

STAT GR5206 Homework 3 [100 pts]

Due 8:00pm Monday, October 30th on Canvas

Your homework should be submitted on Canvas using RMarkdown. Please submit both a knitted .pdf file and a raw .Rmd file. (If you are having trouble knitting to .pdf come to office hours and we'll try to sort it out, but for the homework, knit to .html and then convert to .pdf before handing it in). We will not (and cannot) accept any other formats. Please clearly label the questions in your responses and support your answers by textual explanations and the code you use to produce the result. Note that you cannot answer the questions by observing the data in the “Environment” section of RStudio or in Excel – you must use coded commands.

Goals: writing functions to automate repetitive tasks and using them as larger parts of code, some practice with ggplot, working with data frames and manipulating data from one form to another.

This homework uses the World Top Incomes Database and the Pareto distribution, as in this week's lab. The following notes are a repeat from the lab assignment:

In this lab we look at dataset containing information on the world's richest people from the World Top Incomes Database (WTID) hosted by the Paris School of Economics [<http://wid.world>]. This is derived from income tax reports, and compiles information about the very highest incomes in various countries over time, trying as hard as possible to produce numbers that are comparable across time and space.

For most countries in most time periods, the upper end of the income distribution roughly follows a Pareto distribution, with probability density function

$$f(x) = \frac{(a-1)}{x_{\min}} \left(\frac{x}{x_{\min}} \right)^{-a}$$

*for incomes $X \geq x_{\min}$. (Typically, x_{\min} is large enough that only the richest 3%-4% of the population falls above it.) As the **Pareto exponent**, a , gets smaller, the distribution of income becomes more unequal, that is, more of the population's total income is concentrated among the very richest people.*

The proportion of people whose income is at least x_{\min} and whose income is also at or above any level $w \geq x_{\min}$ is thus

$$\Pr(X \geq w) = \int_w^\infty f(x)dx = \int_w^\infty \frac{(a-1)}{x_{\min}} \left(\frac{x}{x_{\min}} \right)^{-a} dx = \left(\frac{w}{x_{\min}} \right)^{-a+1}.$$

We will use this to estimate how income inequality changed in the US over the last hundred years or so. (Whether the trends are good or bad or a mix is beyond our scope here.) WTID exports its data sets as .xlsx spreadsheets. For this lab session, we have extracted the relevant data and saved it as *wtid-report.csv*.

Part 1: Estimating a on US data

- i. In lab we use the fact that we can estimate the exponent using the following formula:

$$(1) \quad a = 1 - \frac{\log 10}{\log \left(\frac{P99}{P99.9} \right)}.$$

The logic leading to (1) also implies that

$$\left(\frac{P99.5}{P99.9} \right)^{-a+1} = 5$$

Write a function which takes $P99.5$, $P99.9$, and a , and calculates the left-hand side of that equation. Plot the values for each year using ggplot, using the data and your estimates of the exponent from lab (using the `exponent.est_ratio()`). Add a horizontal line with vertical coordinate 5. How good is the fit?

- ii. By parallel reasoning, we should have $(P99/P99.5)^{-a+1} = 2$. Repeat the previous step with this formula. How would you describe this fit compared to the previous one?
- iii. We have shown that if the upper tail of the income distribution followed a perfect Pareto distribution, then

$$(2) \quad \left(\frac{P99}{P99.9} \right)^{-a+1} = 10$$

$$(3) \quad \left(\frac{P99.5}{P99.9} \right)^{-a+1} = 5$$

$$(4) \quad \left(\frac{P99}{P99.5} \right)^{-a+1} = 2$$

We could estimate the Pareto exponent by solving any one of these equations for a ; we did this in lab and in the previous two questions. Because of measurement error and sampling noise, we can't find one value of a which will work for all three equations (2) - (4). Generally, trying to make all three equations come close to balancing gives a better estimate of a than just solving one of them. (This is analogous to finding the slope and intercept of a regression line by trying to come close to all the points in a scatterplot, and not just running a line through two of them.)

We will therefore estimate a by minimizing

$$\left(\left(\frac{P99}{P99.9} \right)^{-a+1} - 10 \right)^2 + \left(\left(\frac{P99.5}{P99.9} \right)^{-a+1} - 5 \right)^2 + \left(\left(\frac{P99}{P99.5} \right)^{-a+1} - 2 \right)^2.$$

Write a function, `percentile_ratio_discrepancies`, which takes as inputs `P99`, `P99.5`, `P99.9` and a , and returns the value of the expression above. Check that when `P99=1e6`, `P99.5=2e6`, `P99.9=1e7` and $a = 2$, your function returns 0.

- iv. Now we'd like to write a function, `exponent_multi_ratios_est`, which takes as inputs the vectors `P99`, `P99.5`, `P99.9`, and estimates a . It should minimize the function `percentile_ratio_discrepancies` you wrote above.

Recall that in class we used gradient descent to minimize the mean squared error (MSE) of a model fit depending on a parameter β to data. For the gradient descent algorithm, we approximated the derivative of the MSE, and adjusted our estimate of β by an amount proportional (and opposite) to that approximation. We stopped the algorithm when the derivative became small (assuming, then, that we were near a minimum). For this homework we will use a built-in R optimization function to do the minimization, which essentially does a fancier version of what we did in class.

R has several built-in functions for optimization, and one of the simplest,

which we use today, is `nlm()`, or non-linear minimization. `nlm()` takes two required arguments, a function to minimize and an initial value for the parameter.

For this problem write a function, `exponent.multi_ratios_est`, which takes as inputs the vectors `P99`, `P99.5`, `P99.9`, and estimates a by minimizing the function `percentile_ratio_discrepancies`. The initial value for the minimization should come from the estimate of a given by (1). Check that when `P99=1e6`, `P99.5=2e6` and `P99.9=1e7`, your function returns an a of 2.

- v. Write a function which uses `exponent.multi_ratios_est` to estimate a for the US for every year from 1913 to 2015. (There are many ways you could do this, including loops.) Plot the estimates using `ggplot`; make sure the labels of the plot are appropriate.
- vi. Use (1) to estimate a for the US for every year. Make a scatter-plot of these estimates against those from problem (v) using `ggplot`. If they are identical or completely independent, something is wrong with at least one part of your code. Otherwise, can you say anything about how the two estimates compare?

Part 2: Data for Other Countries

We're now going to look at this same data for some other countries: Canada, China, Colombia, Germany, India, Italy, Japan, South Africa, and Sweden. This data is in the file `wtid-homework.csv`.

- vii. We're now going to look at this same data for some other countries: Canada, China, Colombia, India, Italy, Japan, and Sweden. This data is in the file `wtid-homework.csv`. The WTID website also has data on the average income per “tax unit” (roughly, household) for the US and the other countries. This info is stored in the `AverageIncome` column.

Use your function from problem (v) to estimate a over time for each

of them. Note that the size of the dataset is different for each of these countries, and there may be some NA values.

- viii. Plot your estimates of a over time for all the countries using `ggplot`. Note that the years covered by the data are different for each country. You may either make multiple plots, or put all the series into one plot. Either way, make sure that the plots are clearly labeled.
- ix. Plot the series of average income per “tax unit” for the US and the countries against time in `ggplot`.
- x. The most influential hypothesis about how inequality is linked to economic growth is the “U-curve” hypothesis proposed by the great economist Simon Kuznets in the 1950s. According to this idea, inequality rises during the early, industrializing phases of economic growth, but then declines as growth continues.

Make a scatter-plot of your estimated exponents for the US against the average income for the US in `ggplot`. Qualitatively, can you say anything about the Kuznets curve? (Remember that smaller exponents indicate more income inequality.)

- xi. For a more quantitative check on the Kuznets hypothesis, use `lm()` to regress your estimated exponents on the average income, including a quadratic term for income. Are the coefficients you get consistent with the hypothesis? Hint: the following will regress y on both x and x^2 :

```
lm(y ~ x + I(x^2))
```

- xii. Do a separate quadratic regression for each country. Which ones have estimates compatible with the hypothesis? Hint: Write a function to fit the model to the data for an arbitrary country.

(If we were doing a more rigorous check of the Kuznet hypothesis, we would want to control for other factors, and not just assume that a quadratic was the right functional form for the curve.)

Please submit the knitted .pdf file!