

Stimulus Checks and Consumer Spending

Aaron Davis, Navya Sonti, Rujula Nadipi, Swapnil Sethi, and Ujas Shah

9/2/2021

Question of Interest

We want to know whether or not the stimulus checks sent out by the US government have had a positive impact on the economy (using consumer spending as a proxy for how healthy the economy is).

Why Should YOU Care How Healthy the US Economy Is?

Why should we care about consumer spending? It seems like a big picture idea that won't really effect any of us specifically, right? Wrong. When the economy is healthy, there are more jobs available. Those jobs also pay more. As graduate students, we all want to get good paying jobs as Data Scientists, and that will be more likely to happen more quickly if the economy is healthy.

This analysis will help us understand how stimulus checks effect the US economy, and therefore, indirectly, the analysis will also help us understand whether or not stimulus checks will help us get good-paying jobs quickly after graduation.

There is also a long-standing debate between the liberal and conservative political ideologies in the US over whether or not the increase in consumer spending caused by stimulus checks is actually worth the debt incurred by the government when sending out the checks. While this analysis doesn't cover that, this would be an interesting topic for future research using this dataset.

Data Source

All of our data was aggregated by Opportunity Insights at <https://github.com/OpportunityInsights/EconomicTracker>. In this analysis, we use spending data provided by Affinity Solutions, job postings data from Burning Glass Technologies, COVID data from the CDC, GPS mobility reports from Google, unemployment claims from the Department of Labor, and employment levels from Paychex, Intuit, Earnin and Kronos.

Primary Reference:

"The Economic Impacts of COVID-19: Evidence from a New Public Database Built Using Private Sector Data", by Raj Chetty, John Friedman, Nathaniel Hendren, Michael Stepner, and the Opportunity Insights Team. November 2020. Available at: https://opportunityinsights.org/wp-content/uploads/2020/05/tracker_paper.pdf

Read in Data from GitHub Repository

Data Selection and Cleaning

In this code chunk, we select the columns that we're interested in from the dataframe that we joined a couple of steps ago. We could use all of the data in the dataframe, but we choose not to, since not all of the features will be helpful in answering the question of whether or not the stimulus checks have boosted the economy. We also combined "month", "day", and "year" columns into a "date" column.

```
consumer_spending <- affinity_daily_df %>%
  select(year, month, day, statefips, spend_all)
visits_to_retail <- move_daily_df %>%
  select(year, month, day, statefips, gps_retail_and_recreation)
employment <- employment_daily_df %>%
  select(year, month, day, statefips, emp)
```

Join the datasets we're interested in into one dataset.

Here we join the datasets of interest based on a shared date and state of measurements. We are joining weekly datasets (job listings and unemployment insurance claims) with the daily datasets. This will leave a bunch of **NA** values. We'll come back to fix that later.

```
regression1_df <- left_join(consumer_spending, visits_to_retail, by = c("year"
  , "month"
  , "day", "statefips"))
regression1_df <- left_join(regression1_df, employment, by = c("year", "month"
  , "day"
  , "statefips"))
```

Adding Stimulus Check Data

Here we add in the data for the COVID stimulus checks. We do this by creating a feature encoding for each check, where the value for that check is **0** before the check is sent out, then \$1200 for two weeks after the first check, then zero, then \$600 for two weeks after the second check, then zero, then \$1400 for two weeks after the third check, and then zero.

```
regression1_df <- regression1_df %>%
  unite("date", day:month:year, remove = FALSE, sep = "-")

## Warning in x:y: numerical expression has 2 elements: only the first used

regression1_df$date <- dmy(regression1_df$date)

regression1_df <- regression1_df %>%
  mutate(stimulus_checks = ifelse(date < ymd("2020-04-15"), 0, ifelse(date <
    ymd("2020-05-01"),
    1200, ifelse(date < ymd("2021-01-04"), 0, ifelse(date < ymd("
    2021-01-19"),
    600, ifelse(date < ymd("2021-03-17"), 0, ifelse(date < ymd("
    2021-04-01"),
    1400, 0)))))))
```

To fix the **NA** values found in our dataset, we replace them all with **0**. We could have handled them many other ways, like replacing them with the column mean, median, mode, or by training a regression model to fill them based on the rows that were not missing those values. For this dataset, though, we found that **NAs** are frequently used when there was no interesting data to report (in other words, the value for the feature was zero). This can be seen particularly in columns like `new_death_count` and `new_case_count` from the CDC COVID dataset.

We also drop rows that have 0 spending data because this value is not realistic. Also, we need the `spend_all` column to be clean because we will be using it for plotting and training a regression model later on.

```
regression1_df$spend_all <- ifelse(regression1_df$spend_all == ".", NA,
  regression1_df$spend_all)
regression1_df$spend_all <- as.double(regression1_df$spend_all)

regression1_df$emp <- ifelse(regression1_df$emp == ".", NA, regression1_df$emp)
regression1_df$emp <- as.double(regression1_df$emp)

glimpse(regression1_df)

## Rows: 31,008
## Columns: 9
## $ date <date> 2020-01-01, 2020-01-01, 2020-01-01,
## 2020-01-~
## $ year <dbl> 2020, 2020, 2020, 2020, 2020, 2020, 2020,
## 20~
## $ month <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
## 1,~
## $ day <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
## 1,~
## $ statefips <dbl> 1, 2, 4, 5, 6, 8, 9, 10, 11, 12, 13, 15,
## 16,~
## $ spend_all <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA,
## NA, ~
## $ gps_retail_and_recreation <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA,
## NA, ~
## $ emp <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA,
## NA, ~
## $ stimulus_checks <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
## 0,~
```

Training the First Linear Model on Our Data

In this code chunk, we fit a linear model to the `gps_retail_and_recreation`, `emp`, and `stimulus_checks` features, with the goal of predicting the `spend_all` variable. The working assumption here is that `spend_all` is a dependent variable, and the others are independent variables.

The reason we chose these three input variables is that `stimulus_checks` are directly relevant to our analysis. We are trying to identify if these features had a positive or negative effect on the economy. These two features, `gps_retail_and_recreation` and `emp`, were added to help account for change in spending that could not necessarily be accounted for by the `stimulus_checks`. We chose not to include any other input features because we believe that `gps_retail_and_recreation` and `emp` account for much of the possible variance without over-complicating our linear model.

```
spending_regression <- lm(spend_all ~ gps_retail_and_recreation + emp +
  stimulus_checks,
```

```

    regression1_df)

summary(spending_regression)

##
## Call:
## lm(formula = spend_all ~ gps_retail_and_recreation + emp + stimulus_checks,
##     data = regression1_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.33348 -0.07187 -0.00643  0.06591  0.57687
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    7.981e-02   1.003e-03   79.54  <2e-16 ***
## gps_retail_and_recreation 2.000e-01   7.927e-03   25.23  <2e-16 ***
## emp            8.467e-01   1.409e-02   60.10  <2e-16 ***
## stimulus_checks 5.551e-05   2.139e-06   25.95  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1038 on 22536 degrees of freedom
## (8468 observations deleted due to missingness)
## Multiple R-squared:  0.4067, Adjusted R-squared:  0.4067
## F-statistic: 5150 on 3 and 22536 DF, p-value: < 2.2e-16

visits_to_retail_weekly <- move_daily_df %>%
  select(year, month, day, statefips, gps_retail_and_recreation)
visits_to_retail_weekly <- visits_to_retail_weekly %>%
  unite("date", day:month:year, remove = FALSE, sep = "-")

## Warning in x:y: numerical expression has 2 elements: only the first used

visits_to_retail_weekly$date = dmy(visits_to_retail_weekly$date)
visits_to_retail_weekly <- mutate(visits_to_retail_weekly, week = isoweek(date))
visits_to_retail_weekly <- visits_to_retail_weekly %>%
  group_by(year, week, statefips) %>%
  summarise(gps_retail_and_recreation = mean(gps_retail_and_recreation),
            date = max(date))

## 'summarise()' has grouped output by 'year', 'week'. You can override using
## the '.groups' argument.

unemployment_claims <- ui_claims_weekly_df %>%
  select(year, month, day_endofweek, statefips, contclaims_rate_combined)
unemployment_claims <- unemployment_claims %>%
  unite("date", day_endofweek:month:year, remove = FALSE, sep = "-")

## Warning in x:y: numerical expression has 2 elements: only the first used

unemployment_claims$date <- dmy(unemployment_claims$date)
unemployment_claims$contclaims_rate_combined <- as.double(unemployment_claims$
  contclaims_rate_combined)

```

```

## Warning: NAs introduced by coercion

glimpse(unemployment_claims)

## Rows: 4,488
## Columns: 6
## $ date          <date> 2020-01-04, 2020-01-04, 2020-01-04,
##   2020-01-~
## $ year          <dbl> 2020, 2020, 2020, 2020, 2020, 2020, 2020,
##   202~
## $ month         <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
##   1, ~
## $ day_endofweek <dbl> 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4,
##   4, ~
## $ statefips     <dbl> 1, 2, 4, 5, 6, 8, 9, 10, 11, 12, 13, 15,
##   16, ~
## $ contclaims_rate_combined <dbl> 0.826, 2.930, 0.507, 1.120, 1.850, 0.679,
##   2.4~

job_postings <- job_listings_weekly_df %>%
  select(year, month, day_endofweek, statefips, bg_posts)
job_postings <- job_postings %>%
  unite("date", day_endofweek:month:year, remove = FALSE, sep = "-")

## Warning in x:y: numerical expression has 2 elements: only the first used

job_postings$date <- dmy(job_postings$date)
job_postings$date <- job_postings$date + 1
job_postings$day_endofweek <- job_postings$day_endofweek + 1
glimpse(job_postings)

## Rows: 4,539
## Columns: 6
## $ date          <date> 2020-01-11, 2020-03-14, 2021-05-01, 2020-08-01,
##   2020-07~
## $ year          <dbl> 2020, 2020, 2021, 2020, 2020, 2021, 2020, 2021, 2020,
##   20~
## $ month         <dbl> 1, 3, 4, 7, 7, 1, 2, 4, 4, 4, 9, 5, 8, 10, 5, 6, 9,
##   6, 3~
## $ day_endofweek <dbl> 11, 14, 31, 32, 4, 2, 29, 24, 11, 18, 18, 15, 7, 10,
##   9, ~
## $ statefips     <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
##   1,~
## $ bg_posts      <dbl> 0.0777, -0.0841, 0.6190, -0.0527, -0.1440, -0.2940,
##   0.01~

regression2_df <- left_join(unemployment_claims, job_postings, by = c("date",
  "statefips",
  "month", "year", "day_endofweek"))
regression2_df <- regression2_df %>%
  mutate(week = isoweek(date))
regression2_df <- left_join(regression2_df, visits_to_retail_weekly, by = c("
  week",
  "statefips", "year"))

```

```

regression2_df <- left_join(regression2_df, state_id, by = c("statefips"))
regression2_df <- rename(regression2_df, c(visits_to_retail_and_recreation = "
  gps_retail_and_recreation",
  date = "date.x"))

states_which_stopped_early_benefits <- c("Arizona", "Indiana", "Maryland", "
  Tennessee",
  "Alaska", "Iowa", "Mississippi", "Missouri", "Alabama", "Idaho", "Nebraska",
  "New_Hampshire", "North_Dakota", "West_Virginia", "Wyoming", "Arkansas", "
  Florida",
  "Georgia", "Ohio", "Oklahoma", "South_Carolina", "South_Dakota", "Texas",
  "Utah",
  "Montana")

regression2_df <- regression2_df %>%
  mutate(unemployment_checks = ifelse((date >= ymd("2020-03-29") & date <=
    ymd("2020-07-31")),
    600, ifelse((date >= ymd("2020-07-26") & date <= ymd("2020-09-05")),
    300,
    ifelse((date >= ymd("2020-12-26") & date <= ymd("2021-03-13")),
    300,
    ifelse((date >= ymd("2021-03-14") & date <= ymd("2021-06-21")
      & (statename %in%
        states_which_stopped_early_benefits)), 300, ifelse((date >=
        ymd("2021-03-14") &
        date <= ymd("2021-09-04") & !(statename %in% states_which_
        stopped_early_benefits)),
        300, 0))))))

regression2_df$contclaims_rate_combined <- ifelse(regression2_df$contclaims_
  rate_combined ==
  ".", NA, regression2_df$contclaims_rate_combined)
regression2_df$contclaims_rate_combined <- as.double(regression2_df$contclaims
  _rate_combined)

```

Training the First Linear Model on Our Data

In this code chunk, we fit a linear model to the *gps_retail_and_recreation*, *emp*, and *stimulus_checks* features, with the goal of predicting the *spend_all* variable. The working assumption here is that *spend_all* is a dependent variable, and the others are independent variables.

The reason we chose these three input variables is that *stimulus_checks* are directly relevant to our analysis. We are trying to identify if these features had a positive or negative effect on the economy. These two features, *gps_retail_and_recreation* and *emp*, were added to help account for change in spending that could not necessarily be accounted for by the *stimulus_checks*. We chose not to include any other input features because we believe that *gps_retail_and_recreation* and *emp* account for much of the possible variance without over-complicating our linear model.

```

unemployment_regression <- lm(contclaims_rate_combined ~ bg_posts + visits_to_
  retail_and_recreation +
  unemployment_checks, regression2_df)
summary(unemployment_regression)

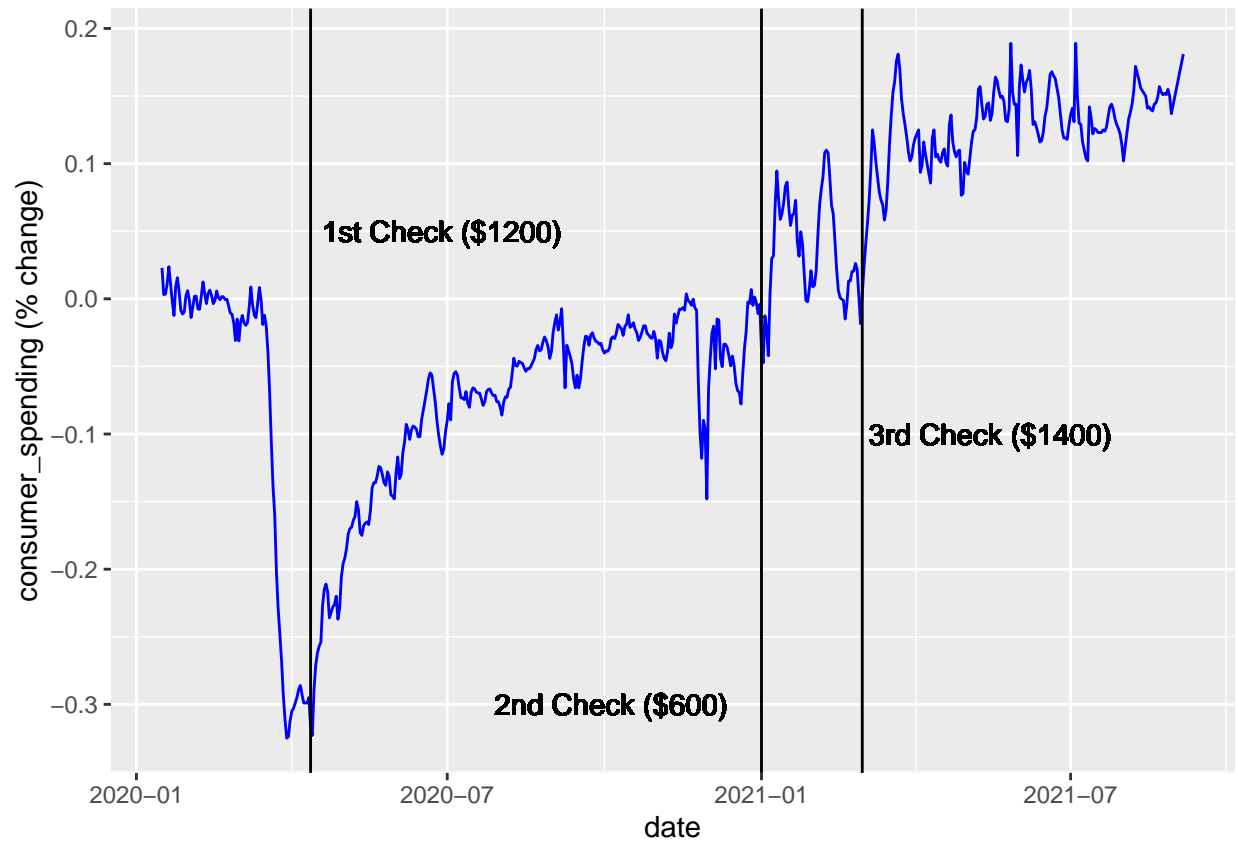
```

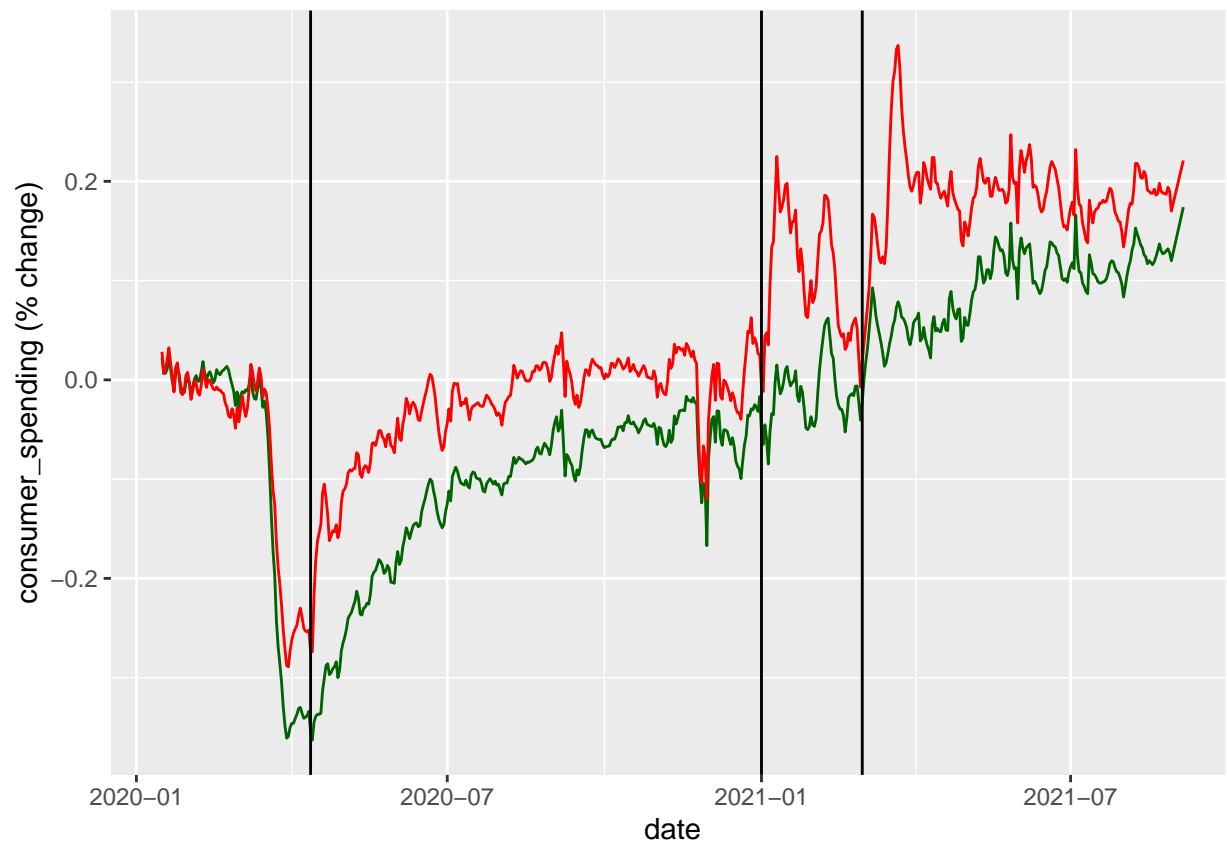
```
##
## Call:
## lm(formula = contclaims_rate_combined ~ bg_posts +
##     visits_to_retail_and_recreation +
##     unemployment_checks, data = regression2_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -16.908  -4.027  -1.171   2.924  57.707
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      5.744725    0.181844    31.59  <2e-16 ***
## bg_posts         -4.937473    0.488912   -10.10  <2e-16 ***
## visits_to_retail_and_recreation -12.011443    0.828006   -14.51  <2e-16 ***
## unemployment_checks    0.008454    0.000505    16.74  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.608 on 3872 degrees of freedom
## (612 observations deleted due to missingness)
## Multiple R-squared:  0.2401, Adjusted R-squared:  0.2395
## F-statistic: 407.9 on 3 and 3872 DF, p-value: < 2.2e-16
```

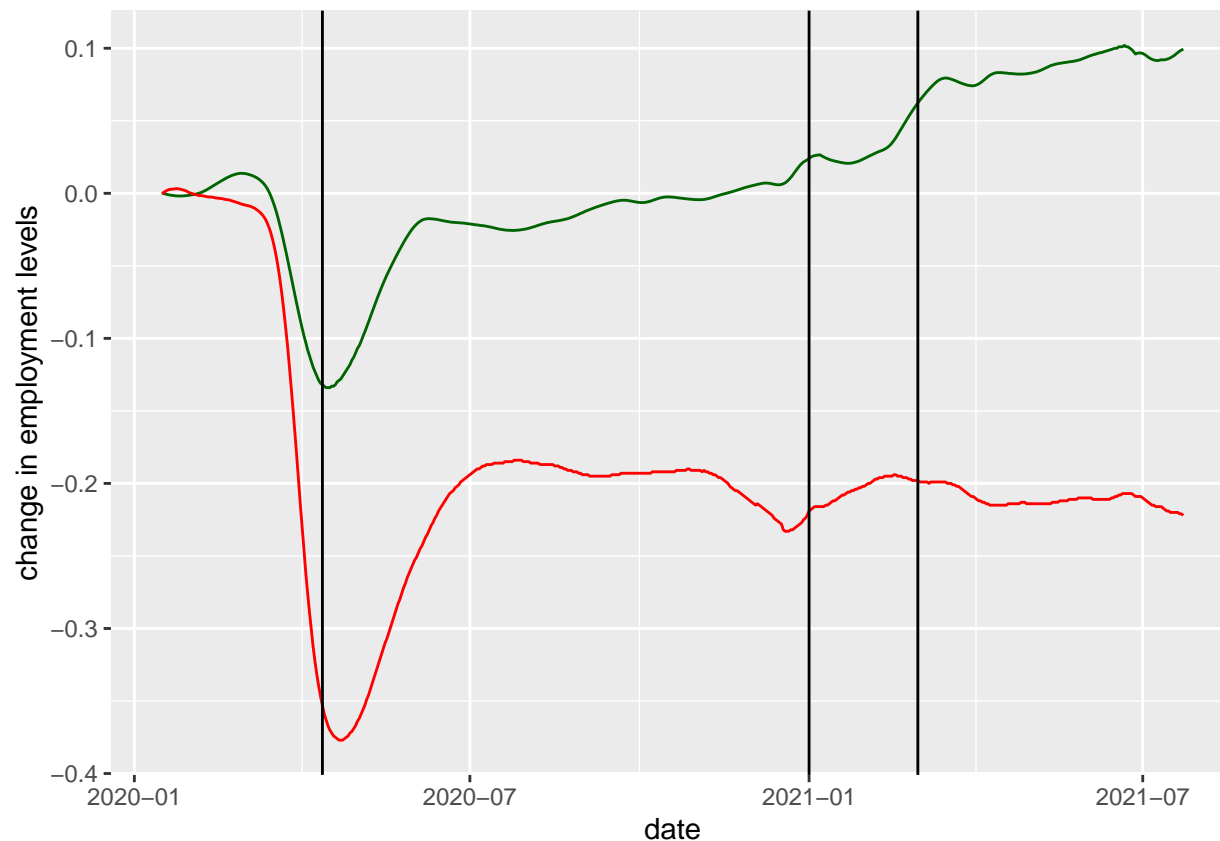
Plotting spending over time for all states and categories

(The dates for the stimulus checks were approximated from this article.)

```
## Warning in x:y: numerical expression has 2 elements: only the first used
## Warning in x:y: numerical expression has 2 elements: only the first used
## Warning in mask$eval_all_mutate(quo): NAs introduced by coercion
## Warning in mask$eval_all_mutate(quo): NAs introduced by coercion
## Warning in mask$eval_all_mutate(quo): NAs introduced by coercion
```







Summarize the Models

Based on the information displayed above, all of the variables we are regressing on are statistically significant for predicting the overall spending.

Our model is fairly limited in how accurately it predicts overall consumer spending because it is a linear model, and because we've limited ourselves to regressing on five variables for the sake of model simplicity, rather than the 30+ that we could've regressed on.

Future analysis could be done using neural networks, gradient boosted decision trees, or recurrent neural networks. We believe that all of these would be able to learn the nuances of our dataset better than a linear model could.

Summarize Our Biases

We were expecting to find that the stimulus checks all had a positive impact on consumer spending, because this is the only outcome that makes any economic sense, as far as we can tell. This effected how we encoded our stimulus check data into feature variables, and it effected how accurate we thought any given model was. If a model said that a stimulus check caused a decrease in consumer spending, we considered that model to be almost certainly highly inaccurate.

Summarize Biases in the Data

Our data contains several potential biases that could skew our results. For example, we use movement data from Google. This data was likely pulled from android phones, not Apple products. Imagine, then, that

more affluent people tend to purchase Apple products. If this is true, our movement data will be skewed away from affluent people who may be more likely to spend more, more casually.

Conclusion

Result Summary

The coefficients of our linear model seem to indicate that, unsurprisingly, giving out stimulus checks is correlated with an increase in overall consumer spending.

Therefore, *if* high consumer spending is correlated with a healthier economy, then it seems reasonable to conclude that giving out stimulus checks is also good for the economy, and if they are good for the economy, they are also good for our chances of getting jobs that pay well post-graduation.

The “if” at the beginning of the previous sentence is a big one, though. If the government has to go further into debt to send out stimulus checks, then is the net effect of stimulus checks really good for the economy? We don’t know, and this analysis has not considered the effect of increased government debt on the economy. This would be an interesting topic for future analysis.