One for the Money: Understanding How Demographics Impacts Personal Income in American Counties

Caleb Pena, cpea@andrew.cmu.edu

October 18, 2021

Abstract

This paper analyzes the effect of demographic factors such as age, education, and crime rates on personal income. We use data from a 1990 report from the University of Virginia's Geospatial and Statistical Data Center. We fit several multiple regression models using advanced variable selection techniques and compared the results. We found that crime rates are not a significant predictor of income after controlling for population. However, our analysis is limited by the fact the data only comes from relatively urban counties. Further study is needed to understand how these relationships play out in rural settings

Introduction

Identifying and interpreting key indicators of economic health is an important task for policy makers. Personal income is an especially important metric since it serves as a useful proxy for standard of living. Understanding what demographic factors influence per capita income gives economists and other analysts a better idea of what kind of policy changes can help improve quality of life.

In this paper we will analyze historical data from 1990 to better understand how variables such as poverty, unemployment, or education impact income level. For this exercise we have limited our toolset to methods discussed in Sheather (2009). In particular, we have been asked to focus our research by answering the following questions:

- Which pairs of variables are clearly related and what is the nature of those relationships?
- Is a county's crime rate associated with its per capita income? Does this association vary regionally?
- What is the best model for predicting per capita income?
- Has missing data skewed the dataset and our subsequent analysis?

Data

We studied these questions using data collected by the Geospatial and Statistical Data Center of the University of Virginia as cited in Kutner et. al. (2013). The dataset includes 3 categorical and 13 numeric variables measuring demographic features of 440 of the most populous counties in the US. The data only tracks counties with at least 100,000 residents and, due to incomplete data, a handful of counties that do meet this requirement are still omitted. Because of this requirement, three states are unrepresented. The two states with the lowest population density (citation needed), Alaska and Wyoming, do not make the cut. Suprisingly, Iowa is also missing. Although it is near the middle of state population rankings, it spreads its population fairly evenly across 99 counties. It also has no major cities so none of its counties meet the strict requirement of our dataset. So long as our analysis is contextualised as pertaining to urban not rural settings this should not significantly affect our conclusions.

Table 1 provides a brief summary of what kind of information was tracked.

Table 1	: D	escription	of	Data
---------	-----	------------	----	------

Variable Name	Description
county	Name of county
state	Name of state
land_area	Area in square miles
pop	Estimated 1990 population
pop_18_34	Percent of population aged 18-34
pop_65_plus	Percent of population aged 65+
doctors	Number of active physicians in 1990
hosp_beds	Number of hospital beds
crimes	Number of serious crimes
pct_hs_grad	Percent of adult population who graduated high school
pct_bach_deg	Percent of adult population with bachelor's degrees
pct_below_pov	Percent of population living below the poverty line
pct_unemp	Percent of population that is unemployed
per_cap_income	Per capita income in dollars
tot_income	Total personal income in dollars
region	Geographic region (S, W, NE, NC)

Figure 1 shows the pairwise Pearson correlations between the variables. Most of these relationships are unsurprising. Total income has a positive relationship with total crimes, total hospital beds, and total doctors since these measures are primarily determined by total population. In addition, counties that are more educated score higher in both percentage of high school graduates and percent with a bachelor's degree. And higher per capita income is associated with a lower unemployment and poverty rate.



 $\mathbf{2}$

A less obvious result is the relationship between age and education. Commentators and pundits commonly talk about how jobs that previous generations could hold down with just a high school diploma now require a college degree. Our data suggests this trend may have impacted how the younger generation has valued pursuing higher education. There is a strong linear relationship between the percentage of the population that is between 18 and 34 and the percent with a bachelor's degree. The opposite relationship holds for the percent of the population greater than 65. These relationships are plotted in Figure 2. Interestingly, despite the strong relationship between education and economic health, neither of the age variables show are closely linearly related to per capita income or poverty. Further graphical exploration of the data can be found in the technical appendix (Part 1).



Relationship between age and higher education levels

Methods

We broke our analysis into two main goals. The first was to identify whether or not crime rates have a direct impact on per capita income in different regions of the country. To answer this question we built a multiple linear model with and without interaction terms between the categorical region variable and crime rate. This was attempted using both the total number of crimes and the per capita crime rate. We also visually inspected these regression lines using the R package ggplot2 (citation needed).

Next, we incorporated the remaining variables into our analysis and compared multiple models to find the one that best fit the shape of the data. We used stepwise regression and lasso regularization to select the most predictive variables. Then we compared these models using a combination of diagnostic plots, measures of fit, and tools to identify multi-colinearity. The code for this section may be found in Technical Appendix part 2.

Results

Modeling using just crime rate

As is often the case with variables closely related to population, both the total crimes and per capita income variables are heavily right-skewed. To address this we took their logarithms and fit the following model:

$$log(PerCapitaIncome) = log(Crimes) + Region + error$$
(1)

Model 1 assumes that a change in the crime rate has the same impact regardless of region. We relaxed this assumption model 2 by adding interaction terms:

$$log(PerCapitaIncome) = log(Crimes) * Region + error$$
(2)

The resulting regression lines can be seen in figure 3. Modeling regional variation does appear to be important but the slopes are all very similar. The model with the interaction terms does not appear to add any meaninful information to the model. This was confirmed when we conducted a series of t-tests on the coefficients of the interaction terms. None of them had values that were significantly different from zero. Thus, we maintain that model 1 better captures the relationship between crime rate and per capita income.



(Figure 3)

The estimated slope for log(crimes) in model 1 is 0.067. Since this coefficient is small, Sheather tells us we can interpret it using percentages (Sheather 2009). That is, a 1% change in crimes correponds with a roughly 0.067% change in per capita income. Even though this impact is statistically significant, it is still small in an absolute sense. More importantly, the direction of the association goes against our intution as well as established research on crime rates and income levels. This is likely because model 1 does not take into account population. We will see below that once we control for this omitted variable crime rates no longer predict income well.

Modeling with all variables

(Note for reviewers: Didn't get the code quite how I want it yet. Hopefully this still gives you an idea of the kind of content that will eventually be here)

Next we tried to find the best model including a wider subset of the variables. Most the variables in the data are right-skewed. The amount of skewing varies so using Box-Cox transformations would likely results in a more formally valid model but for the sake of interpretability we chose to apply simple log transformations to each. The high school graduation variable shows left-skewing so we applied a power 2 transformation.

Two variable selection methods were used. Lasso regression and bidirectional stepwise regression begining with the saturated model. The results from these methods are displayed in table 2 below.

"This will be Table 2"

[1] "This will be Table 2"

Discussion

Now that we have discussed the results of our analysis, we are prepared to answer the questions set out at the beginning. We learned through studying correlation that most economic indicators, including poverty level and unemployment, relate closely to personal income. We also learned that crime rates have a very weak association with income level contrary to common perceptions. Lasso produced the most powerful model as measured by AIC. And due to the limited scope of our research, the missing data has not significantly impacted our analysis.

However, before generalizing our results the reader should remember two things. First, the data our study relies on is nearly 30 years old. Even in 1990 we identified a changing environment with respect to the essential education variables. This effect might be more significant today and it would be helpful to see more recent data. Secondly, urbanization is not well accounted for our in research. Although it is trivial to add a population density to our model, this would have limited value until data is collected on the nearly 2600 rural counties not included in the dataset.

References

Kutner, M. H., Nachtsheim, C. J., Neter, J., & Li, W. (2013). Applied Linear Statistical Models. McGraw-Hill Education (India) Private Limited.

Sheather, S. J. (2009). A modern approach to regression with R. Springer.

Technical Appendix

a) Data description.

Make a table or tables showing appropriate summary statistics for each variable in the data set. Note that summary statistics for continuous variables will be different from the summary statistics for categorical variables.

Answer:

The county column can be omitted since it is (when combined with state) a key column with a unique value for every row in the data. The other two categorical variables (state and region) have their frequencies displayed in tables 1 and 2. We also include a series of numerical summaries of the data in table 3.

count(cdi, state) %>% arrange(desc(n)) %>% head(10) %>%
knitr::kable(caption = "Number of Counties in the Top 10 States")

rabio 2, realiser of countries in the rop to state	Table 2:	Number	of	Counties	in	the	Top	10	States
--	----------	--------	----	----------	----	-----	-----	----	--------

state	n
CA	34
FL	29
PA	29
TX	28
OH	24
NY	22
MI	18
NC	18
NJ	18
IL	17

```
count(cdi, region) %>% arrange(desc(n)) %>%
knitr::kable(caption = "Number of Counties in Each Region")
```

Table 3: Number of Counties in Each	Region
-------------------------------------	--------

region	n
S	152
NC	108
NE	103
W	77

 Table 4: Numerical Summaries

variable	mean	median	sd	min	max
crimes	27111.62	11820.50	58237.51	563.0	688936.0
doctors	988.00	401.00	1789.75	39.0	23677.0
hosp_beds	1458.63	755.00	2289.13	92.0	27700.0
land_area	1041.41	656.50	1549.92	15.0	20062.0
pct_bach_deg	21.08	19.70	7.65	8.1	52.3
pct_below_pov	8.72	7.90	4.66	1.4	36.3
pct_hs_grad	77.56	77.70	7.02	46.6	92.9
pct_unemp	6.60	6.20	2.34	2.2	21.3
per_cap_income	18561.48	17759.00	4059.19	8899.0	37541.0
рор	393010.92	217280.50	601987.02	100043.0	8863164.0
pop_18_34	28.57	28.10	4.19	16.4	49.7
pop_65_plus	12.17	11.75	3.99	3.0	33.8
tot_income	7869.27	3857.00	12884.32	1141.0	184230.0

Indicate where (in which variables) there is missing data (NA's), if any, how much there is (in each variable) and why it might be there.

Answer:

The code below shows that the dataset has the same number of rows before and after removing rows with missing data. From this we can infer there are no NAs in the dataset.

nrow(cdi)

[1] 440

nrow(na.omit(cdi))

[1] 440

Make some appropriate descriptive EDA plots to illustrate any important features of the variables or possible important relationships among them.

Answer:

We begin by examining the marginal distributions of the numeric variables. Figure 1 shows significant right skewing on most of the variables. This is unsurprising since we anticipate many of the "count" variables like the numeber of crimes or doctors will be closely correlated the right-skewed population variable. The percentage of high school graduates is the only column that appears left skewed.

```
cdi %>%
select(where(is.numeric)) %>%
pivot_longer(cols = everything(), names_to = "variable", values_to = "value") %>%
ggplot() +
geom_histogram(aes(value)) +
labs(x = "", y = "", title = "Histograms of Important Variables", caption = "(Figure 1)") +
theme_bw() +
theme(axis.text.x = element_text(angle = 90)) +
facet_wrap(vars(variable), scales = "free")
```



```
(Figure 1)
```

Figure 2 confirms our suspicions that population, total income, number of crimes, doctors, and hospital beds are all closely positively correlated. We also learn that the education variables tend to be negatively correlated with the economic health variables. That is, a county with better educated residents tends to perform better economically.

```
cdi %>%
select(where(is.numeric)) %>%
cor() %>%
ggcorrplot::ggcorrplot(type = "lower", lab = T, digits = 1) +
labs(caption = "(Figure 2)")
```



Most of these relationships are unsurprising. What is less intuitive is the relationship between counties with high population between 18 and 34 and the percentage of residents with bachelors degrees. These variables have a pearson correlation of 0.46. One possible takeaway is that younger generations are expected to have a higher education level than their predcessors. The plot in Figure 3 displays this relationship graphically.

```
cdi %>%
ggplot() +
geom_point(aes(pop_18_34, pct_bach_deg)) +
labs(x = "Population Between 18 and 34", y = "% with a Bachelors Degree",
    title = "Relationship between younger populations and higher education levels",
    caption = "(Figure 3)") +
theme_bw()
```



b) Build a regression model that predicts per-capita income from crime rate and region of the country. Should there be any interactions in the model? What does your model say about the relationship between per-capita income and crime rate? Do your answers change, depending on whether you use number of crimes, or "per-capita crime" = (number of crimes)/(population) as a crime rate measure? If so, which one best answers the question? Why? Show the fitted model results and explain your answer to these questions in terms of those results.

Answer:

As we saw in part A, crimes and per_cap_income are right-skewed so we applied a log transform to both and fit a linear model with region as a potential confounder. Figures 4 and 5 show what these regressions would look like with and without interaction effects. In the first figure, we can see the lines are mostly parallel. The NC region has the most distinct slope but this is possibly caused by the influential points off to the left.

```
mod_wo_interactions <- lm(log(per_cap_income) ~ log(crimes) + region, data = cdi)
mod_w_interactions <- lm(log(per_cap_income) ~ log(crimes)*region, data = cdi)</pre>
```

```
cdi %>%
```

```
## `geom_smooth()` using formula 'y ~ x'
```





The summary output below confirms that the difference in slopes is negligible. None of the interaction terms are significant and the p-value from the partial F-test is 0.5266 much higher than $\alpha = 0.05$. Thus, we are unable to detect a significant improvement to the model by including interaction terms.

```
summary(mod_w_interactions)
```

```
##
## Call:
## lm(formula = log(per_cap_income) ~ log(crimes) * region, data = cdi)
##
## Residuals:
##
       Min
                  1Q
                      Median
                                    ЗQ
                                            Max
##
   -0.68552 -0.10418 -0.01444 0.08302
                                       0.79755
##
## Coefficients:
##
                        Estimate Std. Error t value Pr(>|t|)
## (Intercept)
                         9.33677
                                    0.14579 64.044 < 2e-16 ***
## log(crimes)
                         0.05064
                                    0.01566
                                              3.233 0.00132 **
## regionNE
                        -0.18407
                                    0.21515 -0.856 0.39272
## regionS
                        -0.19717
                                    0.21211 -0.930 0.35312
## regionW
                        -0.31439
                                    0.24465 -1.285
                                                    0.19947
                                    0.02311
## log(crimes):regionNE 0.03122
                                                    0.17749
                                              1.351
                                    0.02228
## log(crimes):regionS
                                                     0.58696
                         0.01211
                                              0.544
## log(crimes):regionW
                         0.02727
                                    0.02523
                                              1.081 0.28028
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1855 on 432 degrees of freedom
## Multiple R-squared: 0.2073, Adjusted R-squared: 0.1945
## F-statistic: 16.14 on 7 and 432 DF, p-value: < 2.2e-16
```

```
summary(mod_wo_interactions)
##
## Call:
## lm(formula = log(per_cap_income) ~ log(crimes) + region, data = cdi)
##
## Residuals:
##
       Min
                 10
                     Median
                                   30
                                           Max
## -0.68757 -0.10557 -0.01422 0.08905 0.78946
##
## Coefficients:
               Estimate Std. Error t value Pr(>|t|)
##
## (Intercept) 9.188431
                         0.079812 115.125 < 2e-16 ***
## log(crimes) 0.066695
                         0.008421
                                    7.920 2.00e-14 ***
               0.104458 0.025531 4.091 5.11e-05 ***
## regionNE
              -0.086983 0.023618 -3.683 0.00026 ***
## regionS
                         0.028167 -1.963 0.05033 .
## regionW
              -0.055280
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1854 on 435 degrees of freedom
## Multiple R-squared: 0.2032, Adjusted R-squared: 0.1959
## F-statistic: 27.74 on 4 and 435 DF, p-value: < 2.2e-16
anova(mod_w_interactions, mod_wo_interactions)
## Analysis of Variance Table
##
## Model 1: log(per_cap_income) ~ log(crimes) * region
## Model 2: log(per_cap_income) ~ log(crimes) + region
    Res.Df
              RSS Df Sum of Sq
##
                                    F Pr(>F)
## 1
       432 14.872
## 2
       435 14.949 -3 -0.076778 0.7434 0.5266
```

We also learn from the above output that, counter intuitively, per capita income is weakly positively related to the crime rate. Specifically, a 1% increase in the crime rate is associated with a 0.067% increase in per capita income.

Just to be thorough, we also attempted substituting per capita crime rate for total crimes and received very similar results - β_1 declined slightly from 0.067 to 0.042. Unfortunately, the explanatory power of the model was almost halved in the process. The per capita rate is likely better for our model since it allows us to adjust for the potential confounder population we identified in part A.

```
model_3 <- lm(log(per_cap_income) ~ log(per_capita_crimes) + region, data = mutate(cdi, per_capita_crimes = as
summary(model_3)
```

```
##
## Call:
## lm(formula = log(per_cap_income) ~ log(per_capita_crimes) + region,
##
       data = mutate(cdi, per_capita_crimes = as.numeric(crimes)/as.numeric(pop)))
##
## Residuals:
##
       Min
                  1Q
                      Median
                                    30
                                            Max
## -0.65832 -0.11431 -0.01548 0.10838 0.75657
##
## Coefficients:
##
                          Estimate Std. Error t value Pr(>|t|)
                                      0.06934 143.303 < 2e-16 ***
## (Intercept)
                          9.93628
## log(per_capita_crimes) 0.04243
                                      0.02148 1.975 0.04885 *
                                      0.02760 4.151 3.99e-05 ***
## regionNE
                          0.11457
```

```
-0.07456
## regionS
                                      0.02624
                                               -2.841
                                                       0.00471 **
                          -0.02426
                                      0.03002
## regionW
                                               -0.808
                                                       0.41952
##
## Signif. codes:
                          0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
                   0
##
## Residual standard error: 0.1974 on 435 degrees of freedom
## Multiple R-squared: 0.09645,
                                    Adjusted R-squared: 0.08814
## F-statistic: 11.61 on 4 and 435 DF, p-value: 5.776e-09
```

Figure 7 shows us more work will need to be done. The residuals vs fit plot shows linearity but it also hints that there might be omitted variable bias since the residuals form two distinct clusters. The variance appears constant by the residuals deviate significantly from the normal distribution.



- c) Use methods we have discussed in class and/or methods from Sheather Chapters 5, 6 & 7 (including, as needed: transformations, interactions, variable selection, residual analysis, fit indices, etc.) to find the multiple regression model predicting per-capita income from the other variables, that makes the "best" tradeoff between the following criteria:
- Reflects the social science and the meaning of the variables
- Satisfies modeling assumptions
- Clearly indicated by the data
- Can be explained to someone who is more interested in social, economic and health factors than in mathematics and statistics.

Answer:

In order to avoid making transformations that will be difficult to explain to a non-expert, we will limit ourselves to just log and power 2 transforms. We square the pct_hs_grad column to address the left-skewing and take the log of all the others except pct_below_pov, pop_18_34, and pop_65_plus. Income is dropped since it is a simply function of population and per capita income. The resulting saturated model has the residual patterns shown in Figure 8. Some non-linear relationship has not been

captured by our model as shown in the residual vs fit plot. The residuals are also not perfectly normally distributed - there is some evidence of overdispersion.



To address some of these shortcomings, and to deal with the multi-collinearity we are confident is present, we use stepwise regression to winnow down the predictor space. This eliminates the log_hosp_beds and log_crimes variables. None of the remaining variables have VIFs greater than 10 so we can be fairly confident multicollinearity is no longer a pressing concern.

```
best_cdi_model <- step(full_model_cdi, direction = "both")</pre>
```

vif(best_cdi_model)

##	log_doctors	log_land_area	log_pop	log_pct_bach	log_pct_ue
##	9.781512	1.180230	7.572984	4.650859	1.868772

##	<pre>pct_below_pov</pre>	sq_pct_hs	pop_18_34	pop_65_plus
##	2.387327	4.099777	1.988428	2.043196

Unfortunately, the same concerns are present in the new diagnostic plots below. This suggests we will need to consider other approaches in the future, perhaps lasso or another regularization technique.

