. For data visualization purposes a scatter plot was generated comparing the actual target attribute values versus the model's predicted values.

To validate these results a cross-validation approach was used for each of the models, where the data matrix was split into five randomly generated folds of 20% testing/80% training data. For each fold a least-squares solution predicting the target attribute was computed, and the RMS error was calculated. For data visualization purposes a scatter plot was generated comparing the actual target attribute values versus the model's predicted values, per model and fold. A total of 10 models were created. The final four models were used to predict the years of education and Human Development Index (HDI) columns using the controllable and uncontrollable attribute categories.

Data, Figures, and Tables

Discovery Models

The purpose of creating the first six models utilized in this analysis is to understand how the different attributes provided in the data are related to each other and can potentially be used

to predict years of education. The logic behind how to pair each set of attributes was decided based on arbitrary categories that eventually developed into the “controllable” and “uncontrollable” categories used in later models.

The models below utilize the following paired attributes to predict years of education:

1. Model 1: Age and Hours Worked Per Week.
2. Model 2: Work Class and Occupation.
3. Model 3: Marital Status and Native Country.
4. Model 4: Race and Sex.
5. Model 5: Gross Domestic Product and Income.
6. Model 6: Human Development Index and Income.

Discovery Models: Initial Least Squares

The least squares model is effective for predicting missing values. The accuracy of these values can be evaluated based on the root-mean-square (RMS) error. The ability to predict years of education of an individual could be a valuable asset to census data, so it was decided to attempt to predict this attribute in the discovery models. Additionally, visualizing the prediction accuracy for each attribute pair allows for easy identification of the most accurate pairing.

The models below resulted in the following RMS errors predicting years of education:

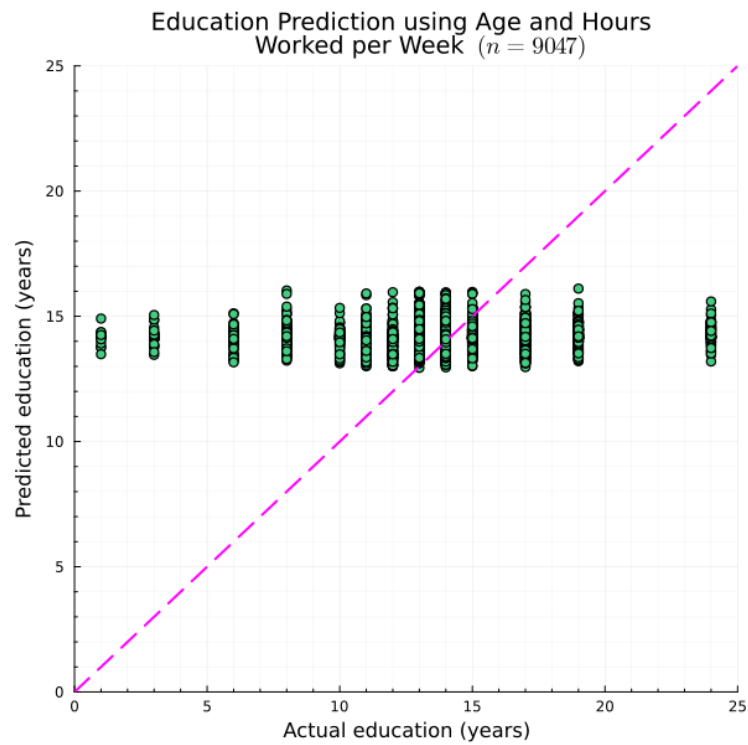
	Training RMS	Testing RMS
Model 1	2.6751650113399617	2.679998454219119
Model 2	2.502494041374889	2.512325126447655
Model 3	2.685490258184635	2.683340402831584
Model 4	2.68532773666033	2.6798834191143026
Model 5	2.53986984246196	2.5391713577388058
Model 6	2.5438250170376357	2.5457038692366583

Based on the results above, we can see that none of the models resulted in significantly different predictions of years of education. Even the models that utilized attributes that were added to the raw data (i.e. GDP and HDI for native country) did not result in large differences in prediction.

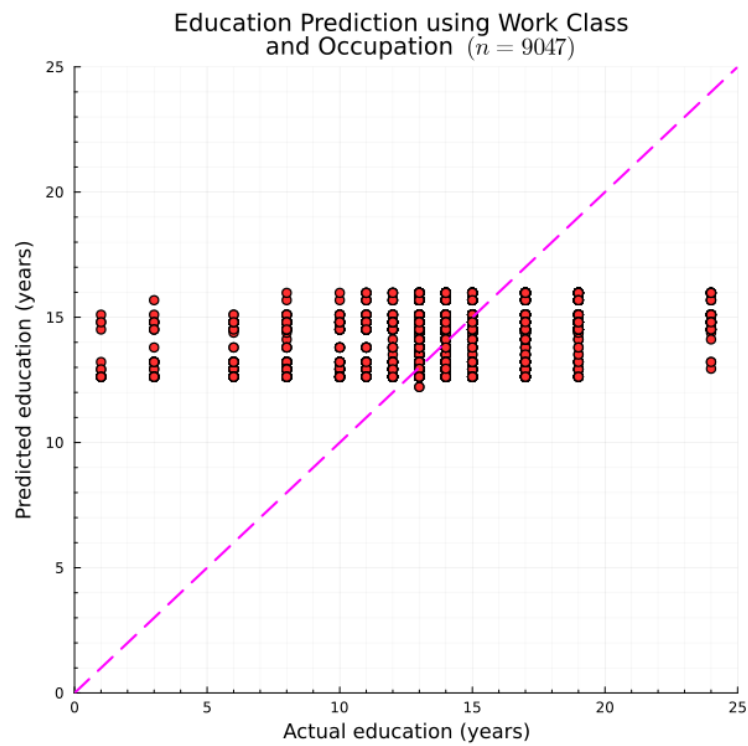
While the RMS errors show an average prediction error of 2 to 3 years, based on the plots below the model overestimates individuals' education levels for low values and underestimates

individuals' education levels for high values, indicated by the low "slope" of the plots. These observations gave justification to create the new categories of controllable and uncontrollable attributes to see if new models based on these categories would lower the RMS error of the predictions to be more significant.

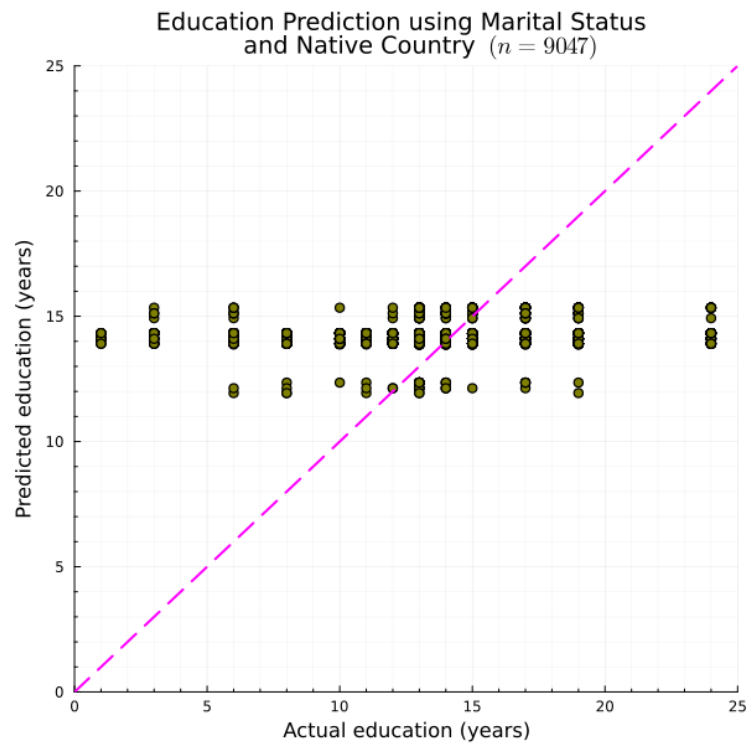
Least Squares Model 1: Education Prediction Using Age and Hours Worked Per Week



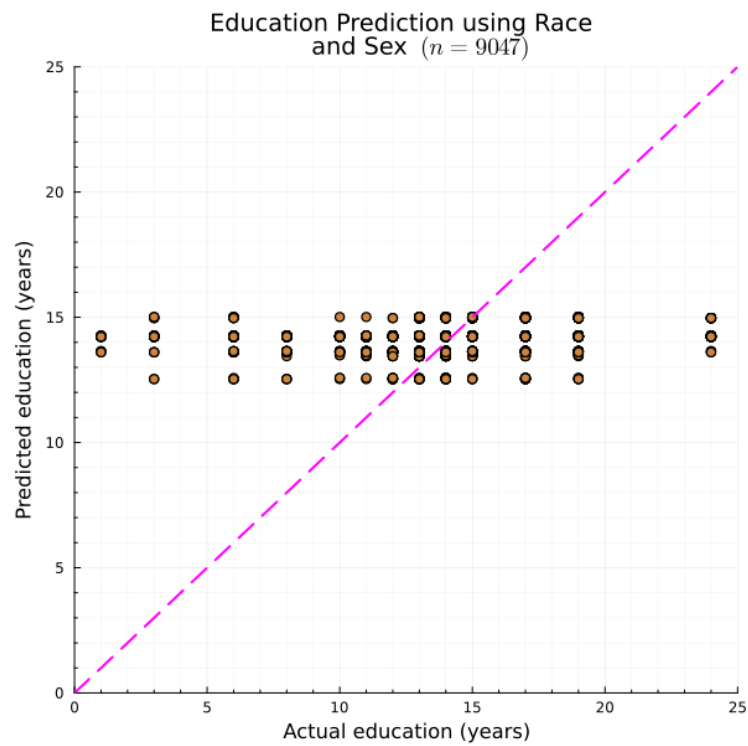
Least Squares Model 2: Education Prediction Using Work Class and Occupation



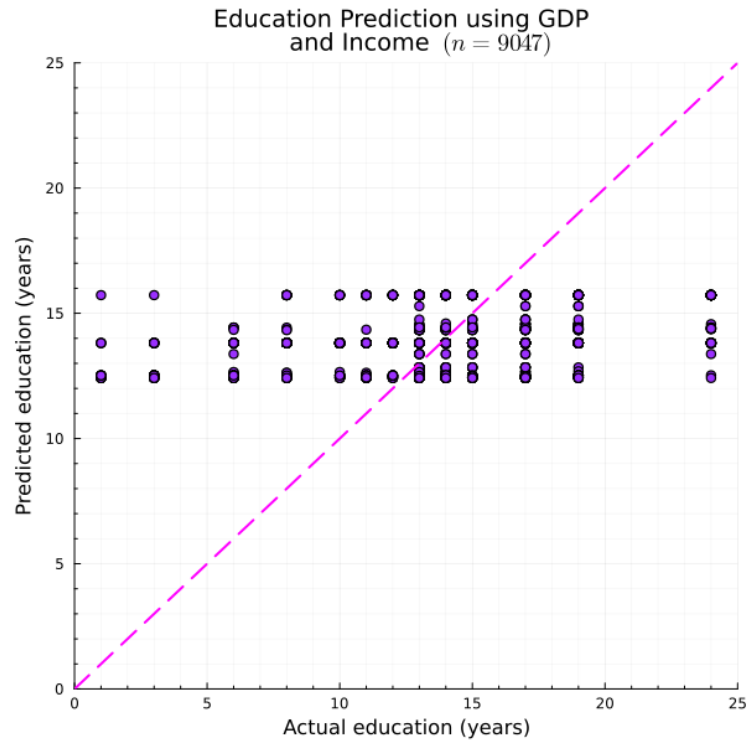
Least Squares Model 3: Education Prediction Using Marital Status and Native Country



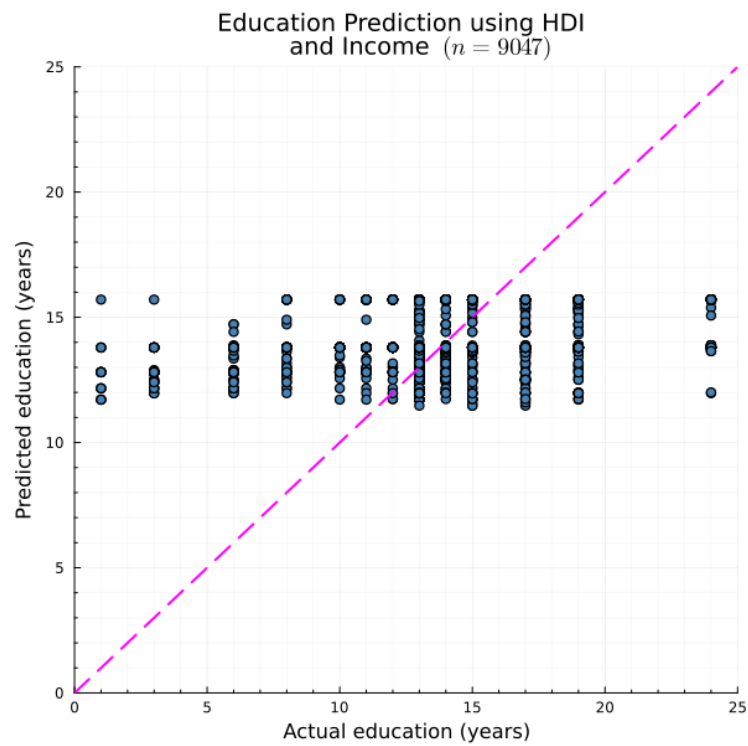
Least Squares Model 4: Education Prediction Using Race and Sex



Least Squares Model 5: Education Prediction Using Gross Domestic Product and Income



Least Squares Model 6: Education Predicting Using Human Development Index and Income



Discovery Models: Cross-Validation

The purpose of cross-validation is to separate data into folds and use multiple iterations to validate the model calculated previously. In this case, cross-validation will be used to validate the discovery models' solutions. A total of five folds were used in the cross-validation of these models. Only the plot from the first fold for each model is included below, but the other four folds can be generated using the Jupyter Notebook included with this report or by viewing the included PowerPoint for a slideshow of the different graphs in a graphics interchange format (GIF).

Below are the RMS error results from each cross-validation:

Model 1: Age and Hours Worked Per Week

Training RMS	2.67745	2.68094	2.67959	2.6779	2.66435
Testing RMS	2.67061	2.65658	2.66205	2.66884	2.72266

Model 2: Work Class and Occupation

Training RMS	2.50425	2.50096	2.51185	2.50032	2.50428
Testing RMS	2.5053	2.51863	2.47474	2.52122	2.50516

Model 3: Marital Status and Native Country

Training RMS	2.69031	2.69045	2.68677	2.67987	2.6775
Testing RMS	2.66411	2.66357	2.67825	2.7057	2.71531

Model 4: Race and Sex

Training RMS	2.69404	2.68734	2.68132	2.68417	2.67377
Testing RMS	2.64442	2.67151	2.69568	2.68441	2.7255

Model 5: Gross Domestic Product and Income

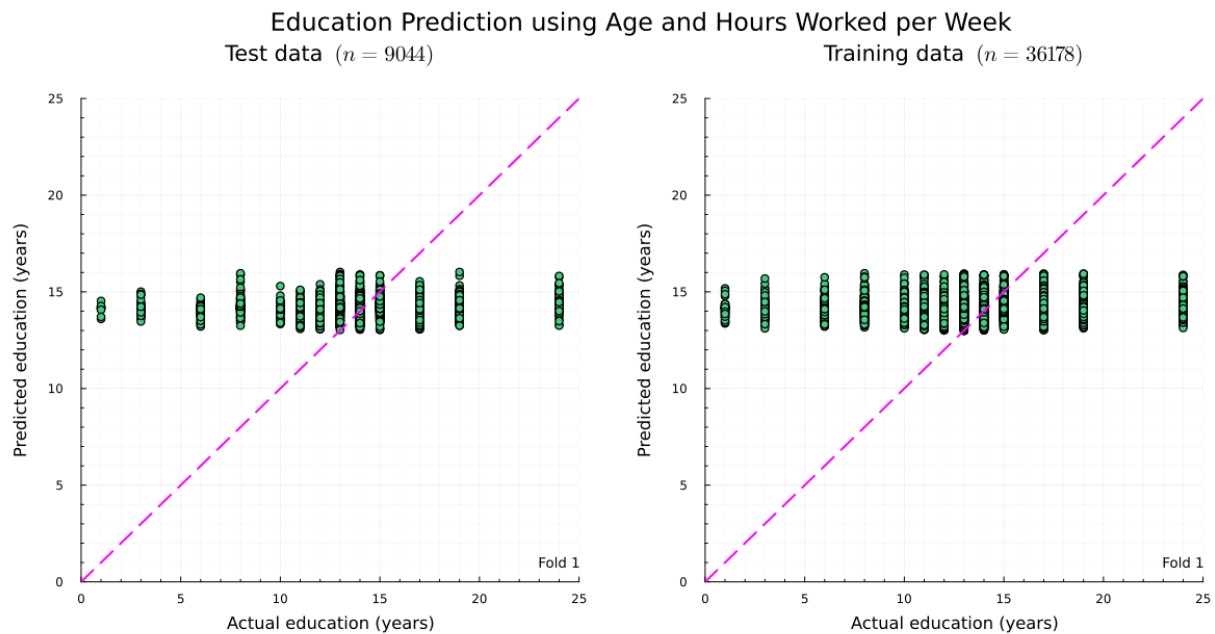
Training RMS	2.53674	2.54578	2.53636	2.52809	2.55153
Testing RMS	2.55167	2.51548	2.55322	2.58578	2.4921

Model 6: Human Development Index and Income

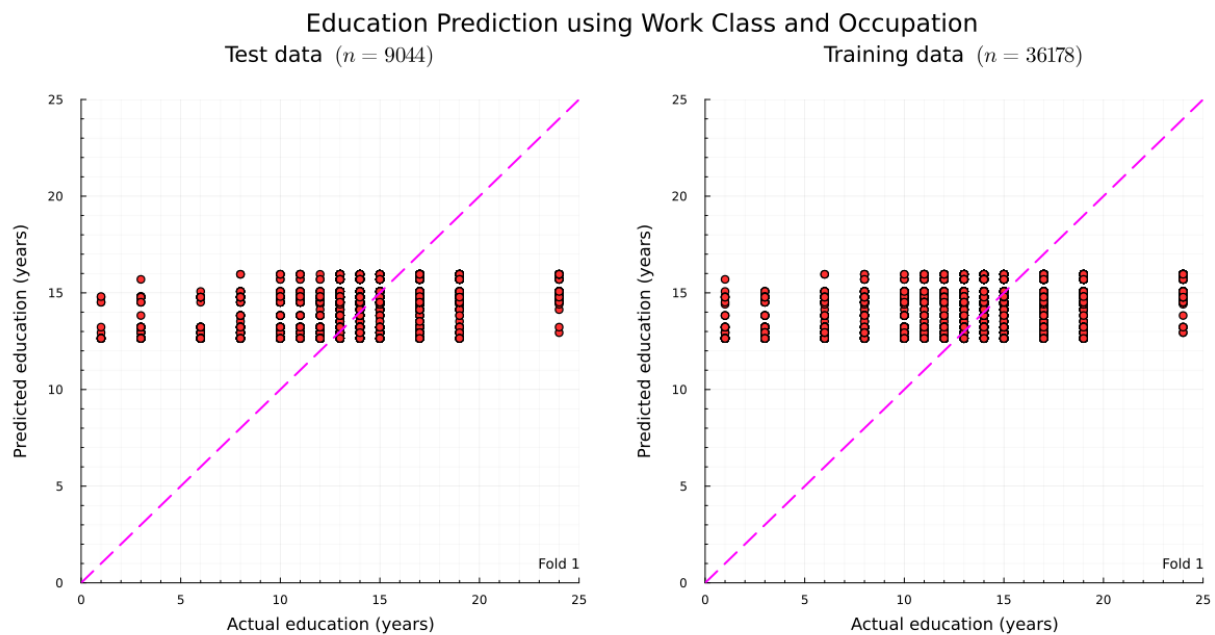
Training RMS	2.68376	2.68492	2.68126	2.68909	2.68584
Testing RMS	2.69063	2.68577	2.70024	2.66901	2.68212

Based on the results above, the conclusion can be made that the least squares models are valid because all the resulting RMS values are similar to each other for each given model. Since the RMS errors for testing data are similar to those of the training data, this indicates each of these models are not overfit to the training data, further increasing the models' validity.

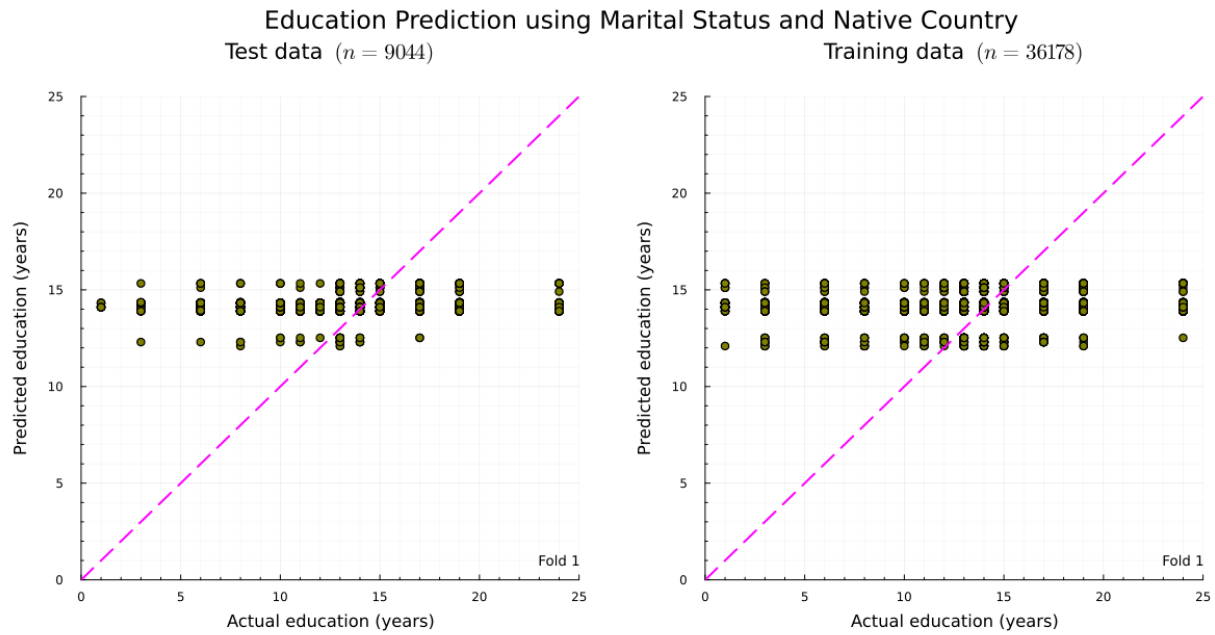
Cross-Validation Model 1: Education Prediction Using Age and Hours Worked Per Week



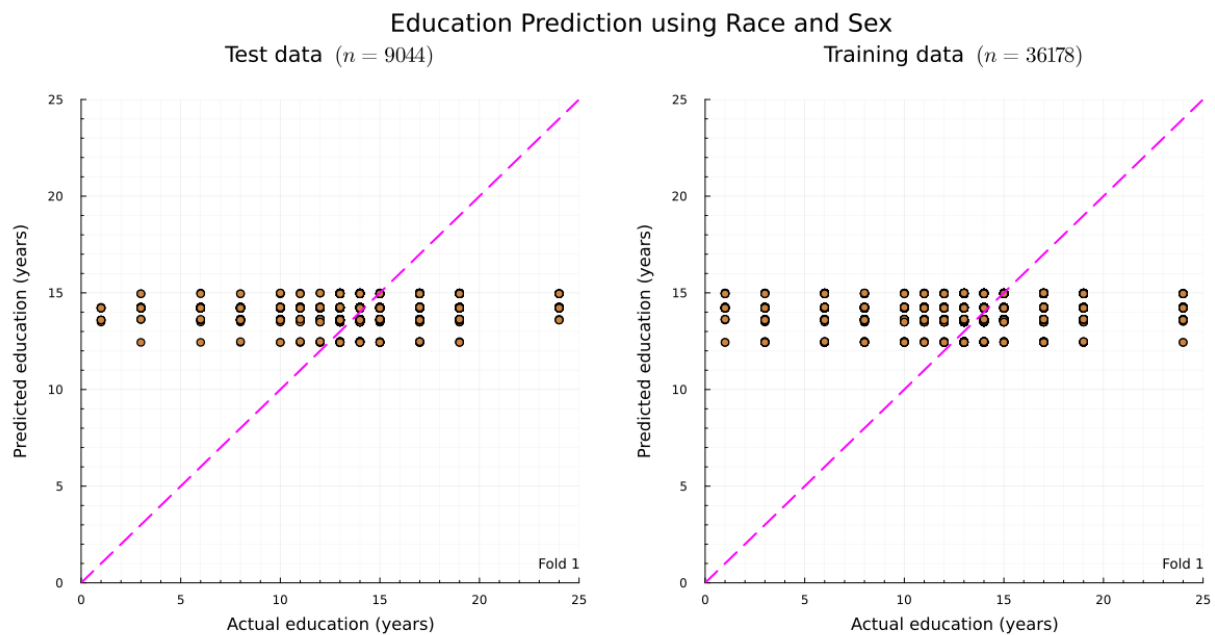
Cross-Validation Model 2: Education Prediction Using Work Class and Occupation



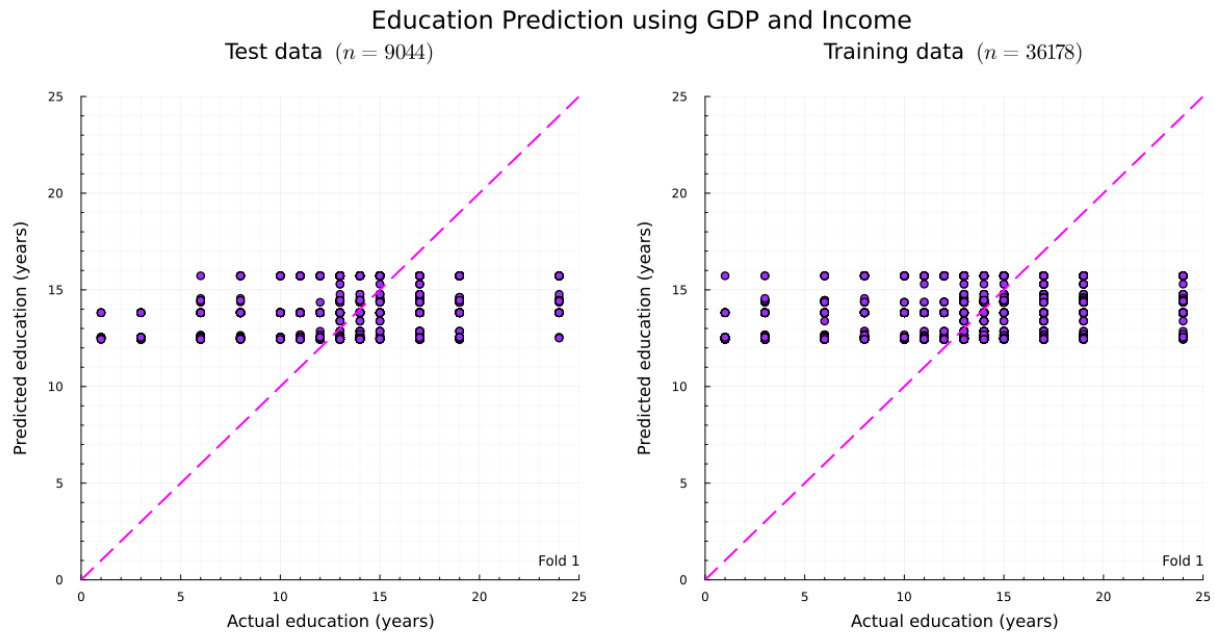
Cross-Validation Model 3: Education Prediction Using Marital Status and Native Country



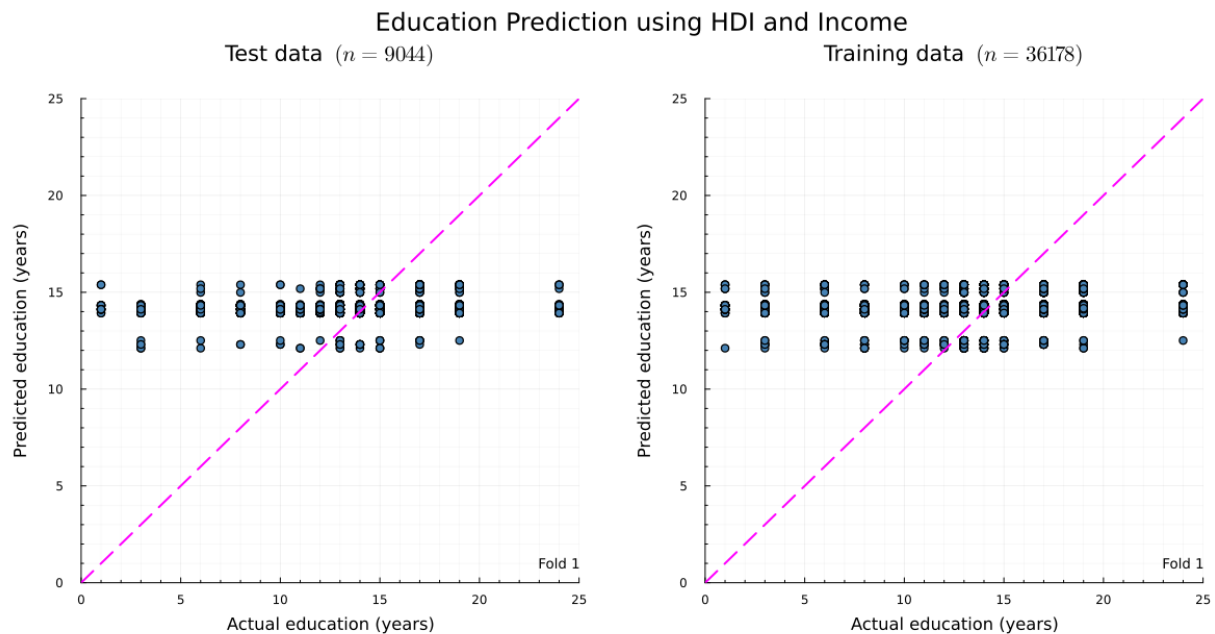
Cross-Validation Model 4: Education Prediction Using Race and Sex



Cross-Validation Model 5: Education Prediction Using Gross Domestic Product and Income



Cross-Validation Model 6: Education Predicting Using Human Development Index and Income



Predicting Years of Education with Controllable and Uncontrollable Attributes

Following the results of the Discovery Models, the attributes were re-grouped into the controllable and uncontrollable categories. These categories were used to make further predictions about the years of education using least squares solutions and then validated with a cross-validation test. The controllable category is defined by the attributes work class, occupation, marital status, income, and hours worked per week. The uncontrollable category is defined by the native country, gross domestic product, race, sex, and age attributes.

Predicting Years of Education: Initial Least Squares

Two new models were created using the controllable and uncontrollable categories to see if more accurate predictions could be made regarding the years of education. Initially a least-squares solution was calculated for these models, without cross-validation. The following \hat{x} values resulted from the models:

X-Hat Values Model 7: Education Prediction Using Uncontrollable Attributes

Attribute		X-Hat Value
Y-intercept		13.042
Native Country:	North America	-2.72529
	South America	-1.54074
	Asia	0.915771
Native GDP (in \$ billions)		0.000450635
Race:	White	0.531306
	Asian/Pacific Islander	0.877301
	Native American	-0.165588
	Black	-0.0172389
Sex:	Female	0.02444
Age:		0.00543944

X-Hat Values Model 8: Education Prediction Using Controllable Attributes

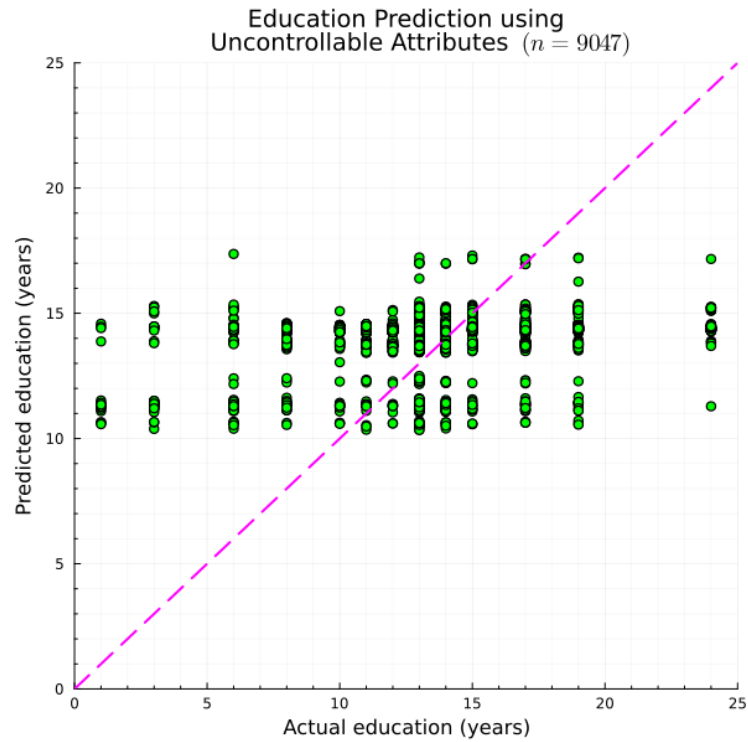
Attribute		X-Hat Value
Y-intercept		13.66
Work Class:	Private	0.0269376
	Self-employed	0.206359
	Government	1.12609
Occupation:	Engineering	0.35586
	Business	1.44653
	Technical	1.68272
	Non-degree	-0.123143
Marital Status:	Single	0.449082
	Married	-0.325552
Income		-1.71803
Hours Worked per Week		0.0151695

The following RMS errors resulted from the models:

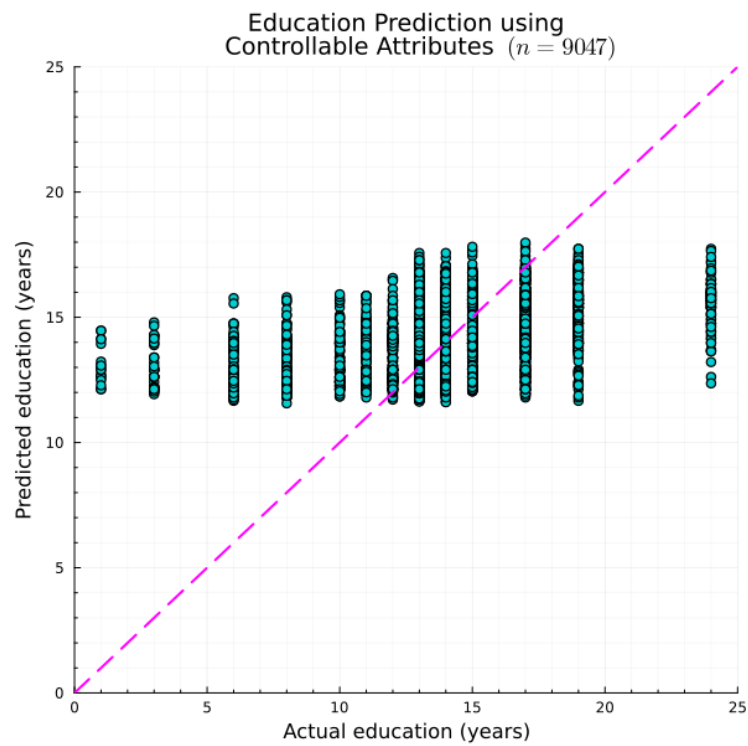
	Training RMS	Testing RMS
Uncontrollable (Model 7)	2.6137331343608308	2.6037659893776155
Controllable (Model 8)	2.398386831440762	2.4171925225529947

Based on the results above, the conclusion is that the controllable category results in the lowest RMS error of about 2.4 years. These results make sense considering that work class and occupation are included in the controllable category and resulted in the lowest RMS error out of all the pairs in the Discovery Models section (Model 2). This may suggest that factors relating to an individual's career (i.e. work class, occupation, income, hours worked per week) most accurately predict how many years of education they have completed.

Least Squares Model 7: Education Prediction Using Uncontrollable Attributes



Least Squares Model 8: Education Prediction Using Controllable Attributes



Predicting Years of Education: Cross-Validation

As explained in the previous cross-validation section, cross-validation is being used to validate the models' least-squares solutions to predict years of education. Only the plot from the first fold for each model is included below, but the other four folds can be generated using the Jupyter Notebook included with this report or by viewing the included PowerPoint for a slideshow of the different graphs in a graphics interchange format (GIF).

Below are the RMS error results from each cross-validation:

Model 7: Years of Education Predicted Using Uncontrollable Attributes

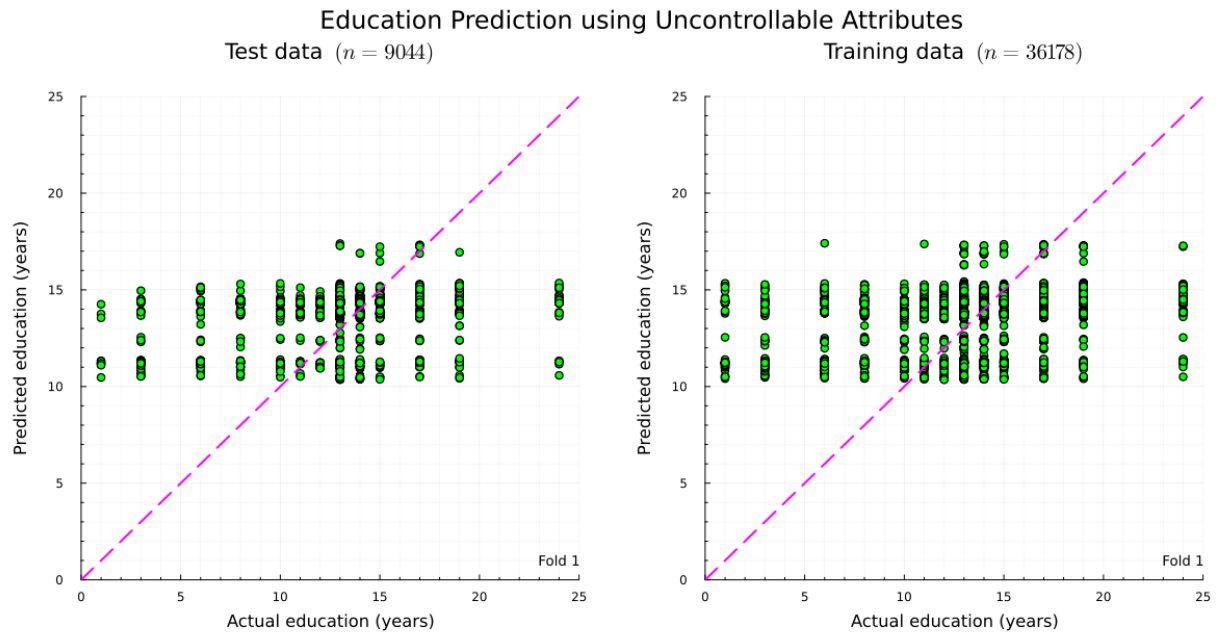
Training RMS	2.61364	2.61541	2.62344	2.59976	2.60386
Testing RMS	2.6025	2.59521	2.56259	2.65769	2.64146

Model 8: Years of Education Predicted Using Controllable Attributes

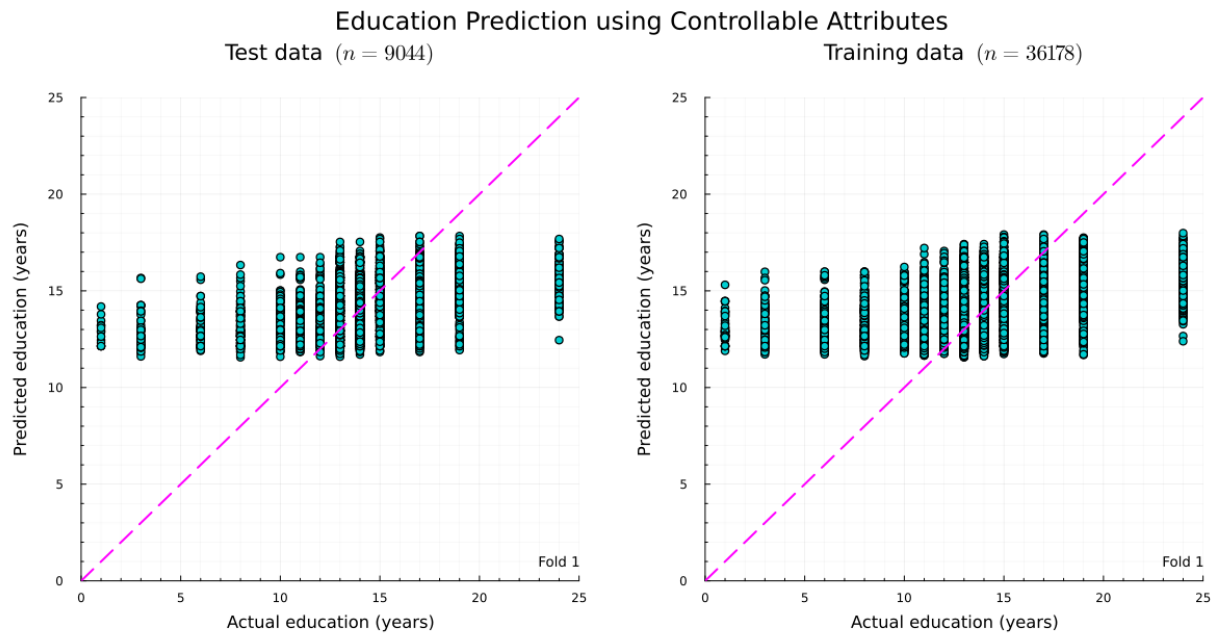
Training RMS	2.38612	2.41069	2.40545	2.39829	2.40919
Testing RMS	2.46485	2.36761	2.38852	2.41742	2.37347

Based on the results above, the conclusion can be made that the least squares models are valid because all the resulting RMS values are similar to each other for each given model. In addition, the testing RMS values were similar to the training RMS values, indicating that models 7 and 8 were not overfit.

Cross-Validation Model 7: Education Prediction Using Uncontrollable Attributes



Cross-Validation Model 8: Education Prediction Using Controllable Attributes



Predicting Human Development Index with Controllable and Uncontrollable Attributes

The Human Development Index (HDI) was added as a secondary indicator of education. This is because the HDI for a country is calculated according to aspects such as life expectancy from birth, expected years of schooling, and the country's GDP (United Nations). As a form of validation (in addition to the cross validation), predicting the HDI of an individual's native country from the sample of census data could help indicate the accuracy of the years of education predictions. If an individual's HDI is easier to predict than their years of education, this may suggest that the diversity of the United States population makes it difficult to predict values of every individual since they come from all over the globe.

Predicting Human Development Index: Initial Least Squares

As previously established, the least squares solution produces models that help predict unknown values. This test was used with the same controllable and uncontrollable attribute categories to observe how accurately predictions could be made for the HDI of each individual's native country. The following \hat{x} values resulted from the models:

X-Hat Values Model 9: Human Development Index Prediction Using Uncontrollable Attributes

Attribute		X-Hat Value
Y-intercept		0.734605
Native Country:	North America	-0.109331
	South America	-0.162632
	Asia	-0.176781
Native GDP (in \$ billions)		0.0000333779
Race:	White	0.00261042
	Asian/Pacific Islander	-0.0068521
	Native American	0.00351155
	Black	-0.00155952
Sex:	Female	0.000928655
Age:		0.000035595

X-Hat Values Model 10: Human Development Prediction Using Controllable Attributes

Attribute		X-Hat Value
Y-intercept		0.861265
Work Class:	Private	-0.00947663
	Self-employed	-0.00360414
	Government	-0.00121374
Occupation:	Engineering	0.00499996
	Business	0.015298
	Technical	0.00959016
	Non-degree	0.00565312
Marital Status:	Single	-0.005171
	Married	-0.0111575
Income		-0.00746344
Hours Worked per Week		0.0000600894

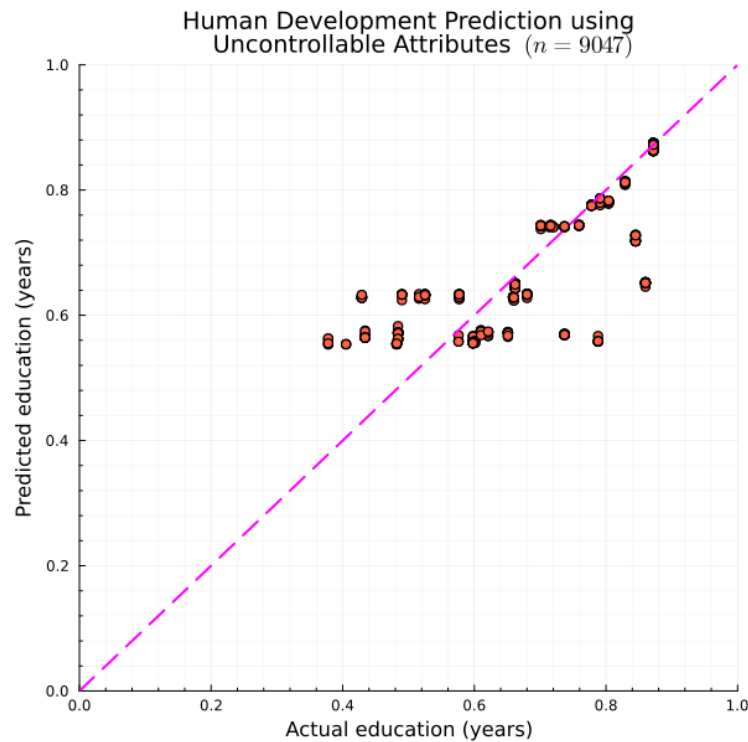
The following RMS errors resulted from the models:

	Training RMS	Testing RMS
Uncontrollable (Model 9)	0.024364083939679447	0.024681003356696894
Controllable (Model 10)	0.07015984984327628	0.06857875224288963

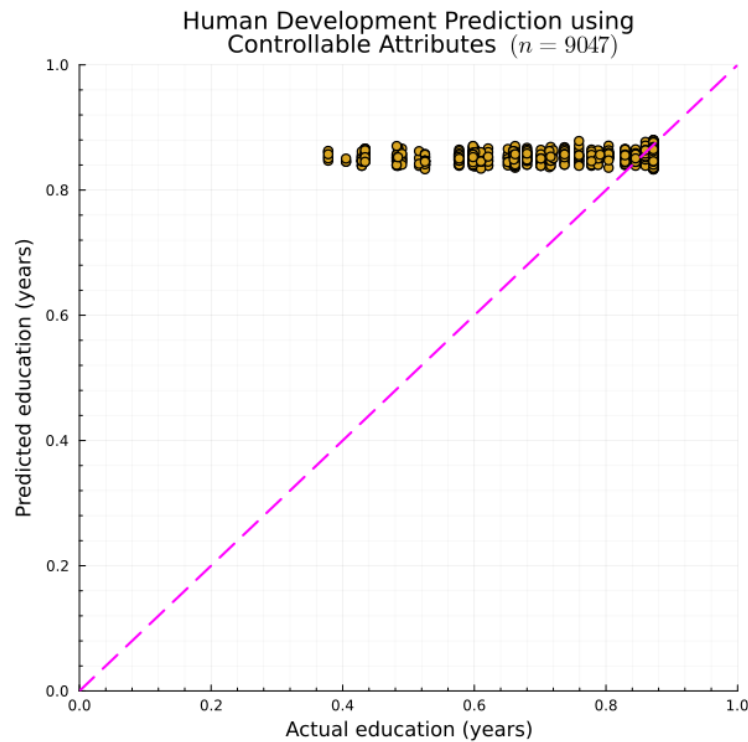
Based on the results above, the conclusion is that the prediction for the HDI based on uncontrollable attributes was significantly better than the prediction based on controllable attributes. While both categories were able to predict the HDI very accurately (less than 10% off in all cases), the model concerning uncontrollable attributes performed much better than the model concerning controllable attributes.

The very small slope of the data in Model 10's graph indicates while the RMS error is low (about 7%), the model seems to predict the sample's HDI as around that of the United States' HDI, 0.872, likely because the United States accounts for most of the data in the native country attribute. This heavily skews the HDI predictions towards the United States' HDI value. This phenomenon can be seen in the Model 10 \hat{x} table above: the \hat{x} values of the different attributes are very small, and the y-intercept, 0.861, is very close to the United States' true HDI, 0.872. This allows for little variation in HDI prediction from the United States' HDI value.

Least Squares Model 9: Human Development Index Prediction Using Uncontrollable Attributes



Least Squares Model 10: Human Development Index Prediction Using Controllable Attributes



Predicting Human Development Index: Cross-Validation

As explained in the previous cross-validation section, cross-validation is being used to validate the least squares solutions created to predict the HDI. Only the plot from the first fold for each model is included below, but the other four folds can be generated using the Jupyter Notebook included with this report or by viewing the included PowerPoint for a slideshow of the different graphs in a graphics interchange format (GIF).

Below are the RMS error results from each cross-validation:

Model 9: HDI Predicted Using Uncontrollable Attributes

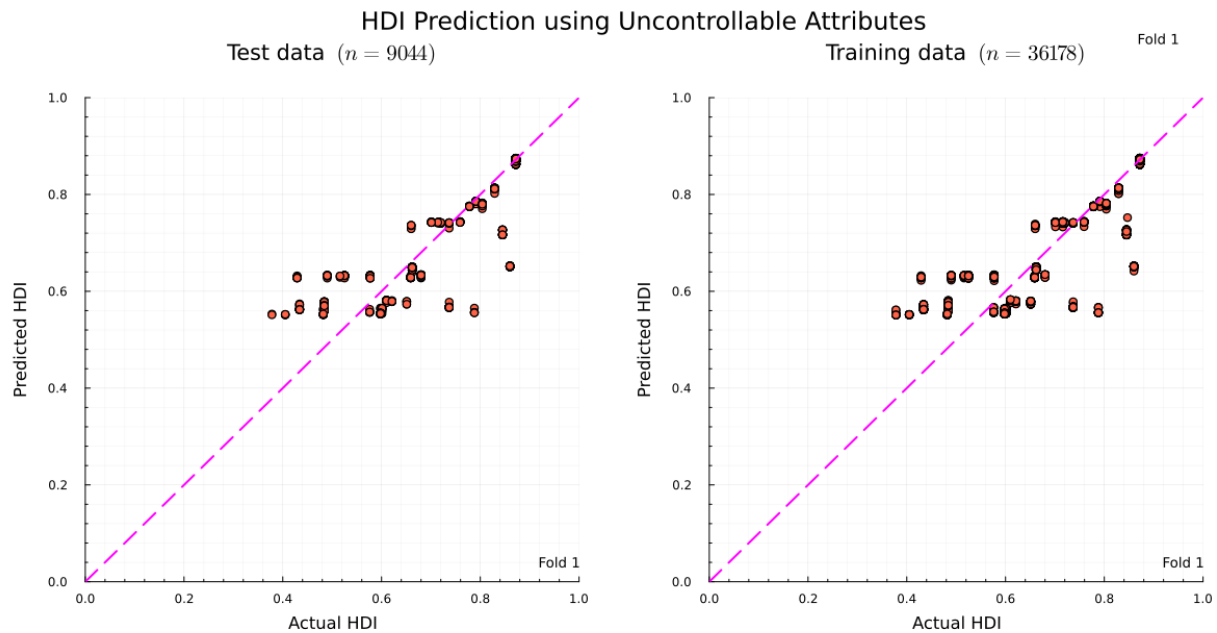
Training RMS	0.0245069	0.0245038	0.0240429	0.0245159	0.0245856
Testing RMS	0.0241563	0.024161	0.0259486	0.0241315	0.0238383

Model 10: HDI Predicted Using Controllable Attributes

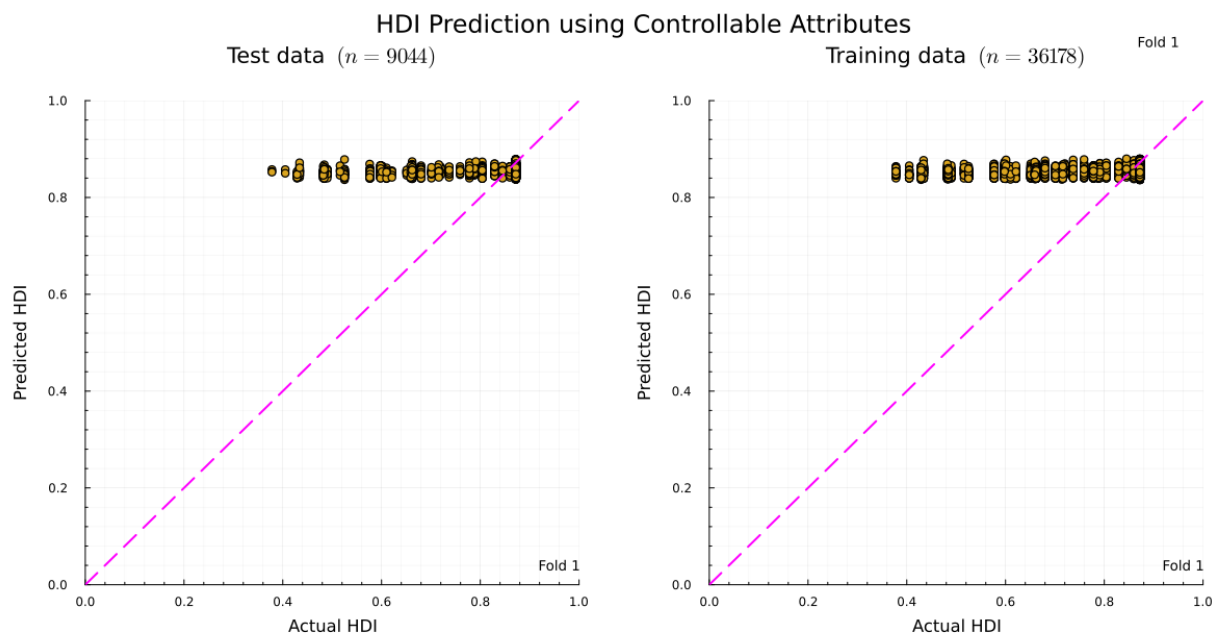
Training RMS	0.0692897	0.0698806	0.0701464	0.0698257	0.0700713
Testing RMS	0.0720302	0.0697094	0.0686374	0.069935	0.0689464

Based on the results above, the conclusion can be made that the least squares models are valid because all the resulting RMS values are similar to each other for each given model.

Cross-Validation Model 9: Human Development Index Prediction Using Uncontrollable Attributes



Cross-Validation Model 10: Human Development Index Prediction Using Controllable Attributes



Results and Future Work

Based on the models produced, the results of the analysis were not necessarily significant. However, the process of refining the groups of attributes into the controllable and uncontrollable categories significantly reduced prediction error. These categories could be further refined or used to generate different types of models that could produce more significant results. Overall, the Discovery Models revealed that splitting up the attributes into controllable and uncontrollable categories reduced the prediction error when predicting years of education.

One conclusion that can be made based on the analysis in this report is that the diversity of the United States population makes it difficult to accurately predict exactly how many years of education an individual has completed. This is further backed by the accuracy in predicting the Human Development Index (HDI) of the individuals' native country even when directly related attributes such as native country was not used in the model to make the prediction.

Another interesting observation is that native country was still a significant predictor of HDI even though countries were identified by continent rather than exact country. This suggests that further work on this project could involve re-analyzing the model to consider each individual country rather than each continent. This could potentially increase the accuracy of the prediction.

Additionally, since most of the entries in the native country column were the United States, this could have reduced the accuracy of models involving native country. This homogenous trend could also have interesting implications for the results of the analysis. What if the heritage of individuals could be used instead of the native country? Would this value allow an analysis of the impact culture has on education?

Future work on this project could follow many different routes. A few ideas include:

- Can we predict where an individual is from based on their census data? This could be significant because of the accuracy in predicting the HDI value for the native country based on controllable attributes.
- Can we apply the models created in this report to 2020 census data? What will the changes in error potentially imply about how the United States population has changed over the last 30 years? Are there new categories of attributes that could be evaluated?
- How would the error or models change when factoring in attributes not included in this sample of census data? Important attributes such as number of children, state of residency, or parental educational levels were not included in this data. These could potentially have a large impact on education predictions. Additionally, having more specific location data could allow future work to utilize the per capita GDP instead of the raw GDP for each country.
- How could a more inclusive data set influence the models in this report? Many attributes in this census data sample lack sufficient categories to describe a global population. Adding more descriptive options to attributes such as native country, race, or even occupation could provide more insight into the people behind the data.

References

- Becker, Barry. *Adult*. UCI Machine Learning Repository, 1996. *DOI.org (Datacite)*, <https://doi.org/10.24432/C5XW20>.
- GDP - Gross Domestic Product 1994 | Countryeconomy.Com*. <https://countryeconomy.com/gdp?year=1994>. Accessed 17 Apr. 2023.
- Home · Plots*. <https://docs.juliaplots.org/stable/>. Accessed 24 Apr. 2023.
- Introduction · DataFrames.Jl*. <https://dataframes.juliadata.org/stable/>. Accessed 24 Apr. 2023.
- Lemon, Chet, et al. *Predicting If Income Exceeds \$50,000 per Year Based on 1994 US Census Data with Simple Classification Techniques*. <https://cseweb.ucsd.edu/classes/sp15/cse190-c/reports/sp15/048.pdf>.
- Nations, United. *Documentation and Downloads*. United Nations. *hdr.undp.org*, <https://hdr.undp.org/data-center/documentation-and-downloads>. Accessed 17 Apr. 2023.
- Nations, United. *Human Development Index*. United Nations. *hdr.undp.org*, <https://hdr.undp.org/data-center/human-development-index>. Accessed 23 Apr. 2023.
- Our Censuses*. United States Census Bureau, *Census.Gov*, <https://www.census.gov/programs-surveys/censuses.html>. Accessed 17 Apr. 2023.
- Roser, Max, and Esteban Ortiz-Ospina. “Global Education.” *Our World in Data*, Aug. 2016. *ourworldindata.org*, <https://ourworldindata.org/global-education>.
- Subramanian, S. V., et al. “Multilevel Perspectives on Modeling Census Data.” *Environment and Planning A: Economy and Space*, vol. 33, no. 3, Mar. 2001, pp. 399–417. *DOI.org (Crossref)*, <https://doi.org/10.1068/a3357>.
- Weighting*. United States Census Bureau, *Census.Gov*, <https://www.census.gov/programs-surveys/sipp/methodology/weighting.html>. Accessed 1 May 2023.