

IPUMS – HigherEd Extraction and Analysis

Exercise 2 - SAS

OBJECTIVE: Gain an understanding of how an IPUMS dataset is structured and how it can be leveraged to explore your research interests. This exercise will use the IPUMS dataset to explore the factors that affect student loan debt and the relatedness between highest degree and work activities.

Research Questions

What are the most common work activities for scientists and engineers in the US?
What factors are most important when predicting remaining student loan debt for college graduates?

Objectives

- Explore a sample of variables from IPUMS-HIGHER ED
- Analyze the data using example code

IPUMS Variables

- GENDER: Respondent's gender
- WAPRI: Primary work activity
- AGEP: Age
- MINRTY: Minority indicator
- GOVSUP: Indicator of government support for work
- DGRDG: Degree type of respondent's highest degree
- BAFLN: Financial support during bachelor's degree included loans
- UGBALPB: Balance left on undergraduate student loans
- UGLOANPB: Total amount of undergraduate student loans

SAS Code to Review

Code	Purpose
proc freq;	Begins a frequency procedure
proc means;	Begins a means procedure, returns the mean value of a variable
tables	Required syntax to display frequencies
where	Selects only specified cases to include in a procedure

Review Answer Key (page 6)

Common Mistakes to Avoid

- 1 Giving the wrong filepath to indicate the dataset
- 2 Forget to close a procedure with "run;"
- 3 Forget to terminate a command with a semicolon ";"

Registering with IPUMS

Go to <http://highered.ipums.org>, click on "Register to Use IPUMS-HIGHER ED" and apply for access. On login screen, enter email address and password and submit it!

Step 1

Make an Extract

- Go back to homepage and go to Select Data
- Click the Select Samples box on the SESTAT tab
- Check the box labeled Select all samples – this will select all years of SESTAT samples. Click on Submit sample selections
- Using the drop down menu or search feature, select the following variables:

GENDER: Gender

MINRTY: Minority indicator

WAPRI: Primary work activity

AGEP: Age

GOVSUP: Indicator of government support for work

DGRDG: Degree type of respondent's highest degree

BAFLN: Loans as financial support for bachelor's degree

UGBALPB: Balance of undergraduate student loans

WEIGHT: SESTAT sample weight variable

UGLOANPB: Total amount of undergraduate student loans

...

Step 2

Request the Data

- Click the green VIEW CART button under your data cart
- Review variable selection. Click the green Create Data Extract button
- Review the 'Extract Request Summary' screen, describe your extract and click Submit Extract
- You will get an email when the data is available to download
- To get to the page to download the data, follow the link in the email, or follow the Download and Revise Extracts link on the homepage

Getting the data into your statistics software

The following instructions are for SAS. If you would like to use a different stats package, see: http://highered.ipums.org/highered/extract_instructions.shtml

Step 1

Download the Data

...

Step 2

Decompress the Data

...

Step 3

Read in the Data

- Go to <http://highered.ipums.org> and click on Download or Revise Extracts
 - Right-click on the data link next to extract you created
 - Choose "Save Target As..." (or "Save Link As...")
 - Save into "Documents" (that should pop up as the default location)
 - Do the same thing for the SAS link next to the extract
-
- Find the "Documents" folder under the Start menu.
 - Right click on the ".dat" file
 - Use your decompression software to extract here
 - Double-check that the Documents folder contains three files starting "highered_000..."
 - Free decompression software is available at <http://www.ironis.net/soft/wingzip/>
-
- Open the "highered_000##.sas" file.
 - In the do file window, change the first line from "libname IPUMS '.'" to "libname IPUMS '\\Documents...;" using the file directory where you saved your data files.
 - After "filename ASCIIDAT", enter the full file location, ending with "highered_000##.dat";
 - Choose Submit under the Run file menu.

Analyze the Sample – Part I Frequencies

Step 1

Analyze the Data

...

Step 2

Weighting the Data

A) On the website, find the codes page for the DGRDG variable and write down each code value, and what category each code represents. _____

B) Describe the kind of information WAPRI provides.

C) What were the top three most common work activities for SESTAT respondents in 2010? _____

```
proc freq; tables wapri; where year eq 2010 and wapri < 90; run;
```

D) What were the top three most common work activities for SESTAT respondents in 2010 whose highest degree was a doctorate? _____

```
proc freq; tables wapri; where year eq 2010 and wapri < 90 and  
dgrdg eq 3; run;
```

Using weights (WEIGHT)

In order to find a nationally representative estimate of doctorate recipients, we need to use a frequency weight. The above analysis finds frequencies specific to the survey samples, but a weight adjusts the analysis to be representative of the US target population for each year. The variable WEIGHT is specifically adjusted to analyze the entire SESTAT sample in a particular year.

A) How many doctorate recipients nationally had an occupation in which the primary work activity was teaching?

B) What proportion of doctorate recipients nationally had an occupation in which the primary work activity was basic research?

```
proc freq; tables wapri; where year eq 2010 and wapri < 90  
and dgrdg eq 3; weight weight; run;
```

Analyze the Sample – Part II Relationships in the Data

Section 1

Crosstabs

A) What is the universe for GOVSUP?

B) What percentage of the science and engineering workforce in the United States received funding from the federal government for their work? Note that we are excluding respondents who were out of universe for GOVSUP. _____

```
proc freq; tables govsup*year; where govsup < 90; weight weight; run;
```

C) Do doctorate or bachelor's degree holders tend to receive federal funding for their work? _____

```
proc freq; tables govsup*dgrdg; where govsup < 90 and year eq 2010; weight weight; run;
```

D) How many individuals in the science and engineering workforce in 1999 used loans from lending institutions to finance their bachelor's degree?

```
proc freq; tables bafln*year; where bafln < 90; weight weight; run;
```

Analyze the Sample – Part II Relationships in the Data

Section 2

Means

A) What is the average amount of remaining undergraduate loans for those who used loans to finance their undergraduate degree in 1999? _____

```
proc means; var ugbalpb; where ugbalpb < 999990 and year eq 1999; weight weight; run;
```

B) What is the average amount of total undergraduate loans borrowed for those who used loans to finance their undergraduate degree in 1999? _____

```
proc means; var ugloanpb; where ugloanpb < 999990 and ugloanpb > 0 and year eq 1999; weight weight; run;
```

Note: We exclude zeros from original loans, but not from the remaining balance. We want to include the people who paid off the loans, not the ones who never took out loans.

Regression

C) Now let's use some common variables to predict the value of the remaining undergraduate student loans of the science and engineering workforce in 1999. _____

```
proc reg; model ugbalpb = agep gender minrty; where ugbalpb lt 999990 and year eq 1999; weight weight; run;
```

...

Complete!
Check
your
Answers!

ANSWERS - Analyze the Sample – Part I Frequencies

Step 1

Analyze the Data

A) On the website, find the codes page for the DGRDG variable and write down each code value, and what category each code represents. 1 = Bachelor's 2 = Master's 3 = Doctorate 4 = Professional 5 = Other

B) Describe the kind of information WAPRI provides. The work activity the respondent performs for the majority of his/her job.

C) What were the top three most common work activities for SESTAT respondents in 2010? Management and Administration, Professional Services, Teaching

```
proc freq; tables wapri; where year eq 2010 and wapri < 90; run;
```

...

Step 2

Weighting the Data

D) What were the top three most common work activities for SESTAT respondents in 2010 whose highest degree was a doctorate? Teaching, Basic research, Management and Administration

```
proc freq; tables wapri; where year eq 2010 and wapri < 90 and dgrdg eq 3; run;
```

Using weights (WEIGHT)

In order to find a nationally representative estimate of doctorate recipients, we need to use a frequency weight. The above analysis finds frequencies specific to the survey samples, but a weight adjusts the analysis to be representative of the US target population for each year. The variable WEIGHT is specifically adjusted to analyze the entire SESTAT sample in a particular year.

A) How many doctorate recipients nationally had an occupation in which the primary work activity was teaching? 234,630

B) What proportion of doctorate recipients nationally had an occupation in which the primary work activity was basic research? 17.7%

```
proc freq; tables wapri; where year eq 2010 and wapri < 90 and dgrdg eq 3; weight weight; run;
```


ANSWERS - Analyze the Sample – Part II Relationships in the Data

Section 1

Crosstabs

A) What is the universe for GOVSUP in 2010? Worked during calendar year 2009

B) What percentage of the science and engineering workforce in the United States received funding from the federal government for their work in 2010? Note that we are excluding respondents who were out of universe for GOVSUP. 16.15%

```
proc freq; tables govsup*year; where govsup < 90; weight weight; run;
```

C) Do doctorate or bachelor's degree holders tend to receive federal funding for their work? Doctorates 30% (vs Bachelors, of which only 14.1% receive funding)

```
proc freq; tables govsup*dgrd; where govsup < 90 and year eq 2010; weight weight; run;
```

D) How many individuals in the science and engineering workforce in 1999 had used loans from lending institutions to finance their bachelor's degree? 48.95%

```
proc freq; tables bafln*year; where bafln < 90; weight weight; run;
```

ANSWERS - Analyze the Sample – Part II Relationships in the Data

Section 2

Means

A) What is the average amount of remaining undergraduate loans for those who used loans to finance their undergraduate degree in 1999? \$7088.90

```
proc means; var ugbalpb; where ugbalpb < 999990 and year eq 1999; weight weight; run;
```

B) What is the average amount of total undergraduate loans borrowed for those who used loans to finance their undergraduate degree in 1999? \$8948.18

```
proc means; var ugloanpb; where ugloanpb < 999990 and ugloanpb > 0 and year eq 1999; weight weight; run;
```

Note: We exclude zeros from original loans, but not from the remaining balance. We want to include the people who paid off the loans, not the ones who never took out loans.

Regression

C) Now let's use some common variables to predict the value of the remaining undergraduate student loans of the science and engineering workforce in 1999. (Coefficients below)

```
proc reg; model ugbalpb = agep gender minrty; where ugbalpb lt 999990 and year eq 1999; weight weight; run;
```

Variable	Coef	t-stat
Age	-134	-127
Gender	-568	-52
Minority	563	39
Constant	12008	319