

Hw_2

Assignment 2

- Overview of Hw 2 using R, Python and STATA

A.1 Import and create folders

A.1.1 Import libraries

```
1 # R
2 library(tidyverse)
3 library(foreign)
4 library(dplyr)
```

```
1 # Python
2 import re
3 import os
4 import sys
5 import numpy as np
6 import pandas as pd
7 import matplotlib.pyplot as plt
```

```
1 * Stata
2 * Not necessary for stata
```

A.1.2 Set or change working directory

```
1 # R
2 path = getwd()
```

```
1 # Python
2 path = os.getcwd() #Main folder
3 files_folder = path+'/hh02dta_bc' #DTA files
4
```

```
1 * Stata
2
```

A.2 Calculation consumption

A.2.1 Weekly variables

A.2.1.1 Import weekly variables data "i_cs"

```
1 # R
2 df1 = read.dta(paste(path, '/hh02dta_b1/i_cs.dta', sep=''))
3 df1 = as_tibble(df1)
4 df1 = df1 %>% rename(Household_id = folio)
```

```
1 # Python
2 path = os.getcwd() #Main folder
3 df1 = pd.read_stata(path+'/hh02dta_b1/i_cs.dta')
4 df1 = df1.rename(columns={'folio': 'Household_id'})
```

```
1 * Stata
2 use "hh02dta_b1\i_cs.dta", clear
3
```

A.2.1.2 Select weekly consumption variables

A.2.1.2.1 Convert weekly variables to monthly vars

```
1 # R
2 #A.2.1.2 Select weekly consumption variables
3 cons1 = df1 %>%
4   select(Household_id,
5           intersect(contains('cs02a_'), ends_with('2'))))
6
7 # A.2.1.2.1 Convert weekly variables to monthly vars
8 cons1 = cons1 %>%
9   mutate_at(vars(cs02a_12:cs02a_82), .funs = funs(. *4.3))
```

```
1 # Python
2 #A.2.1.2 Select weekly consumption variables
3 cons1 = [i for i in df1.columns if i.startswith('cs02a_') and i.endswith('2')]
4 cons1.append('Household_id')
5 cons1 = df1[cons1]
6 cons1.set_index('Household_id', inplace=True)
7
8 # A.2.1.2.1 Convert weekly variables to monthly vars
9 cons1 = cons1*4.3
```

```
1 * Stata
2 local weeklyvar "cs02a_12 cs02a_22 cs02a_32 cs02a_42 cs02a_52 cs02a_62 cs02a_72 cs02a_82"
3 foreach x in `weeklyvar' {
4   generate ics_monthly_`x' = `x' * 4.3
5 }
```

A.2.2 Monthly variables

```
1 # R
2 cons2 = df1 %>%
3   select(Household_id,
4     ..... intersect(contains('cs16'), ends_with('2'))))

1 # Python
2 cons2 = [i for i in df1.columns if i.startswith('cs16') and i.endswith('2')]
3 cons2.append('Household_id')
4 cons2 = df1[cons2]
5 cons2.set_index('Household_id', inplace=True)

1 * Stata
2 local 1month_var "cs16a_2 cs16b_2 cs16c_2 cs16d_2 cs16e_2 cs16f_2 cs16g_2 cs16h_2 cs16i_2"
3 foreach x in `1month_var' {
4   generate ics_monthly_`x' = `x'
5 }
```

A.2.3 3-Month variables

A.2.3.1 Import weekly variables data "i_cs1"

A.2.3.2 Select 3-month variables

A.2.3.2.1 Convert 3-month variables to monthly vars

```
1 # R
2 ##### A.2.3.1 Import weekly variables data "i_cs1"
3 df2 = read.dta(paste(path, '/hh02dta_b1/i_cs1.dta', sep=''))
4 df2 = as_tibble(df2)
5 df2 = df2 %>% rename(Household_id = folio) #Rename
6
7 # A.2.3.2 Select 3-month variables
8 cons3 = df2 %>%
9   select(Household_id,
10     ..... intersect(contains('cs22'), ends_with('2'))))
11
12 # A.2.3.2.1 Convert 3-month variables to monthly vars
13 cons3 = cons3 %>%
14   mutate_at(vars(cs22a_2:cs22h_2), .funs = funs(. /3))

1 # Python
2 # A.2.3.1 Import weekly variables data "i_cs1"
3 df2 = pd.read_stata(path+'/hh02dta_b1/i_cs1.dta')
4 df2 = df2.rename(columns={'folio': 'Household_id' })
5
6 # A.2.3.2 Select 3-month variables
7 cons3 = [i for i in df2.columns if i.startswith('cs22') and i.endswith('2')]
8 cons3.append('Household_id')
```

```

9 cons3 = df2[cons3]
10 cons3.set_index('Household_id', inplace=True)
11
12 # A.2.3.2.1 Select 3-month variables
13 cons3 = cons3/3

1 * Stata
2 * A.2.3.1 Import weekly variables data "i_cs1"
3 merge m:m folio using "hh02dta_b1\i_cs1.dta"
4
5 * A.2.3.2.1 Select 3-month variables
6 local 3month_var "cs22a_2 cs22b_2 cs22c_2 cs22d_2 cs22e_2 cs22f_2 cs22g_2 cs22h_2"
7
8 foreach x in `3month_var' {
9 generate ics1_monthly_`x' = `x' /3
10 }
11

```

A.2.4 Merge consumption variables into one dataframe

```

1 # R
2 merge1 = merge(cons1, cons2, by = 'Household_id', all= TRUE)
3 cons_merge = merge(merge1, cons3, by = 'Household_id', all= TRUE)
4 head(cons_merge,3)

```

```

1 # Python
2 cons_merge = pd.concat([cons1, cons2, cons3], axis=1)

```

```

1 * Stata
2 * Already merged

```

Homework Questions

Q.1.

- Calculate a measure of total consumption and per capita consumption for each household in the 2002 round.
 - To calculate per capita note you will have to calculate the number of individuals in each household.

Q.1.1 Total consumption

1.1.1 Calculate total consumer spending

- Sum the rows from cons_merge (A.2.4)

```

1 # R
2 total_cons = cons_merge %>%
3   mutate(

```

```

4 consumption = rowSums(., cs02a_12:cs22h_2)
5 )%>% select(Household_id, consumption)

```

```

1 # Python
2 cons_merge['total_cons'] = cons_merge.sum(axis=1)

```

```

1 * Stata
2 egen total_cons = rowtotal(ics1_monthly_cs22a_2 ics1_monthly_cs22b_2 ics1_monthly_cs22c_2
  ics1_monthly_cs22d_2 ics1_monthly_cs22e_2 ics1_monthly_cs22f_2 ics1_monthly_cs22g_2
  ics1_monthly_cs22h_2 ics1_monthly_cs02a_12 ics1_monthly_cs02a_22 ics1_monthly_cs02a_32
  ics1_monthly_cs02a_42 ics1_monthly_cs02a_52 ics1_monthly_cs02a_62 ics1_monthly_cs02a_72
  ics1_monthly_cs02a_82 ics1_monthly_cs16a_2 ics1_monthly_cs16b_2 ics1_monthly_cs16c_2
  ics1_monthly_cs16d_2 ics1_monthly_cs16e_2 ics1_monthly_cs16f_2 ics1_monthly_cs16g_2
  ics1_monthly_cs16h_2 ics1_monthly_cs16i_2)

```

1.1.1.1 Basic stats: total consumer spending

```

1 # R
2 summary(total_cons$consumption)

```

```

1 # Python
2 cons_merge['total_cons'].describe()

```

```

1 * Stata
2 sum total_cons

```

1.1.1.2 Graph: total consumer spending

```

1 # R
2 hist(total_cons$consumption)

```

```

1 # Python
2 from scipy import stats
3 #Calculates Z-scores
4 z = np.abs(stats.zscore(cons_merge['total_cons']))
5 threshold = 3
6
7 outlier_removed = cons_merge['total_cons'][(z < 5)].values
8 import seaborn as sns
9 ax = sns.distplot(outlier_removed, hist=True, kde=True,
10                  bins=int(180/5), color = 'darkblue',
11                  hist_kws={'edgecolor':'black'},
12                  kde_kws={'linewidth': 3})
13 ax.set_xlabel('Income')
14 ax.set_ylabel('Pct of households')
15 ax.set_title('Kernal density plot of household total income')
16 plt.show()

```

```

1 * Stata
2 egen total_cons_sd =std(total_cons)

```

```
3 histogram total_cons if total_cons_sd<3, bin(10)
```

Q.1.2 Per capita consumption (Total/house size)

1.2.1 Avg. house size

1.2.1.1 Import house size data set "c_ls"

1.2.1.2 Count family members "ls" in each household "folio"

```
1 # R
2 # 1.2.1.1 Import house size data set "c_ls"
3 df_housesize = read.dta(paste(path, '/hh02dta_bc/c_ls.dta', sep=''))
4 df_housesize = as_tibble(df_housesize) #Set as tibble
5 head(df_housesize,2)
6
7 # 1.2.1.2 Count family members "ls" in each household "folio"
8 family_members = df_housesize %>% group_by(folio) %>% count(folio)
9 family_members = family_members%>% rename(family_members = n, Household_id=folio)
10 head(family_members)
```

```
1 # Python
2 # 1.2.1.1 Import house size data set "c_ls"
3 df_housesize = pd.read_stata(path+'/'hh02dta_bc/c_ls.dta')
4 df_housesize = df_housesize.rename(columns={'folio':'Household_id'})
5 df_housesize.set_index('Household_id', inplace=True)
6 merge = pd.merge(df_housesize, cons_merge, left_index=True, right_index=True)
7
8 # 1.2.1.2 Count family members "ls" in each household "folio"
9 merge['family_members'] = merge.groupby(merge.index.get_level_values(0))['ls'].count()
10 merge.family_members.head(5)
```

```
1 * Stata
2 * 1.2.1.1 Import house size data set "c_ls"
3 drop _merge
4 merge m:m folio using "hh02dta_bc\c_ls.dta"
5
6 *1.2.1.2 Count family members "ls" in each household "folio"
7 egen family_members = count(ls), by(folio)
8 histogram family_members, bin(10) title(Family members) xtitle(# of family members)
```

1.3 Per capita consumption (Total/house size)

1.3.1 Merge family_members and consumption data set

```
1 # R
2 percap_consum = merge(family_members, total_cons, by='Household_id')
```

```
1 # Python
```

```

2 merge['percap_consum'] = (merge.total_cons)/(merge.family_members)
3
1 * Stata
2

```

1.3.1.3 Calculate percapita consumption

```

1 # R
2 percap_consum = percap_consum %>%
3 mutate(percap_consum = percap_consum$consumption/percap_consum$family_members)
4 percap_consum%>% head(3)

1 # Python
2 merge['percap_consum'] = (merge.total_cons)/(merge.family_members)

1 * Stata
2 gen percap_consum = total_cons / family_members
3 egen percap_consum_sd = std(percap_consum)
4 histogram percap_consum if percap_consum_sd <3, bin(100)

```

Q.2

- Calculate the set of poverty rates nationwide using the FGT indicators of poverty:
- Head count
- Average poverty gap
- Average poverty gap squared.
- Assume the poverty line=500 pesos per person. Provide poverty rates based on household consumption per capita.

2.1 Headcount using 1000 as an example

```

1 # R
2 povertyline = 500
3
4 below_poverty = percap_consum %>%filter(percap_consum < povertyline )%>%
5 select(percap_consum)
6 observations = length(percap_consum$percap_consum)
7 head_count = length(below_poverty$percap_consum)/observations
8 print(paste('Headcount: ', round(head_count*100, 2), '%', sep=''), quote=FALSE)

1 # Python
2 povertyline = 500
3 merge['below_poverty_dummy'] = merge.percap_consum< int(povertyline)

```

```

4
5 head_count = merge['below_poverty_dummy'].mean()
6 print('Headcount: {}'.format(head_count))

1 * Stata
2 gen poverty_line = 500
3 gen p0 =percap_consum <poverty_line //Poor households dummy
4 la var p0 "poverty incidence"
5 la de p0 0 "non-poor" 1 "poor"
6 sum p0
7 local p0_tot=r(mean)
8 display `p0_tot'

```

2.2 Avg. poverty gap

```

1 # R
2 (povertyline-below_poverty$percap_consum)*head_count

1 # Python
2 ((povertyline-merge['percap_consum'])*merge['below_poverty_dummy']).mean()

1 * Stata
2 gen p1=(poverty_line-percap_consum)*p0 //calculate the poverty gap for each household
3 la var p1 "poverty gap"
4 gen p1z= p1/poverty_line
5
6 sum p1z
7 local p1z_tot=r(mean)
8 display `p1z_tot'

```

2.3 Avg. poverty gap squared

```

1 # R
2 mean((povertyline-below_poverty$percap_consum)**2)

1 # Python
2 (((povertyline-merge['percap_consum'])*merge['below_poverty_dummy'])**2).mean()

1 * Stata
2 gen p2=(poverty_line-percap_consum)^2*p0
3 la var p2 "poverty gap squared"
4 gen p2z=p2/(poverty_line^2)
5
6 sum p2z
7 local p2z_tot=r(mean)
8 display `p2z_tot'

```

Q.3 Repeat 2) by area of residence.

- How does the poverty rate change by rural/urban residence?

3.1. Import residence data from "c_portad"

3.1.1. Merge residence df with percap_consum df from Q.1 & Q.2

```

1 # R
2 # 3.1. Import residence data from "c_portad"
3 residence_df = read.dta(paste(path, '/hh02dta_bc/c_portad.dta', sep=''))
4 residence_df = residence_df %>% rename(Household_id = folio)
5 head(residence_df)
6
7 # 3.1.1. Merge residence df with percap_consum df from Q.1 & Q.2
8 consum_residence_df = merge(residence_df, percap_consum, by='Household_id')
9 consum_residence_df = as_tibble(consum_residence_df)
10 consum_residence_df %>% head(3)

1 # Python
2 # 3.1
3 residence_df = pd.read_stata(path+'/hh02dta_bc/c_portad.dta')
4 residence_df = residence_df.rename(columns={'folio': 'Household_id'})
5 residence_df.set_index('Household_id', inplace=True)
6
7 # 3.1.1. Merge residence df with percap_consum df from Q.1 & Q.2
8 consum_residence_df = pd.merge(residence_df, merge, left_index=True, right_index=True)
9

1 * Stata
2 * 3.1. Import residence data from "c_portad"
3 drop _merge
4 merge m:m folio using "hh02dta_bc\c_portad.dta"

```

3.1.1.2 Create poverty dummy

```

1 # R
2 consum_residence_df$poverty_dummy <-
  as.numeric(consum_residence_df$percap_consum < poverty_line)
3 consum_residence_df %>% filter(poverty_dummy == 1) #Show observations living in poverty

1 # Python
2 #Created previously

1 * Stata
2 *Created in 2.1 (Note don't run again)
3 gen poverty_line = 500
4 gen p0 = percap_consum < poverty_line //Poor households dummy
5 la var p0 "poverty incidence"
6 la de p0 0 "non-poor" 1 "poor"

```

3.2 Show poverty by area of residence

```
1 # R
2 consum_residence_df %>%
3 group_by(estrato) %>% #Groupby estrato
4 count(poverty_dummy) %>% #Count poverty dummy by each estrato
5 mutate(obs = sum(n))%>%
6 mutate(head_count = (n/obs)*100)%>%
7 filter(poverty_dummy==1) %>%
8 select(estrato, head_count)
```

```
1 # Python
2 consum_residence_df.groupby('estrato')['below_poverty_dummy'].mean()
```

```
1 * Stata
2 tab estrato p0, row
```

Q.4.

Calculate the Gini coefficient overall and by urban and rural areas both using consumption measures. In both R and Stata there are programs which can be downloaded to calculate the Gini. Provide a graph of the Lorenz curve.

4.1 Calculate cumulative sum for population and consumption

Note: Python code runs but does not produce the correct results

```
1 # R
2 gini = percap_consum %>%
3 arrange(percap_consum) %>% #Sort consumption from least to greatest
4 mutate(consum_cumulative = cumsum(percap_consum), consum_total = sum(percap_consum),
5        pop_total = sum(family_members), pop_cumulative = cumsum(family_members))%>% #Calculate
6        total and cumulative sum for variables
7 mutate(consum_pct = (consum_cumulative/consum_total), pop_pct= ((pop_cumulative/
8        pop_total))) #Calculate quintiles
9 gini %>% head(3) #Show new data
```

```
1 # Python
2 var_list = ['percap_consum', 'family_members']
3 for i in var_list:
4     consum_residence_df[str(i)+'_cumulative'] = consum_residence_df[str(i)].cumsum()
5     consum_residence_df[str(i)+'_cumulative_pct'] =
6     consum_residence_df[str(i)+'_cumulative']/sum(consum_residence_df[str(i)])
```

```
1 * Stata
2 sort percap_consum
3 gen rank_percap_consum =_n
4 cumul percap_consum, generate(cum_percap_consum)
```

```

5 egen total_percap_consum =total(percapi_consum)
6 sort percap_consum
7
8
9 xtile q5=percapi_consum, n(5) //create the variable that indicates the quintiles
10 sum percap_consum
11 scalar percap_consum_sum=r(sum)
12 forvalues i = 1/5 {
13     quietly sum percap_consum if q5==`i'
14     scalar percap_consum_`i'_sum=r(sum)
15     scalar percap_consum_`i'_share=percapi_consum_`i'_sum/percapi_consum_sum
16 }
17 scalar list percap_consum_1_share percap_consum_2_share percap_consum_3_share
    percap_consum_4_share percap_consum_5_share
18
19
20 gen cum_total_conspc=0
21 replace cum_total_conspc= percap_consum if _n==1
22 local N=_N
23 forvalues k=2/`N' {
24     quietly replace cum_total_conspc = percap_consum+cum_total_conspc[_n-1] if _n==`k'
25 }
26 gen cuml_conspc=cum_total_conspc/total_percap_consum //create the variable for the
    cumulative consumption share
27
28 gen pop_share = rank_percapi_consum/35677

```

4.1.1 Plot:

```

1 # R
2 library(ggplot2)
3
4 ggplot(data= gini, aes(x=pop_pct, y= consum_pct))+
5 ggtitle("Lorenz curve")+theme(plot.title = element_text(hjust = 0.5))+
6 geom_line()+
7 geom_abline(intercept = 0, slope = 1, color='red')+xlab('cum. % of households')+ylab('cum. %
    consum/percap')

```

```

1 # Python
2 plt.plot(consum_residence_df['percapi_consum_cumulative_pct'],consum_residence_df['family_mem
    bers_cumulative_pct'])

```

```

1 * Stata
2 line cuml_conspc pop_share, title("lorenz curve") xtitle("cum. % of households")
    ytitle("cum. % of consumption per capita")
3

```

4.2 Gini coefficient calculation

```
1 # R
2 cov_consum_V_consum_pct = cov(gini$percap_consum, gini$consum_pct)
3 mean_cons = mean(gini$percap_consum)
4 print((2*cov_consum_V_consum_pct)/(mean_cons))
```

```
1 # Python
2 cov_consum_V_consum_pct =
  consum_residence_df[['percap_consum', 'percap_consum_cumulative_pct']].cov().iloc[0,1:]
3 mean_cons = np.mean(consum_residence_df['percap_consum'])
4 print((2*cov_consum_V_consum_pct)/mean_cons)
```

```
1 * Stata
2 correlate percap_consum cuml_conspc, covariance
3 scalar conspc_cov=r(cov_12) //Calculate the covariance
4 su percap_consum
5 scalar conspc_mean=r(mean) //Calculate the mean
6 scalar gini=2*conspc_cov/conspc_mean //Gini coefficient
7 display gini
```