

# The Influence of Individual Characteristics on Public Transportation Planning\*

Iris Zhong

## Abstract

xx

## Warning: package 'tidyverse' was built under R version 3.5.3

## Warning: package 'ggplot2' was built under R version 3.5.3

## Warning: package 'tibble' was built under R version 3.5.3

## Warning: package 'tidyr' was built under R version 3.5.3

## Warning: package 'readr' was built under R version 3.5.3

## Warning: package 'purrr' was built under R version 3.5.3

## Warning: package 'dplyr' was built under R version 3.5.3

## Warning: package 'stringr' was built under R version 3.5.3

## Warning: package 'forcats' was built under R version 3.5.3

## Warning: package 'lubridate' was built under R version 3.5.3

## Warning: package 'stargazer' was built under R version 3.5.2

---

\*xx

```
## Warning: package 'corrplot' was built under R version 3.5.3

## Warning: package 'Hmisc' was built under R version 3.5.3

## Warning: package 'survival' was built under R version 3.5.3

## Warning: package 'Formula' was built under R version 3.5.2

## Warning: package 'dlookr' was built under R version 3.5.3

## Warning: package 'mice' was built under R version 3.5.3

## Warning: package 'car' was built under R version 3.5.3

## Warning: package 'carData' was built under R version 3.5.3

## Warning: package 'patchwork' was built under R version 3.5.3
```

## 1 Literature Review

Allen et al. (2016) study the reasoning of the failure of a referendum on a congestion charging scheme in Edinburgh. Instead of using direct voting data, they conduct a survey after the referendum, which allows them to ask more specific questions. Researchers can gain detailed data by surveying, because the unit of measurement is each individual; however, a possible disadvantage of surveying is that respondents who turn in the questionnaire tend to have stronger attitudes towards the proposal, generating sampling bias. They conclude that people who use cars as the primary transportation mean, demonstrate a misconception of the pricing plan, or question the effectiveness of the scheme at reducing congestion are more likely to oppose it. Their findings can give insights to the similar failure in the Gwinnett referendum. Voters against the proposal could be those who rarely use public transportation and those who are not convinced by the effectiveness of expanding public transit in alleviating the traffic.

Another crucial factor is the accessibility of the proposed transit system. Kinsey et al. (2010) examine the relationship between the distance to the scheduled railway station and voter turnout by studying the Seattle monorail referendum. They introduce the concept of diffused and concentrated benefit/cost. People who live far from the monorail enjoy the diffused benefit of less traffic congestion, and bear the diffused cost of

increased tax. People living close to the rail experience the same diffused benefit and cost, but they also gain the concentrated benefit of easily accessing the public good. Finally, those who live very close to the railway have the same benefits and costs, but they also face the concentrated cost such as inconvenience during construction. Since “people are more strongly motivated to avoid losses than to approach gains,” they expect a higher turnout rate in farther places with votes for “no,” which is verified from their analyses. Besides distance, they also find out precincts with a higher percentage of people of lower socioeconomic status or young people have a lower turnout rate. Interestingly, there is a significant interaction between partisanship and distance, which would be also tested in my study. In essence, the effect of distance on turnout is weakened by partisanship, and vanishes beyond a threshold of distance. Even though my dependent variable is voters’ responses rather than turnout, it can be inferred from Kinsey et al.’s findings that people farther away from the transit system would vote against the referendum more. However, the relationship might be non-linear and requires some form of transformation. Regarding the methods, they utilize the spatial lag model to correct for autocorrelation, which is proper to use in my project as well since both studies use precinct-level data.

## **2 Background**

### **2.1 current transportation**

### **2.2 Connect Gwinnett: Transit Plan**

#### **2.2.1 Plan overview**

#### **2.2.2 Financing**

### **2.3 referendum**

## **3 Data & Methods**

### **3.1 Conceptual model**

According to previous research, sociodemographic elements can influence people’s voting decisions in the referendum. For example, the effect of income is mixed: on the one hand, people with higher income will pay a smaller portion of their earnings for the implementation of the plan; on the other hand, they will pay

a larger amount of tax. Bollino (2008) finds a positive correlation between income and people’s willingness to pay for renewable resources. Burkhardt and Chan (2017) separate the influence of income from tax, and discover their opposite effects on voting. Therefore, it is worth considering the relationship between income and percentage of supporters in this referendum. Voters’ partisanship attachment is found to be a significant factor as well in Burkhardt and Chan’s (2017) paper. Areas with higher proportions of Republicans are less supportive of fiscally costly propositions. In my project, it can be hypothesized that tracts that have a higher proportion of Trump supporters tend to have a lower percentage of agreement to the proposal.

In addition, some factors related to transportation can intuitively shape people’s attitudes towards public transit. For example, the areas in which people do not use public transit at all might have a higher percentage of refusal of the proposal. People who have to travel a long time to work are more likely to support the extension plan if it helps save time. These two factors serve as controls in the models.

Finally, people favor the proposition if it benefits them. Specifically, tracts that are not covered by public transport at present but will be covered in the expansion plan are predicted to support the proposal more.

## 3.2 Data

The ballot results of the Gwinnett County referendum on Mar 19th, 2019 is obtained from a website powered by Scytl, a trusted source of election outcomes. The cross-sectional dataset contains the voting information of all 157 precincts in the county. The dependent variable – the proportion of supporters of the referendum, and the voter turnout rate are calculated from this data source.

The result of the 2016 Presidential Election is chosen to reflect partisanship. A cross-sectional precinct-level election data is obtained from the MIT Election Data and Science Lab website. The number of votes for Trump at each precinct in Gwinnett county comes from this dataset.

Next, cross-sectional sociodemographic characteristics at the census tract level are found in Census Bureau via an R package called `tidycensus`. The data comes from the American Community Survey 5-year estimate published in 2018. The median age and median income are collected here. In addition, the proportion of white, the percentage of people who go to work by public transportation, and the percentage of people who travel more than an hour to work are calculated by dividing the relevant variables by the total population or the survey sample size.

The information on whether the tract enjoys the proximity of public transportation now and future can be acquired from spatial maps and analyses. First, a precinct-level shapefile of Gwinnett County made in 2018 is obtained from the Georgia General Assembly. Gwinnett County maps with current and proposed future

**Table 1:** Variable definitions

Variable name	Description
GEOID	The geographic identifier of the census tract
medage	The median age of the population in the tract
medincome	The median income of the population in the tract
white_pct	The percentage of white population in the tract
public_pct	The percentage of people who go to work by public transportation (excluding taxi or cab)
time_pct	The percentage of people who travel more than an hour to work
trump_pct	The estimated percentage of votes for Donald Trump in that tract
voter_turnout	The estimated percentage of voters who voted in this referendum in the tract
yes_pct	The estimated percentage of voters who voted yes in this referendum in the tract
plan_yes	Whether the tract is covered by the public transportation now and in the short-range, defined by whether any transportation is available within 500 meters. 1 stands for the tract doesn't have transit both now and in the short-range plan. 2 stands for the tract has transit now but not in the short-range plan. 3 stands for the tract that doesn't have transit now and will have in the future. 4 stands for the tract that has public transit both now and in the future

public transit systems are available on the Gwinnett County government website. I select the short-range expansion plan (Y2020 – 2025) because the cost and benefit of the expansion in the far future are discounted more. After modifications in QGIS software, the maps are then transformed into spatial data readable in R. Five-hundred-meter buffer zones are created around bus and railway routes. The categorical variable *current\_plan* has four levels: 1 represents no transportation near this tract at present and in the future; 2 represents the tracts that are accessible to public transit currently but not in the future; 3 stands for the tracts that do not have transit at present but will do in the expansion plan; 4 represents the tracts that have and will have public transit for now and for the future.

As noted above, both referendum and 2016 election data are collected at the precinct level. However, the other datasets are performed at the census tract level. Therefore, referendum and election data are redistributed by the areas shared by the precinct and the tract. See the data appendix for detailed steps of transformation.

The final dataset joins the datasets above by GEOID. It is cross-sectional, measured with the unit of the census tract. Since there are 113 census tracts in Gwinnett, the dataset has 113 observations, with no missing data. A description of the variables can be found in *Table 1*.

```
data_numeric <- final_data %>%
  select(-c(current_plan,GEOID))

data_cor = cor(data_numeric)
```

**Table 2:** Summary statistics

Statistic	N	Mean	St. Dev.	Min	Pctl(25)	Pctl(75)	Max
medage	113	35.56	4.58	26	32.8	38.8	52
medincome	113	69,439.24	24,358.44	33,020	51,429	82,845	156,136
white_pct	113	0.48	0.15	0.17	0.38	0.61	0.89
public_pct	113	0.01	0.01	0	0.002	0.02	0
time_pct	113	0.16	0.05	0.04	0.12	0.20	0.31
trump_pct	113	0.40	0.15	0.11	0.27	0.52	0.69
voter_turnout	113	0.16	0.06	0.05	0.13	0.18	0.37
yes_pct	113	0.53	0.14	0.27	0.42	0.61	0.84

```

data_cor_1 <- rcorr(as.matrix(data_numeric))
M <- data_cor_1$r
p_mat <- data_cor_1$P
corrplot(M, type = "upper", order = "hclust",
          p.mat = p_mat, sig.level = 0.05)

```

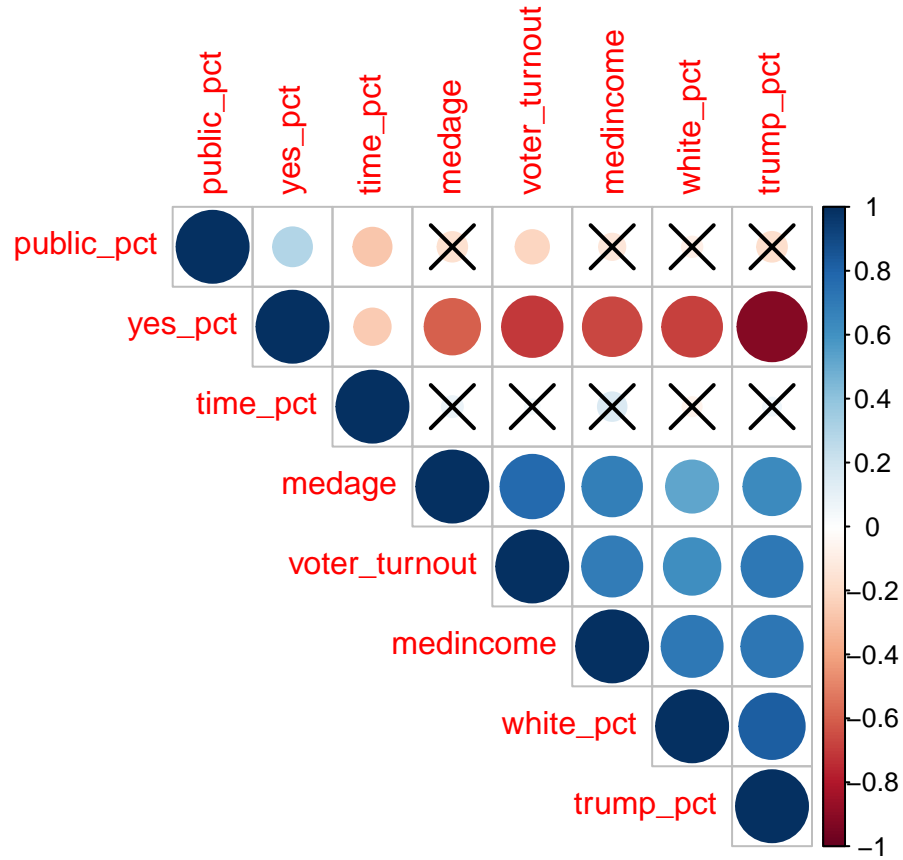


Table 2 lists the summary statistics. Median income is skewed to the right, with a few tracts demonstrating high income levels. Such a pattern of inequality is universally observed.

The usage of public transit is low. On average, only one percent of the population relies on public transportation to go to work. As proved by later analyses, the distribution is highly skewed, and will be transformed in some models. **public\_pct max needs to be fixed**

The variable *current\_plan* is not in *Table 2* because it is a categorical variable. In sum, 31 tracts do not have access to public transport within 500 meters, both now and in the short-range. There is one tract that is categorized as 2, indicating that it has public transit at present, but not in the short-range proposition. It could be due to the plan of reducing circuitous routing. Twenty-two tracts do not enjoy the proximity of public transport now, but will in the short-range. Lastly, 59 tracts have public transit available both at present and in the short-range plan. Since a level of this variable contains only one value (“2”), problems with degrees of freedom will potentially arise.

A correlation matrix is created from **corrplot** package to investigate the correlation between each pair of factors (see *Figure ???*). Positive correlations are in blue colors, and negative correlations are in red colors. The magnitude is reflected by the color intensity and the size of the dot. Non-significant ( $p > 0.05$ ) correlations are marked with a cross. *Current\_plan* is omitted in the matrix because it is categorical. The dependent variable *yes\_pct* is significantly correlated with every independent variable. *Medage*, *voter\_turnout*, *medincome*, *white\_pct*, and *trump\_pct* are all strongly positively correlated with each other, and all of them are negatively associated with *yes\_pct*. Given the number of strong, significant correlations among the factors, it is essential to check collinearity in the regression model.

The major limitation of the data is the unit conversion from precinct to census tract. Such a method assumes that residents in one precinct have the same characteristics, and population density is identical. Clearly, the assumptions cannot be satisfied in a real dataset.

### 3.3 Model specification

Four models are tested in the analysis.

Model 1 adopts linear regression model with all the original variables:

$$yes\_pct_t = \beta_0 + \beta_1 medincome_t + \beta_2 white\_pct_t + \beta_3 trump\_pct_t + \beta_4 current\_plan_t + \beta_5 medage_t + \beta_6 public\_pct_t + \beta_7 time\_pct_t + \beta_8 voter\_turnout_t + \epsilon_t$$

where  $t$  indexes census tract. Since all variables are measured in the same scope, no fixed effects are tested. The dependent variable is *yes\_pct*, the proportion of supporters of the referendum in a tract. Independent variables include *medincome*, *white\_pct*, *trump\_pct*, and *current\_plan*. For more information about these

variables, refer to *Table ???*. I hypothesize that the effect of *medincome* is ambiguous. As explained in the Conceptual Model section, people with higher income pay a larger amount of tax, but the tax takes up a smaller portion of their earnings compared to those with lower income. *White\_pct* is hypothesized to have a negative effect on *yes\_pct*, based on the negative correlation between the two variables. A higher percentage of Trump supporters is expected to predict a lower percentage of votes for “yes” in the referendum, as evidenced by common sense and previous literature (???). Finally, since tracts that do not have public transit now and will have it in the short-range benefit the most from the proposition, I predict that the tracts with this feature will have a higher level of *yes\_pct* than the others.

*Medage*, *public\_pct*, *time\_pct*, and *voter\_turnout* add to the model as controls. For example, tracts with higher voter turnout are hypothesized to have a lower *yes\_pct*, according to the rule of loss aversion – people who believe the referendum incur losses to them are inclined to participate in the referendum actively and vote against it, but people who like the proposition are less motivated to vote in nature.

*Medincome*, *public\_pct* and *voter\_turnout* are found to be positively skewed. Thus, Model 2 uses multiple linear regression after log or square root transformation on these three variables. Among them, since many of the values in *public\_pct* are 0, a constant has to add to the variable before taking log transformation.

$$yes\_pct_t = \beta_0 + \beta_1 \log(medincome_t) + \beta_2 white\_pct_t + \beta_3 trump\_pct_t + \beta_4 current\_plan_t + \beta_5 medage_t + \beta_6 \log(public\_pct_t + 0.01) + \beta_7 time\_pct_t + \beta_8 \sqrt{voter\_turnout_t} + \epsilon_t$$

## 4 Results & Discussion

### 4.1 Regression results

*Table ???* provides the results of Model 1. The equation is statistically significant ( $F(10, 102) = 113.22$ ,  $p < 0.05$ ). Adjusted  $R^2$  is 0.909, indicating that the independent variables in this specification explain a large portion of the variance in the dependent variable. Consistent with the hypothesis, *medincome* is not a significant predictor of *yes\_pct*. *Trump\_pct* significantly predicts *yes\_pct* as expected, and results in a large coefficient: a one percent increase in the proportion of Trump supporters decreases the percentage of voters in favor of the referendum by 0.886%, holding other variables constant.

However, unlike the outcome from the correlation matrix (see *Figure ???*), *white\_pct* significantly predicts *yes\_pct* in a positive direction. A higher percentage of the white population in the tract is associated with higher support of the proposition when holding other explanatory variables constant.



**Table 3:** Primary regression

	<i>Dependent variable:</i>
	yes_pct
medincome	0.00000 (0.00000)
white_pct	0.114** (0.053)
trump_pct	−0.886*** (0.057)
current_plan2	0.019 (0.044)
current_plan3	−0.032*** (0.012)
current_plan4	0.016 (0.012)
medage	0.003* (0.001)
public_pct	0.817** (0.338)
time_pct	−0.306*** (0.090)
voter_turnout	−0.326** (0.126)
Constant	0.804*** (0.047)
Observations	113
R <sup>2</sup>	0.917
Adjusted R <sup>2</sup>	0.909
Residual Std. Error	0.042 (df = 102)
F Statistic	113.224*** (df = 10; 102)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

Another unanticipated result is *current\_plan*. In contrast to baseline tracts that do not access public transit service both now and in the short-range (*current\_plan 1*), only *current\_plan 3* tracts differ from them significantly, in the opposite direction to the hypothesis. That is to say, tracts planned to be newly added to the public transit service have a significantly higher rejection of the proposition than the other tracts. The control variables are all significant factors. *Medage* and *public\_pct* positively predict the level of *yes\_pct*, while *time\_pct* and *voter\_turnout* negatively predict *yes\_pct*.

Model 2 uses multiple linear regression as well, except that the skewed variables *medincome*, *public\_pct*, and *voter\_turnout* are transformed. The overall outcome is identical to Model 1 (see *Table ???*). After log transformation, the effect of median income is still insignificant.

A series of assumptions are examined. First, linearity and homoscedasticity assumptions are met, as illustrated by the residual plots. Next, because several factors correlate with each other (see *Figure ???*), VIFs are calculated to detect multicollinearity. Since all of the independent variables exhibit a VIF below 5 in both models, no collinearity issue is detected. Lastly, there is a potential fault of using linear regression in this dataset. The value of the dependent variable *yes\_pct* is restricted to the range of 0 to 1, because it represents a percentage. On the other hand, the predicted outcome from linear regression is unbounded. As a result, I calculate the predicted values from the two models with the actual datasets. The predicted outcome is close to the actual value of *yes\_pct*, all between 0 and 1. Then, I continue testing by finding the minimum and maximum values of each independent variable in the dataset (to see the actual minimum and maximum values, see *Table ???*), and plug them into the models accordingly to gain the extreme predicted *yes\_pct*. The extreme values for Model 1 are 0.09 and 0.96, and 0.11 and 0.95 in Model 2, all within the interval [0,1]. Therefore, the range of the dependent variable in these two regression models meets the assumption.

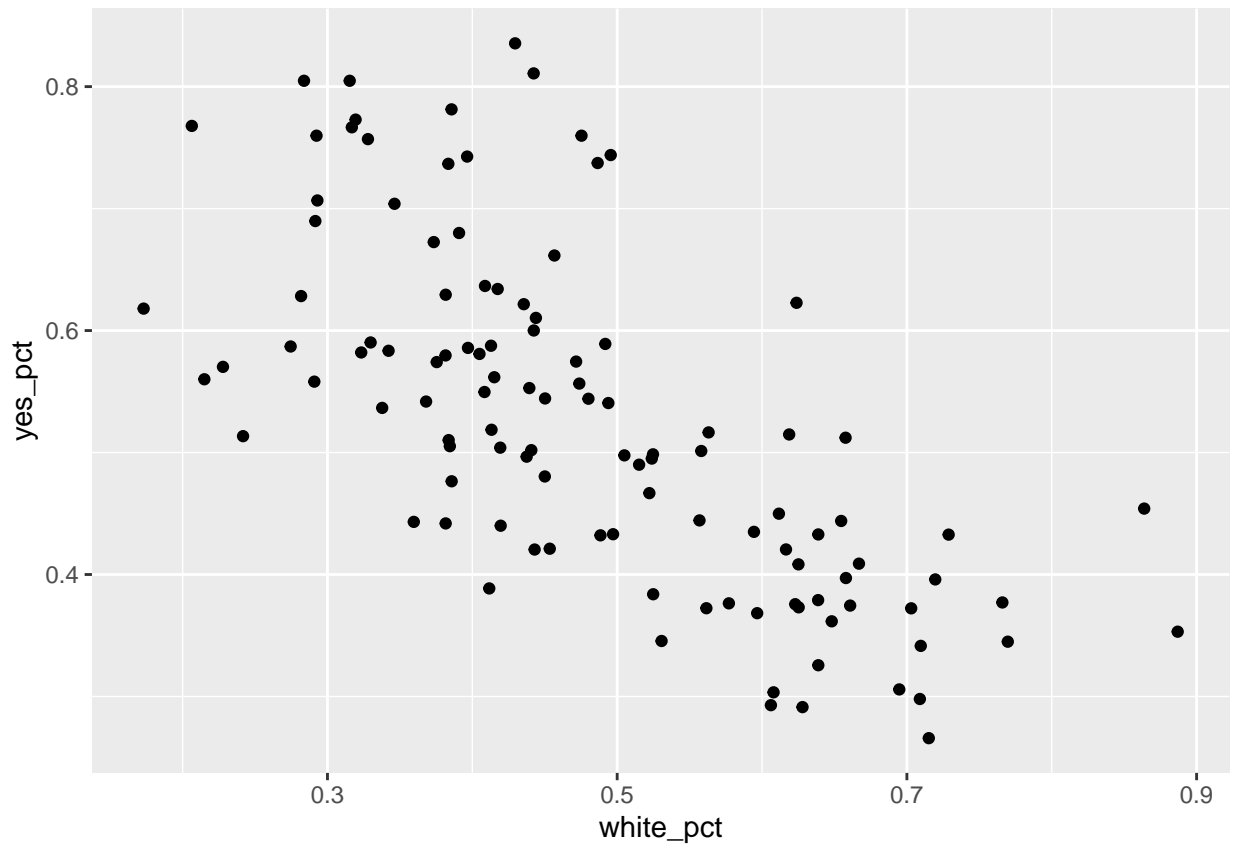
## 4.2 Discussion

## 4.3 Mediating factor

```
ggplot(final_data, aes(x = white_pct, y = yes_pct))+
  geom_point()
```

**Table 4:** Transformed regression

	<i>Dependent variable:</i>
	yes_pct
log_medincome	−0.008 (0.022)
white_pct	0.117** (0.051)
trump_pct	−0.857*** (0.059)
current_plan2	0.021 (0.043)
current_plan3	−0.033*** (0.012)
current_plan4	0.014 (0.011)
medage	0.004** (0.002)
log_public_pct	0.022*** (0.008)
time_pct	−0.281*** (0.090)
sqrt_voter_turnout	−0.319*** (0.103)
Constant	1.032*** (0.219)
Observations	113
R <sup>2</sup>	0.920
Adjusted R <sup>2</sup>	0.912
Residual Std. Error	0.041 (df = 102)
F Statistic	117.561*** (df = 10; 102)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01



```
library(psych)
```

```
## Warning: package 'psych' was built under R version 3.5.3
```

```
##
```

```
## Attaching package: 'psych'
```

```
## The following object is masked from 'package:car':
```

```
##
```

```
##   logit
```

```
## The following object is masked from 'package:dlookr':
```

```
##
```

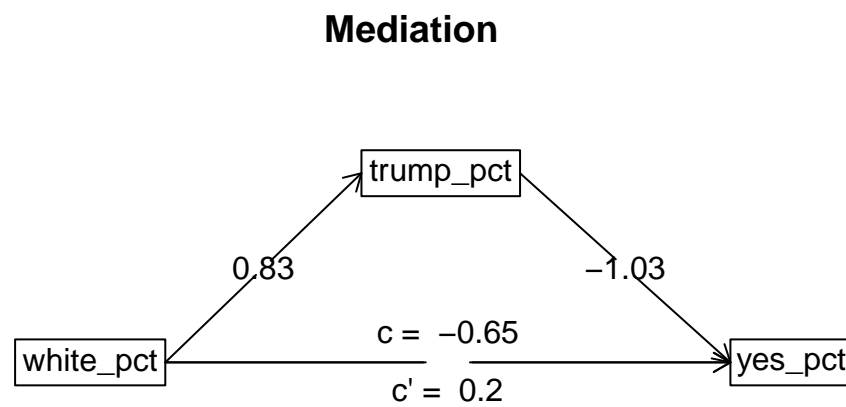
```
##   describe
```

```
## The following object is masked from 'package:Hmisc':
```

```
##
## describe

## The following objects are masked from 'package:ggplot2':
##
## %+%, alpha

mediate(yes_pct ~ white_pct + (trump_pct), data = final_data, n.iter = 10000)
```



## 5 References