

Awesome Paper

Zachary Bisson^{a,*}, Anna Elfstrum^{a,*}, Shawn Graven^{a,*}, Jack Kuivanen^{a,*},
Caleb Olson^{a,*}, Kate Soukkala^{a,*}

^aUniversity of Wisconsin-River Falls, 410 South Third Street, River Falls, WI, 54022

Abstract

This is the abstract.

It consists of two paragraphs.

1. Introduction

Legislation has been passed and social norms challenged to help level the playing field for different races and genders. Yet, according to Williams (Williams, 1987), “risk-averse employers believe and act as if black workers are on average less productive than their white counterparts; employers thus hire blacks at a wage discount or not at all.” Williams (Williams, 1987) goes on to say that there is a second case that “presumes blacks and white are equally productive on average, but black display a greater variance in ability; hence risk-averse employers’ hiring decision could precipitate a racial wage gap.” Due to this, it holds that business owners are more likely to make productivity and skill-based decisions based on race rather than incur the cost of acquiring and interpreting statistically significant data. This is witnessed by looking at historical wage data amongst seemingly disparate groups to see how over time, wages have increased but not at the same rate. The wage gap remains.

2. Literature Review

(Keating2019?)

Here are two sample references: (Pais, 2011; bob?).

3. Theoretical Analysis

$$H_0 : WhiteIncome \propto AllIncome$$

*Equal contribution, author order is alphabetical

Email addresses: zachary.bisson@my.uwrf.edu (Zachary Bisson),
anna.elfstrum@my.uwrf.edu (Anna Elfstrum), shawn.graven@my.uwrf.edu (Shawn Graven),
jack.kuivanen@my.uwrf.edu (Jack Kuivanen), caleb.olson@my.uwrf.edu (Caleb Olson),
katelyn.soukkala@my.uwrf.edu (Kate Soukkala)

Our null hypothesis is that there is no significant difference between the median income of individuals aged 16 and older in the United States and the median income of white individuals who are 16 and older in the United States.

$$H_A : \text{WhiteIncome} \not\propto \text{AllIncome}$$

Our alternate hypothesis is that the median income of white individuals aged 16 and older in the United States is significantly lower than the median income of individuals aged 16 and older.

$$H_0 : \text{BlackIncome} \propto \text{AllIncome}$$

Our null hypothesis is that there is no significant difference between the median income of individuals aged 16 and older in the United States and the median income of black individuals who are 16 and older in the United States.

$$H_A : \text{BlackIncome} \not\propto \text{AllIncome}$$

Our alternate hypothesis is that the median income of black individuals aged 16 and older in the United States is significantly lower than the median income of individuals aged 16 and older.

Our model is simple, `mod1` predicting `Median` only depends on the `Date`. `mod2` performs the $\log(\text{Median})$ as it makes the data more linear as shown above. These will represent *AllIncome*. `modBlack` includes the `Black` factor and will represent the Black population and non-black's. The same is done with `modWhite` and `modHispanic` with the white's and the hispanics respectively.

$$\text{mod1} : \text{Median} = \beta_1 \text{Date} + \beta_0 + e$$

$$\text{mod2} : \log(\text{Median}) = \beta_1 \text{Date} + \beta_0 + e$$

$$\text{modWhite} : \log(\text{Median}) = \beta_1 \text{Date} * \text{White} + \beta_0 + e$$

$$\text{modBlack} : \log(\text{Median}) = \beta_1 \text{Date} * \text{Black} + \beta_0 + e$$

$$\text{modHispanic} : \log(\text{Median}) = \beta_1 \text{Date} * \text{Hispanic} + \beta_0 + e$$

4. Empirical Analysis

4.1. Data

4.2. Import Data

4.2.1. About the Data

We have two sources of data, one from U.S. Bureau of Labor Statistics (BLS) and the majority of data from Economic Policy Institute (EPI).

BLS maintains a data set called `cpsaat`, this data summarizes the wage earnings per type of job, based on race and gender. To access the data in R we use a `curl_download` to retrieve the `.xlsx` file off the internet. To read the file we use the function `readxl::read_excel`.

EPI hosts a lot of data on wage statistics including, minimum wage, the participation, and earnings of each race, gender, education level, and much more. Due to the way EPI presents the data, it cannot be downloaded with `curl`. Instead, I have accessed the data with the package `epidata`, this simple package interfaces with EPI so that you don't have to manually download the data. EPI does not contain individual observations for wage, instead it provides 2 summarizations of the data grouped by race, age, gender, and education. This is the median, 50% of people make more and 50% of people make less than this value. The other one is mean, or they call average, this is the sum of wages added up and divided by the amount.

$$\bar{x} = \frac{\sum_{i=0}^{n-1} x_i}{n}$$

To reduce the effect of the highest earners we will be using the median, like they use in the housing market as a high outlier will only add one rather than a lot more.

4.2.2. Import cpsaat Data

Make sure we have internet and if not abort if not

`cpsaat` data is provided online at [bls.gov](https://www.bls.gov). As it is a direct link we can download it and save it to a temporary file and process the data with `readxl::read_excel()`

4.2.3. Import EPI Data

Get the data at EPI. As there is no direct link available we cannot use `curl`, instead there is a package that we can use to access the data, `epidata`. This will download data in the background.

4.3. Clean Data

As with most data, it will have to be cleaned. This includes pivoting the tibble into a longer tibble, as it will work better for `ggplot2`. This current format is called wide format as it has many columns. To fix this we can convert it into long format, as there are many rows, with `pivot_longer`. When we do this sometimes the new column we create contains more than one value, to remedy this issue we can use `separate` and mutate if necessary to get the values in the right column. Another inconsistency we should be aware of is that the currency values are in different years, not a large difference, but something that should be corrected.

4.3.1. Clean cpsaat11

Looks fine.

4.3.2. Clean Labor_force_participation

4.3.3. Clean Medianaverage_hourly_wages

4.3.4. Clean Minimum_wage

This data has data in terms of 2018, the other data is in 2019 USD. As it will be easiest and the latest data, we will be using 2019. Although small, there will be a difference and we need to adjust for inflation. The package `priceR` allows us to convert those monetary values into other ones using online inflation data.

4.3.5. Fix inconsistant case

As the data was imported with `epidata`, the column names have been changed from what the csv has. So we need to fix that to conform to consistency.

4.4. Export the data as .csv files

To backup our data we will export the cleaned tibbles.

```
if(!dir.exists("../data"))dir.create("../data")

cpsaat11%>%
  write_csv("../data/cpsaat11.csv")

Minimum_wage%>%
  write_csv("../data/Minimum_wage.csv")

Participation%>%
  write_csv("../data/Participation.csv")

Wages%>%
  write_csv("../data/Wages.csv")
```

4.5. Methodology

4.6. Results

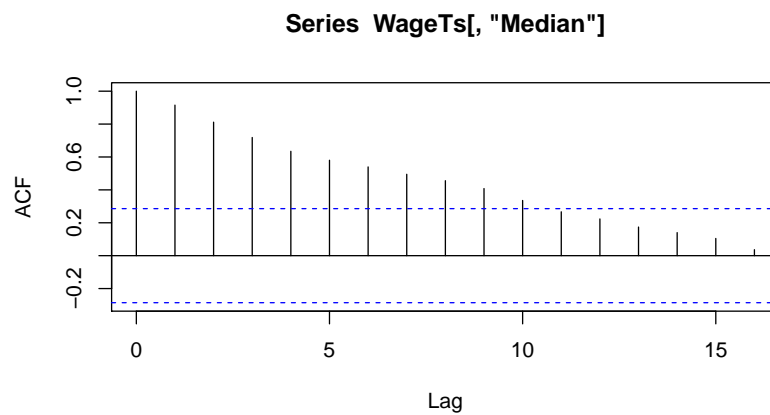
4.7. Estimators

$$mod0 : Median = Date\beta_1 + \beta_0 + e$$

Data is already aggregated

```
WagesAll=Wages%>%
  filter(is.na(Race),is.na(Gender))
WageTs=ts(WagesAll, start = min(WagesAll$Date), end = max(WagesAll$Date), frequency = 1)

acf(WageTs[, "Median"])
```



There is a lot of autocorelation

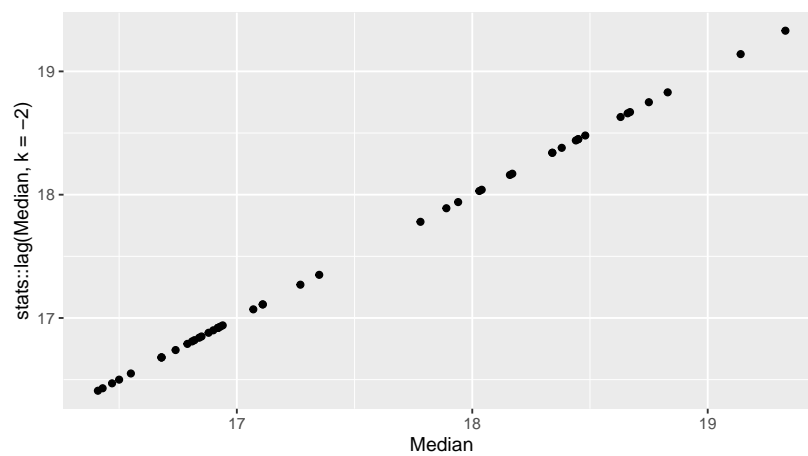
```
mod0=dynlm(Median~Date, data = Wages)

bgtest(mod0, order = 1, type = "F", fill = NA)
```

Breusch-Godfrey test for serial correlation of order up to 1

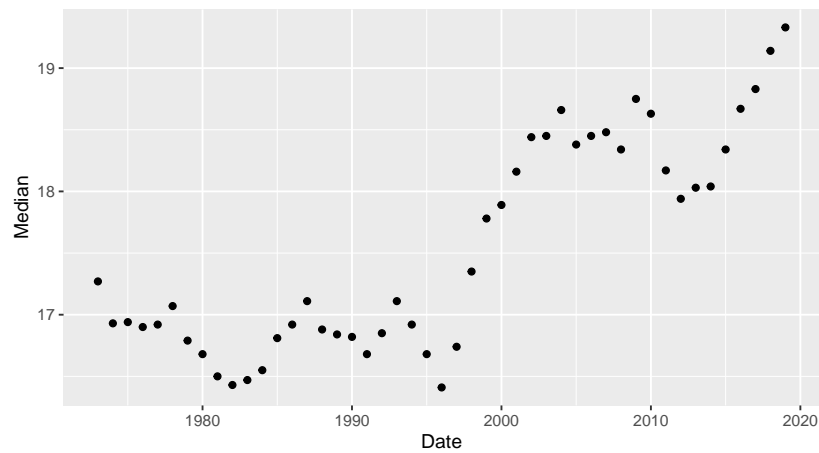
```
data: mod0
LM test = 10.624, df1 = 1, df2 = 560, p-value = 0.001184
```

```
WagesAll%>%
  ggplot(aes(x=Median, y=stats::lag(Median, k=-2)))+
  geom_point()
```



```
mod1=lm(Median~Date, data = Wages)
mod2=lm(log(Median)~Date, data = Wages)
```

```
WagesAll%>%
  ggplot(aes(x=Date, y=Median))+
  geom_point()
```



```
chow=function(racestr){
  WagesRace=Wages%>%
    mutate(R=if_else(Race==racestr, 1, 0))%>%
    filter(!is.na(Race),is.na(Gender))

  mod2=lm(log(Median)~Date,
          data=WagesRace
  )

  modRace=lm(log(Median)~Date*R,
             data=WagesRace
  )

  stargazer(mod2, modRace,
            header=FALSE,
            type=knittype,
            title="Model comparison, 'wage' equation",
            keep.stat="n",digits=2, single.row=TRUE,
            intercept.bottom=FALSE
  )
}
```

```

    anova(mod2, modRace)%>%
      kable()
}
chow("white")

```

```

\begin{table}[!htbp] \centering
  \caption{Model comparison, 'wage' equation}
  \label{}
\begin{tabular}{@{\extracolsep{5pt}}lcc}
\\[-1.8ex]\hline
\hline \\[-1.8ex]
& \multicolumn{2}{c}{\textit{Dependent variable:}} \\
\cline{2-3}
\\[-1.8ex] & \multicolumn{2}{c}{\log(Median)} \\
\\[-1.8ex] & (1) & (2)\\
\hline \\[-1.8ex]
Constant &  $-\$3.21^{*}$  &  $-\$1.62^{**}$  \\
Date &  $0.003^{***}$  &  $0.002^{***}$  \\
R &  $-\$4.77^{***}$  & \\
Date:R &  $0.003^{***}$  & \\
\hline \\[-1.8ex]
Observations & 141 & 141 \\
\hline
\hline \\[-1.8ex]
\textit{Note:} & \multicolumn{2}{r}{ $^{*}p<0.1$ ;  $^{**}p<0.05$ ;  $^{***}p<0.01$ } \\
\end{tabular}
\end{table}

```

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
139	2.5887721	NA	NA	NA	NA
137	0.2989179	2	2.289854	524.7428	0

After performing a chow test we can reject our null hypothesis that there is no significant difference between the median income of individuals aged 16 and older in the United States and the median income of white individuals who are 16 and older in the United States, since our p-value is less than 0.01. We conclude that the median income of white individuals aged 16 and older in the United States is significantly higher than the median income of individuals aged 16 and older.

```
chow("black")
```

```

\begin{table}[!htbp] \centering
  \caption{Model comparison, 'wage' equation}
  \label{}
\begin{tabular}{@{\extracolsep{5pt}}lcc}
\\[-1.8ex]\hline

```

```

\hline \[-1.8ex]
& \multicolumn{2}{c}{\textit{Dependent variable:}} \\\
\cline{2-3}
\[-1.8ex] & \multicolumn{2}{c}{\log(Median)} \\\
\[-1.8ex] & (1) & (2)\\
\hline \[-1.8ex]
Constant &  $-\$3.21^{*}$  (1.69) &  $-\$3.10$  (1.99) \\\
Date &  $0.003^{***}$  (0.001) &  $0.003^{***}$  (0.001) \\\
R &  $-\$0.34$  (3.45) \\\
Date:R &  $0.0001$  (0.002) \\\
\hline \[-1.8ex]
Observations & 141 & 141 \\\
\hline
\hline \[-1.8ex]
\textit{Note:} & \multicolumn{2}{r}{ $^{*}p < \$0.1$ ;  $^{**}p < \$0.05$ ;  $^{***}p < \$0.01$ } \\\
\end{tabular}
\end{table}

```

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
139	2.588772	NA	NA	NA	NA
137	2.361682	2	0.2270902	6.586694	0.0018567

After performing another chow test we can reject our null hypothesis that there is no significant difference between the median income of individuals aged 16 and older in the United States and the median income of black individuals who are 16 and older in the United States, since our p-value is less than 0.01. We conclude that the median income of black individuals aged 16 and older in the United States is significantly lower than the median income of individuals aged 16 and older.

```
chow("hispanic")
```

```

\begin{table}[!htbp] \centering
\caption{Model comparison, 'wage' equation}
\label{}
\begin{tabular}{@{\extracolsep{5pt}}lcc}
\[-1.8ex]\hline
\hline \[-1.8ex]
& \multicolumn{2}{c}{\textit{Dependent variable:}} \\\
\cline{2-3}
\[-1.8ex] & \multicolumn{2}{c}{\log(Median)} \\\
\[-1.8ex] & (1) & (2)\\
\hline \[-1.8ex]
Constant &  $-\$3.21^{*}$  (1.69) &  $-\$4.92^{***}$  (1.59) \\\
Date &  $0.003^{***}$  (0.001) &  $0.004^{***}$  (0.001) \\\
R &  $5.11^{*}$  (2.75) \\\
Date:R &  $-\$0.003^{*}$  (0.001) \\\

```



```

\hline \[-1.8ex]
Observations & 141 & 141 \\\
\hline
\hline \[-1.8ex]
\textit{Note:} & \multicolumn{2}{r}{\textit{*}}$p$<$0.1; \textit{**}}$p$<$0.05; \textit{***}}$p$<$0.01} \\\
\end{tabular}
\end{table}

```

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
139	2.588772	NA	NA	NA	NA
137	1.498359	2	1.090413	49.85007	0

After performing our final chow test, we can reject our null hypothesis that there is no significant difference between the median income of individuals aged 16 and older in the United States and the median income.

Some limitations to the experiment are the data collection. This is because we are unable to collect everyone's income in the united states to test this. However, the data we do have gives a good representation of the income of people as we currently know it in the United States. Another major issue would be the voluntary data used. People who volunteer to give out this data may not participate due to their current financial status. This would skew the data and ultimately change the outcome.

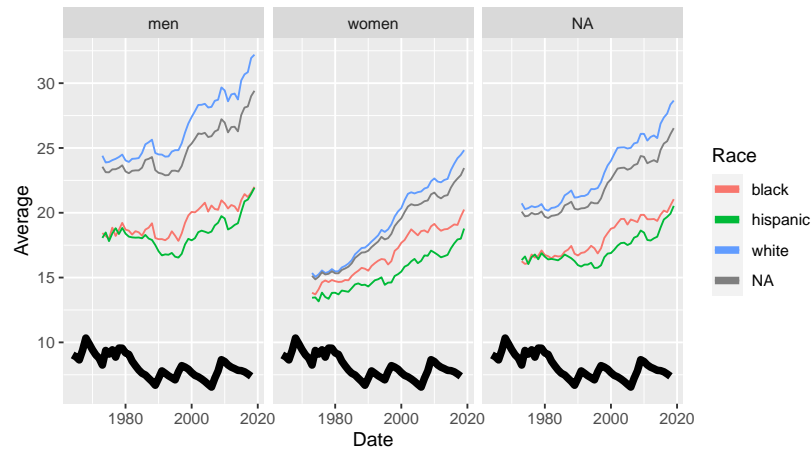
4.8. Wage over Time by Race and Gender

4.8.1. Average and Medium Wage over Time by Race and Gender

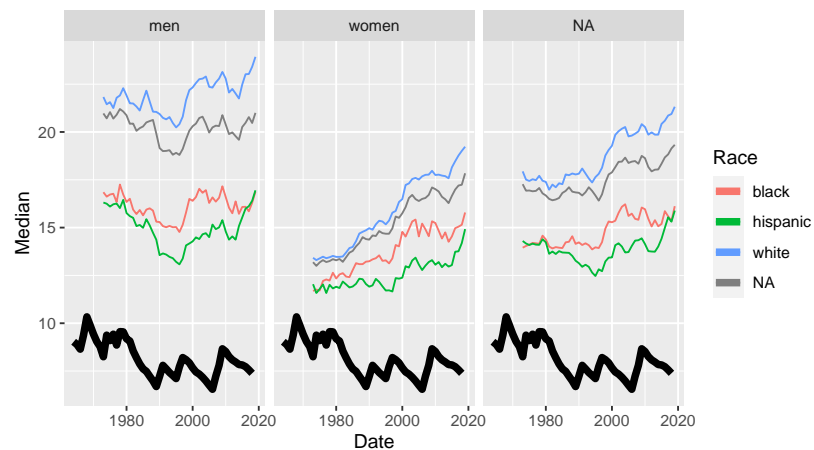
```

g=Wages%>%
  ggplot(aes(col=Race, x=Date))+
  geom_line(aes(y=Average))+
  geom_line(aes(y=Min2019, col=NULL), data=Minimum_wage, size=2)+
  facet_wrap(~Gender)
ggdisp(g)

```



```
g=Wages%>%
  ggplot(aes(col=Race, x=Date))+
  geom_line(aes(y=Median))+
  geom_line(aes(y=Min2019, col=NULL), data=Minimum_wage, size=2)+
  facet_wrap(~Gender)
ggdisp(g)
```



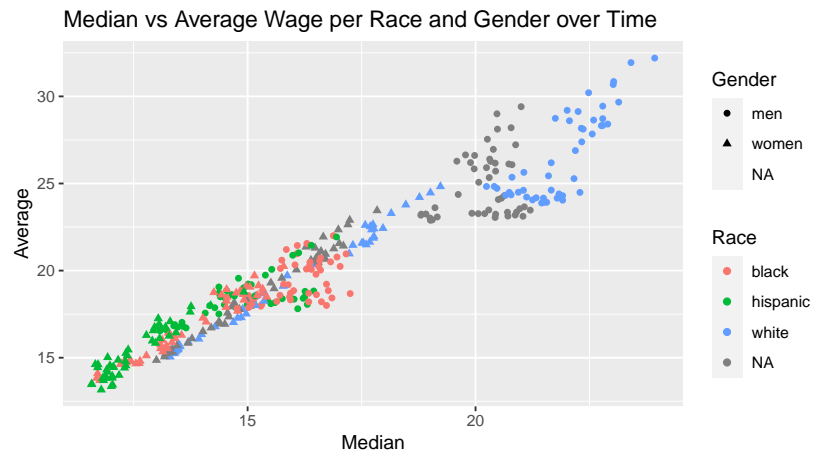
4.8.2. Scatter Plot over Time

```
g=Wages%>%
  ggplot()+
  geom_point(aes(x=Median, y=Average, col=Race, shape=Gender, frame=Date))+
  ggtitle("Median vs Average Wage per Race and Gender over Time")
```

Warning: Ignoring unknown aesthetics: frame

```
ggdisp(g)
```

Warning: Removed 188 rows containing missing values (geom_point).



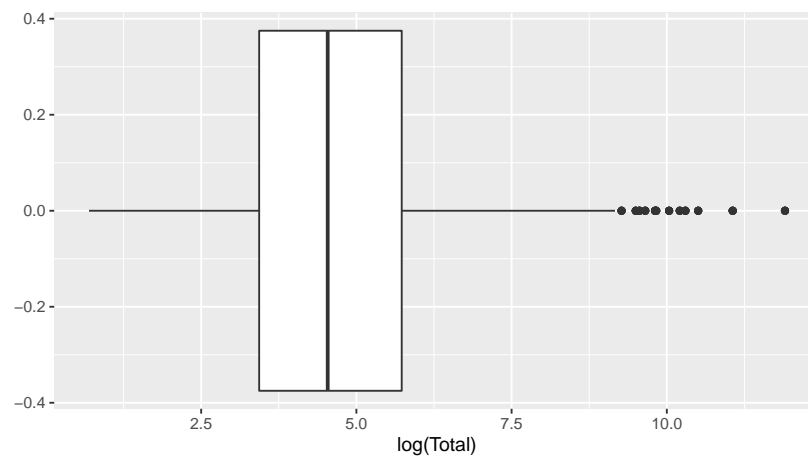
4.9. Wages according to Jobs

4.9.1. Sumarise data according to income of jobs

This data is currently unusable as there is only one opservation per type of job, we don't have over time statistics. We do however, have a snapshot of the diverse earnings, we don't care what the job is, but the average wage of each race per earning bracket.

```
cpsaat11%>%
  ggplot(aes(x=log(Total)))+
  geom_boxplot()
```

Warning: Removed 10 rows containing non-finite values (stat_boxplot).



```

# Generate the percentiles
se=quantile(log(cpsaat11$Total), seq(0, 1, by=.1), na.rm=T)

# Add outliers
se["200%"]=Inf

# break into groups and drop NAs
d=cpsaat11%>%
  drop_na(Percentage)%>%
  group_by(gr=cut(Total, breaks=exp(se)), Race)

# Summarize the data and remove women as it is not a race
# This is so it add up to 100% or so
d=d%>%
  summarise(Percentage=mean(Percentage), Total=mean(Total))%>%
  filter(Race!="Women")
d

```

```

# A tibble: 32 x 4
# Groups:   gr [8]
  gr      Race      Percentage Total
<fct>   <chr>      <dbl>   <dbl>
1 (40,60] Asian          3.72   53.7
2 (40,60] Black/African American 10.9   53.7
3 (40,60] Hispanic/Latino    16.9   53.7
4 (40,60] White            82.1   53.7
5 (60,93] Asian          8.60   74.6
6 (60,93] Black/African American 13.1   74.6
7 (60,93] Hispanic/Latino    14.1   74.6
8 (60,93] White           74.6   74.6
9 (93,131] Asian          5.88  110.
10 (93,131] Black/African American 11.9  110.
# ... with 22 more rows

```

4.9.2. Is there missing data

```

cpsaat11%>%
  drop_na(Percentage)%>%
  filter(Total<30)

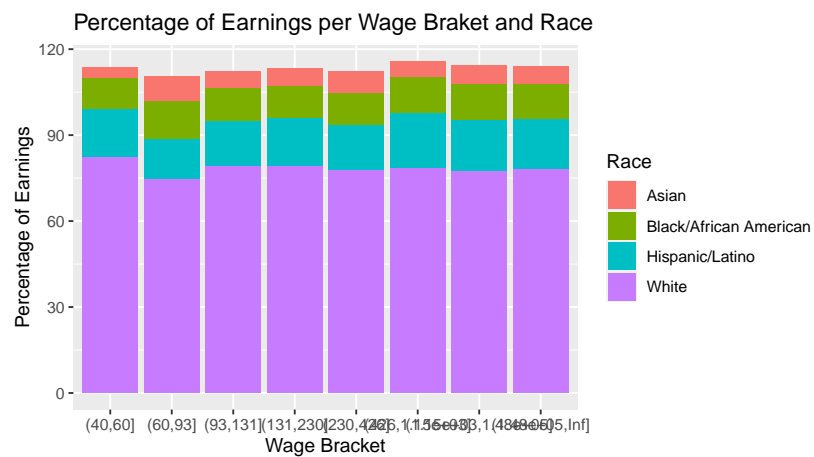
```

Occupation	Total	Race	Percentage
------------	-------	------	------------

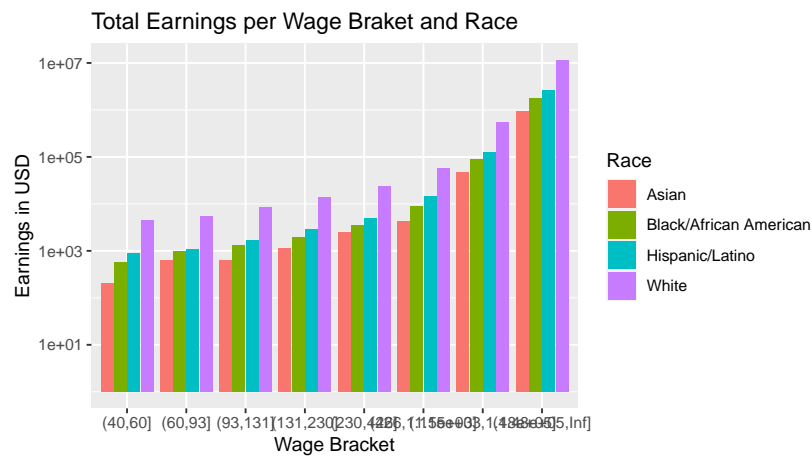
No, we just have a lack of observations for poor paying jobs.

4.9.3. Graph

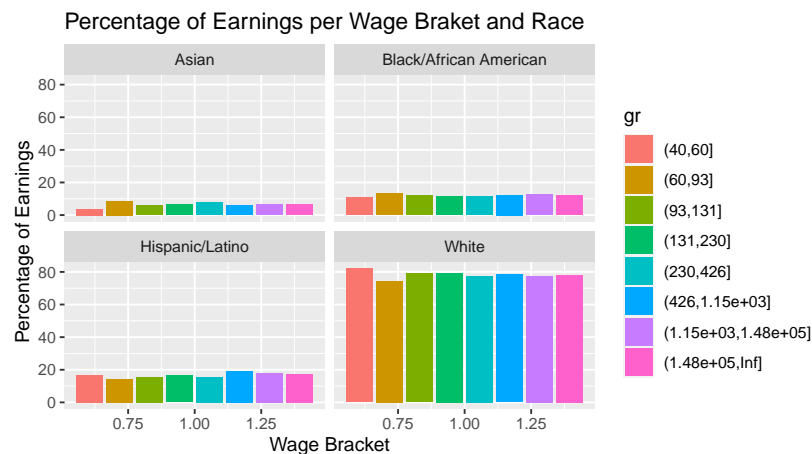
```
g=d%>%
  ggplot(aes(fill=Race, y=Percentage, x=gr))+
  geom_col()+
  xlab("Wage Bracket")+
  ylab("Percentage of Earnings")+
  ggtitle("Percentage of Earnings per Wage Braket and Race")
ggdisp(g)
```



```
g=d%>%
  ggplot(aes(fill=Race, y=Percentage*Total, x=gr))+
  geom_col(position = "dodge2")+
  scale_y_log10()+
  xlab("Wage Bracket")+
  ylab("Earnings in USD")+
  ggtitle("Total Earnings per Wage Braket and Race")
ggdisp(g)
```



```
g=d%>%
  ggplot(aes(fill=gr, x=1, y=Percentage))+
  geom_col(position = "dodge2")+
  facet_wrap(~Race)+
  xlab("Wage Bracket")+
  ylab("Percentage of Earnings")+
  ggtitle("Percentage of Earnings per Wage Bracket and Race")
ggdisp(g)
```

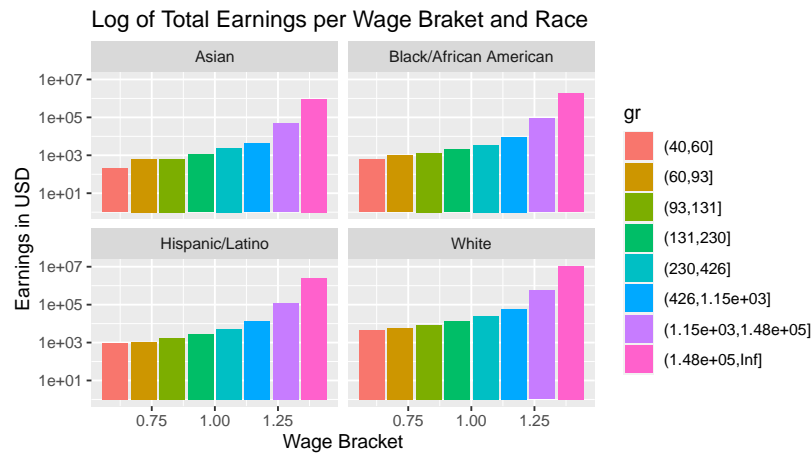


```
g=d%>%
  ggplot(aes(fill=gr, x=1, y=Percentage*Total))+
  geom_col(position = "dodge2")+
  facet_wrap(~Race)+
  scale_y_log10()+
```

```

xlab("Wage Bracket")+
ylab("Earnings in USD")+
ggtitle("Log of Total Earnings per Wage Bracket and Race")
ggdisp(g)

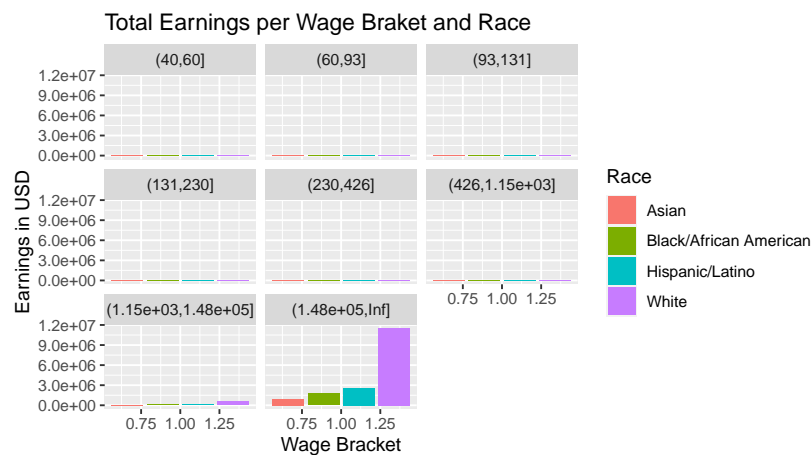
```



```

g=d%>%
ggplot(aes(fill=Race, x=1, y=Percentage*Total))+
geom_col(position = "dodge2")+
facet_wrap(~gr)+
xlab("Wage Bracket")+
ylab("Earnings in USD")+
ggtitle("Total Earnings per Wage Bracket and Race")
ggdisp(g)

```



5. Conclusion

6. References

- Pais, J., 2011. Socioeconomic background and racial earnings inequality: A propensity score analysis. *Social Science Research* 40, 37–49. doi:10.1016/j.ssresearch.2010.06.016
- Williams, R.M., 1987. Capital, competition, and discrimination: A reconsideration of racial earnings inequality. *Review of Radical Political Economics* 19, 1–15.