

LAST UPDATE: OCTOBER 24, 2023

# Principles of Quantitative Project Workflows Management

*Data Management and Analysis*

DOCUMENT CREATED BY:

**USC**Dornsife  
*Center for Economic  
and Social Research*

635 Downey Way, VPD, Los Angeles, CA 90089  
<https://cesr.usc.edu/care>

# Contents

---

	Pg.
<b>1 Introduction</b>	<b>5</b>
<b>2 Overarching Goals and Rationale</b>	<b>5</b>
2.1 Some Inviolable Commandments . . . . .	5
2.2 Replicability . . . . .	6
2.3 Transferrability . . . . .	7
2.4 Accessibility . . . . .	8
<b>3 Expectations and Guiding Principles</b>	<b>8</b>
<b>4 Naming folders and files</b>	<b>9</b>
4.1 Folder Naming and Structure . . . . .	10
4.1.1 Common sub-folders . . . . .	10
4.1.2 Working across Project Years . . . . .	11
4.2 File Naming . . . . .	11
4.2.1 Syntax file naming . . . . .	12
<b>5 Naming variables, macros, and other characteristics in Stata datasets</b>	<b>13</b>
<b>6 USC-CARE's computing environment and tools</b>	<b>16</b>
6.1 Project-related storage . . . . .	16
<b>7 Researcher written Stata programs</b>	<b>17</b>
7.1 <code>preamble.ado</code> . . . . .	17
7.1.1 Author: Marshall Garland . . . . .	17
7.1.2 Description . . . . .	17
7.1.3 Usage and options . . . . .	17
7.2 <code>uas_preamble.ado</code> . . . . .	18
7.2.1 Author: Marshall Garland . . . . .	18
7.2.2 Description . . . . .	18
7.3 <code>varcount.ado</code> . . . . .	18
7.3.1 Author: Marshall Garland . . . . .	18
7.3.2 Description . . . . .	18
7.3.3 Usage and options . . . . .	18
7.4 <code>graphsout.ado</code> . . . . .	18
7.4.1 Author: Marshall Garland . . . . .	18
7.4.2 Description . . . . .	18
7.4.3 Usage and options . . . . .	18
7.5 <code>qualtrics.ado</code> . . . . .	19

7.5.1	Author: <b>Danial Hoepfner</b>	19
7.5.2	Description	19
7.5.3	Usage and options	19
7.6	gtab.ado	20
7.6.1	Author: <b>Danial Hoepfner</b>	20
7.6.2	Description	20
7.6.3	Usage and options	20
7.7	fmtexl.ado	23
7.7.1	Author: <b>Danial Hoepfner</b>	23
7.7.2	Description	23
7.7.3	Usage and options	24
7.8	apchx.ado	24
7.8.1	Author: <b>Danial Hoepfner</b>	24
7.8.2	Description	24
7.8.3	Usage and options	24
7.9	conwithin.ado	25
7.9.1	Author: <b>Danial Hoepfner</b>	25
7.9.2	Description	25
7.9.3	Usage and options	25
7.10	edrop.ado	25
7.10.1	Author: Eric Booth	25
7.10.2	Description	25
7.10.3	Usage and options	25
7.11	latest.ado	25
7.11.1	Author: <b>Danial Hoepfner</b>	25
7.11.2	Description	26
7.11.3	Usage and options	26
7.12	char_finder.ado	26
7.12.1	Author: <b>Danial Hoepfner</b>	26
7.12.2	Description	26
7.12.3	Usage and options	26
7.13	fixtex.ado	26
7.13.1	Author: <b>Danial Hoepfner</b>	26
7.13.2	Description	26
7.13.3	Usage and options	26
7.14	texpdf.ado	27
7.14.1	Author: <b>Danial Hoepfner</b>	27
7.14.2	Description	27
7.14.3	Usage and options	27
7.15	syntax_saver.ado	27
7.15.1	Author: <b>Danial Hoepfner</b>	27
7.15.2	Description	27
7.15.3	Usage and options	27

7.16 latout.ado . . . . .	27
7.16.1 Author: <b>Danial Hoepfner</b> . . . . .	27
7.16.2 Description . . . . .	28
7.16.3 Usage and options . . . . .	28

---

# 1 Introduction

The purpose of this document is to outline and codify some general principles for working with data on projects with a quantitative component. Below, we elaborate upon three principles that should systematically and consistently guide how each analyst approaches quantitative analyses. They are:

- Replicability
- Transferability
- Accessibility.

We explain the meaning and rationale for each of these goals below. And, then, we list some expectations and overarching principles. Finally, we discuss how to best organize information (that is, folder, file, and variable naming conventions to better achieve these principals/goals). At the end, we provide an introduction to the Center for Applied Research in Education's (CARE's) computing environment, and a quick synopsis of Stata packages written by contributing researchers that may automate or simplify some repetitive procedures that commonly arise on our projects.

This document is inspired, and borrows heavily, from similar guides for applied researchers, including Matthew Gentzkow and Jesse Shapiro's [Code and Data for the Social Sciences: A Practitioner's Guide](#) and Julian Reif's [Stata Coding Guide](#).<sup>1</sup> It has benefited from the input and contributions from several prior colleagues from different organizations. The full list of contributors is provided on the back page of the manual.

## 2 Overarching Goals and Rationale

### 2.1 Some Inviolable Commandments

- Never overwrite source data.
- Never store files that contain personally identifiable data on your local machine. In fact, try not to save *any* data on your local machine.
- Try your best to be wrong internally before being wrong externally.
- When output or results do not look correct (e.g., odd counts, ranges, means), always investigate with a skeptical mindset. Recheck your code, examine the raw data, and

---

<sup>1</sup>This document is narrowly focused on principles and rules for developing and maintaining transparent and replicable quantitative workflows in a collaborative environment. For a similar document that codifies analytic choices, see Donald Green's [Standard Operating Procedures](#) handbook for experimental designs.

establish additional examination points within your code. Hold the mindset that everything you have done is wrong, until you have confirmed that this is not the case.<sup>2</sup>

- Be intensely curious and thorough about examining the raw data.
- Be **assertive**!
- Never ignore the results from a merge. Mismatches can be the result of missing observations, ID issues, or other ways in which the data are not structured as you expected<sup>3</sup>
- Always log the results of each session that manipulates a dataset or generates a calculation used in a document that's shared externally. The `preamble .ado` written will help institutionalize this rule in your workflow.
- Never haphazardly drop an observation.
- All calculations in a report **must** have a code footprint.
- Don't copy and paste results from a statistical program to a document that is shared externally.
- Be sure to set the seed (and, in Stata, your `-sortseed-`) at the top of each script.<sup>4</sup>
- In the script preamble, be sure to indicate the creation date, author, and purpose of the syntax file, and maintain an update log with a concise note summarizing the reason and date for a substantive change in the source code.
- Where possible, rely on legacy code to perform procedures that have already been done.
- If you're not sure, **just ask**.

## 2.2 Replicability

- Workflows and results generated by one analyst for one project should be easily and accurately replicated by another analyst. This principle requires that the principal analyst establishes a clear, documented workflow, and each procedure that leads to a transformation of variables, or destruction of data, is recorded in a syntax file.

<sup>2</sup>To facilitate this type of scrutiny, a good habit is to generate tables of summary statistics (e.g., minimum/-maximum values, missing patterns) on key variables, disaggregated by natural grouping structures in your data (e.g., school year, state). This can help identify programming or coding errors before you begin analysis. Stata code for these procedures is available in the code repository.

<sup>3</sup>The R package `tidylog` provides functionality similar to Stata's for summarizing merge results in the `tidyverse`.

<sup>4</sup>Seeds and particularly sortseeds should be set immediately prior to the command requiring a stable random number or sort is used. Since pseudo random numbers come one by one from a list created by an algorithm using the seed, if commands which utilize a random number (for breaking ties in sorts or otherwise) any commands added between setting the seed and the desired command could change the result.

- Although we aim to review each other's code for every project, we typically only adhere to this principle consistently for projects with a complicated data management or analysis component, for projects with a high level of external visibility, and for new researchers. But, it is our goal to continue to make this more systematic, particularly as we add new staff.
- Each analyst should be able to replicate the results another analyst calculates in any software (with some decimal-place types of exceptions, particularly for more sophisticated procedures).
- For code review, the assigned QA researchers should denote sections requiring the attention of the primary analyst with a unique symbol that will allow the analyst to quickly control+F to find sections that require review. Below, is an example from a Stata .do file that was QA'd by Danial Hoepfner.

```
/*
Created by: MWG
Creation date: 4/15/2019
UPDATE LOG:
*!DH Note 5/30/2019: QC by DH look for *!DH for comments
*!DH Note 5/31/2019: More QC
*/
*!DH Note 5/30/2019: I don't have this scheme
if "'c(username)'" != "dhoepfner" preamble, ///
  clientf(rels_2017/sped_transitions) ///
  log mkdir sub(output syntax results documentation ///
  final raw csv converted zip clean) ///
  scheme(plotplainblind_rel) ///
  logpath(output) logname(${rq})
```

- To the extent possible, the methodological lead or primary analyst should enumerate QA steps that should be performed at different waypoints throughout a project workflow.

## 2.3 Transferrability

- Code should be transferable across projects and across individuals. That is, if one analyst left tomorrow, another one should be able to pick up her code from a given project and use it.
- Be sure to annotate any program or package dependencies in the preamble of your code. In R, this means loading all packages used in the preamble.

- Annotate your code richly, throughout each section, so that a future analyst (perhaps you!) can easily understand what was done (e.g., the workflow, processes), and why. Remember that what is apparent to you while you are in the middle of a project may not be apparent to a future analyst (or your future self!).
- If one analyst develops code for preparing files for a project, another analyst should be able to use this code for similar projects and problems, or even borrow components of it for similar procedures for different projects. Clearly structuring and annotating your code will facilitate this type of sharing.
- For multi-year projects, include the date (with year) with the comment so future analysts can see which comments are most recent.
- In R (or any statistical program), never `setwd` to a directory on your local machine.<sup>5</sup>.

## 2.4 Accessibility

- Code should be annotated and grouped by data management and analysis procedure.
- Ideally, some sort of structure should be created for your syntax files so that it flows logically, intuitively, and clearly.
- Sections of code should be clearly demarcated, such as “begin cleaning file X”, “begin cleaning file Y”, “merging Y and X”.
- Use informative names for local and global macros and variables, even if temporary (unless created and deleted within a few lines of code).
- If any sort of cleaning occurs outside of syntax, for example using an excel sheet to organize variable names across files, annotate that process in the script extensively.

# 3 Expectations and Guiding Principles

- Code should have a structure or outline that (roughly) follows the workflow of analytic project (see Section 4.2.1 below for more on this):
  - Data import/conversion from .txt/.csv/.sasbdat/etc. to format of the chosen software
  - Data cleaning and managing
  - Labeling and variable recoding and construction
  - Analysis

<sup>5</sup>In RStudio, it is better to maximize the Project functionality to set working directories, or use the package [here](#). Jenny Bryan has an informative encomium about the [here](#) package...[here](#). The reasons she provides for **never** setting your working directory in a script is a universal principle that applies to all programming languages.



- Presentation and output
- Anything that destroys, manipulates, or changes a file in any way should be recorded in code
  - After it is finalized, the base analytic file should **never** be overwritten.
  - This includes duplicate removal, dropping extraneous cases, variable creation and/or recoding, etc.
  - That is, the file constructed during the data preparation and management phase should always be retained. The inevitable changes to this file that occur during the analysis phase should either not be saved (that is, the analysis portion of your .do file must be re-run each time), or should be saved as a separate file.
  - This perhaps goes without saying, but the raw files should also not be overwritten. Raw files should not be edited, and if they must be, they should be saved with a new name noting the nature of any edits made to the original.<sup>6</sup>
- Analysis commands that produce results that ultimately go in a deliverable or published report should be annotated in the syntax file.
  - Ideally, this will be annotated, i.e., \*Table 2 crosstab.
- For analysts who use Stata, we strongly recommend you use the Gibson preamble command. The front-matter defines global path macros that will facilitate replication and quality assurance inspections.
- In general, syntax or .do files and raw source files should be stored on a network drive (e.g., the **S:/drive**) to ensure they are recoverable in case of some type of failure. If you need to have them accessible locally, please ensure they are backed-up to a network drive as soon as you are able.

## 4 Naming folders and files

When possible, (particularly for data and code stored on the **S:/shared/data/** drive on Azure) please name folders and files in a systematic way that is free of non-alphanumeric characters (save for underscore delimiters). Please do not include spaces, parentheses, or other symbols in folder or filenames as these can cause programming scripts to choke. Typical convention is to use snake\_case rather than CamelCase, though this has varied by researcher (and, consequently, project). Once a project folder name has been set, it should almost never be changed, since doing so will break all code referring to it.

---

<sup>6</sup>This should only be done in rare cases where, say, a file includes a header row above variable names, or includes characters which prevent the accurate import of data. Even still, scripted management of these issues is almost always preferred. **Data values in raw files should never be edited.**

## 4.1 Folder Naming and Structure

The project folder should include some distinctive moniker including the project client or title (like AISD.student\_survey) without project or school year included in the title. Project directories should be clearly and intuitively structured and labeled. Typically, syntax files will be in a folder separate from data files, and raw source files will be in a separate sub-folder from final analytic files.<sup>7</sup>

### 4.1.1 Common sub-folders

We generally use the same sub-folder or sub-directory structure in each project folder. We separate information into folders like **raw**, **converted**, **processed**, and **final** data as well as **syntax**, **output**, **results** (e.g. from models), **temp** (for tempfiles), **archive** (for archived or deprecated code), and other associated files like **latex** (for  $\text{\LaTeX}$  files) or HTML. If further organization of files is useful, nest those distinctions within these base folders. For example, if a large project has district administrative data, observation data, afterschool data, and survey data create those folders within raw, converted, processed, and so forth.

- **raw**: Unaltered files as received from the client or other source.
- **converted**: Raw files converted to format used by analyst's language.
- **processed**: Converted files which are partially processed (variable and value labels for example), but need more work before they are final/analytic (need to be merged or combined across years).
- **final**: Analytic dataset which is used for descriptive tables, analysis, figures and other results.
- **output**: Tables and other files which will not be used in reporting, such as tables describing data issues for the client. Files output at one stage to be imported at another(which are not data files). Examples: a list of students missing test data, an excel sheet you export to organize variable renames and re-import.
- **output/logs/**: Store log files from Stata/R sessions. Logs should always be saved with the date in the name so that the point where errors may have been introduced can be detected.
- **results**: Tables, figures, and other output that will or could be used in reporting.
- **syntax**: Script files.
- **latex/html/temp**: Use as needed to store  $\text{\LaTeX}$  and HTML files or other atypical files needed or produced.

<sup>7</sup>Keep in mind that there limits to folder path lengths (260 characters on our MS Windows server). This is affected by both file name length and the number of sub-folders the file is nested in; so, please keep file nesting to a minimum where possible to avoid hitting this limit.

The top of your script file should contain macros for each folder to make it easy to access, or use the preamble command which does so.

```
1 foreach fol in raw converted processed output final {
2     global 'fol' = ${sf}${fyear}'fol'
3 }
```

### 4.1.2 Working across Project Years

In most cases, projects take place across multiple years.<sup>8</sup> The folder structure silos the data by year in most cases. Commonly, we have the sub-folders shown in the section above are nested within project year folders (named “year\_1”, “year\_2” ... “year\_N”) This notation also makes it easy to access to create macros which can refer to the prior year by parsing the original year\_N macro. This allows for easy inclusion of data elements you want to add to the current year’s data.<sup>9</sup>

```
global fyear "year_7"
global lyfyear "year_'"'substr("${fyear}",-1,1)'"'-1'"

foreach y in "ly" "" {
    global 'y'sf "${stem}project_folder/${'y'fyear}/"
    global 'y'syntax "${'y'sf'}//syntax/"
    global 'y'converted "${'y'sf'}/converted/"
    global 'y'processed "${'y'sf'}/processed/"
    global 'y'final "${'y'sf'}/final/"
    global 'y'results "${'y'sf'}/results/"
    global 'y'output "${'y'sf'}/output/"
    global 'y'raw "${'y'sf'}/raw/"
}
```

The code example above will product raw and lyraw globals which refer to the folder in the relevant years.

## 4.2 File Naming

In general, filenames should:

- NOT contain non-alphanumeric characters (only contain [A-Za-z0-9] and underscores)
- contain underscores rather than spaces
- avoid starting folder and filenames with numbers (e.g., 2017\_projectname\_school)

<sup>8</sup>Even if a project is expected to only last one year, create a year\_1 as projects extending is not uncommon.

<sup>9</sup>preamble can also do this for you.

- use consistent abbreviations
- use camel case (e.g., **FileNameOne**) to compress file name length where underscores are not needed or all lowercase. Do not use arbitrary case (e.g., staffSTUDENT\_YR2012, myFileName\_for2013\_vER2.txt)
- be constructed to encourage looping over file names (that is, consistently named with sequenced elements). Example: surveyresponses\_school1\_2014\_grade8.dta. Underscores can help with looping over sub-elements of a filename!
- If you must use versioning, be consistent in version labeling and meaning. File versions 1 and 2 could be suffixed with “V1” and “V2” but if you make a dataset for Stata version 12 or 14 make this distinct (e.g., school1\_2014\_survey\_stata\_ver14.dta).
- For analysis files, particularly when working collaboratively, it can be helpful to include a **BornDate** variable with the date the file was created. In case of inconsistencies between versions, this can help diagnose the reasons for the discrepancies since different users can refer to this variable and determine whether they are using a deprecated file.

```
//At the end of a .do file, create a new variable containing the last
    saved date/time of the file.
//The last saved date+time is stored in the system macro c(filedate).
gen BornDate=clock("`c(filedate)'", "DMYhm")
format BornDate %tc
```

#### 4.2.1 *Syntax file naming*

These rules apply for naming all file types in our Data folder hierarchy; however, for syntax (e.g., do-files) files we prefer that the file names are prefixed with a numeric 'order' so they can be run in the proper order across the folder.

For example:

```
10_clean_raw.do
11_convert_raw.do
12_process_clean.do
21_analysis_descriptivesRQ1.do
22_analysis_descriptivesRQ2.do
23_analysis_modelsRQ1.do
30_analysis_appendix.do
```

For surplus syntax files that don't neatly fit into the analysis process/sequence, use a distinct (or large) number to set it apart (like 99). Examples include files containing meta code (like labeling), ado-files used to process procedures across multiple syntax files in this sub-directory, or scratchpad files of code you don't want to erase/lose.

Continuing the previous example:

```
10_clean_raw.do
11_convert_raw.do
12_process_clean.do
21_analysis_descriptivesRQ1.do
22_analysis_descriptivesRQ2.do
23_analysis_modelsRQ1.do
30_analysis_appendix.do
99_varlabels.do
99_value_labels.do
99_scratchpad_assigngroups.do
99_adofile_maketables.do
```

## 5 Naming variables, macros, and other characteristics in Stata datasets

Variables (columns), macros, and characteristics names should follow the constraints and conditions from Stata and R. Even in Excel and other program (where these rules don't apply), if we can enter variable names with the same rules that apply in Stata and R, then transferring this data cleanly will help reduce data cleaning effort.

Rules which meet Stata and R standards: variables (and macros) can contain up to 32 characters. A variable name may contain only the digits 0 to 9, upper or lower case English alpha characters (A to Z), and underscores; and the first character cannot be a number or an underscore.<sup>10</sup> Except where necessary, such as folder paths, years, and macros that need to persist when other script files are called from a main script, locals should be used rather than globals.

---

<sup>10</sup>Stata variables can be started with an underscore without issue, locals in Stata can, but should not start with an underscore, globals cannot as Stata uses the underscore internally to differentiate between locals and globals.

We advise creating loopable, consistent variable names that also convey meaning when possible (though you have -variable labels- and characteristics (-char-) to store that information. If the variable attributes are too long to easily store in a variable name then use short abbreviations (that hopefully capture some attribute with loopable structure) and rely on the labels and codebook to decipher meaning.

In terms of having loopable, consistent naming, follow similar patterns as for file naming where the entity is followed by the time or level component and any adjunct details in the suffix. Since space is at a premium, abbreviate any entities as necessary.<sup>11</sup>

Here are some good examples:

**student**\_grade8\_2004

**student**\_grade8\_2008

**student**\_grade8\_2012

**staff**\_grade8\_2004

**staff**\_grade8\_2008

**staff**\_grade8\_2012

**student**\_grade12\_2004

**student**\_grade12\_2006

**student**\_grade12\_2008

**student**\_grade12\_2012

**staff**\_grade12\_2004

**staff**\_grade12\_2008

**staff**\_grade12\_2012

This type of variable naming makes it easy to loop over type (**student**, **staff**),<sup>12</sup> grade level (8, 10, 12), and school year (2004, 2006, 2008, and 2012). Since these variable names are shorter than 32 characters, it leaves plenty of room to add suffixes to these names when generating new versions of these variables (e.g., **staff**\_grade12\_2012\_**meanscore**, **staff**\_grade12\_2012\_**v2**).<sup>13</sup>

---

<sup>11</sup>For example: “grade” can be “gr.”, student can be “stu.”, but remember to label your variables!

<sup>12</sup>We typically indicate the level of analysis (e.g., **stu\_** for **students** or **tch\_** for **teachers**) in the variable name prefix

<sup>13</sup>Beyond looping convenience, (like the folder naming convention above) these variable names are extensible. As new data is added in future years of the project, we can add a **staff\_grade6\_2018** survey with no problem to the existing programming if our code takes advantage of all levels of values contained in the data (e.g., the stata code uses something like `levelsof year, loc(yy)` rather than explicitly listing each year in the dataset).

Macro names should be descriptive and facilitate easy recall of what they are within loops. Descriptive macro names also reduce the chance an analyst will accidentally use the same macro name twice, which can be a difficult bug to figure out.<sup>14</sup> While Stata does not require indentation in loops or if blocks, each loop and if block should be indented to facilitate reading. For long loops or if blocks, the start and particularly the end, should be noted. For example, this may work:

```
local list race course gender grade
levelsof teacher, local(aa)
foreach a in 'aa' {
levelsof student, local(bb)
foreach b in 'bb' {
foreach c in 'list' {
levelsof 'c' if teacher == 'a' & student == 'b'
foreach d in 'r(levels)' {
count if teacher == 'a' student == 'b' & 'c' == 'd'
}
}
}
}
```

But this will work *and* be easier to read, especially if the manipulations or calculations or are more extensive, or if there are if blocks within it. Particularly with complex nested loops, labeling or numbering each loop eases de-bugging and future modifications.

```
local varlist race course gender grade
levelsof teacher, local(teachers)
foreach t in 'teachers' { //Teacher
levelsof student, local(students)
foreach s in 'students' { //Students
foreach v in 'varlist' { //Variable
levelsof 'v' if teacher == 't' & student == 's', local(varlevs)
foreach vl in 'varlevs' { //Levels
count if teacher == 't' student == 's' & 'v' == 'vl'
} //End level of variable
} //End of variable
} //End of student
} //End of teacher
```

<sup>14</sup>Or worse, not be detected and produce the wrong result!

Procedures and lists should be unified. This is to say, if a similar non-standard procedure is going to be executed multiple times, it should be done in a loop, or using a program/function written at the top of the syntax file.<sup>15</sup> Lists of variables which will be used repeatedly should be defined once.<sup>16</sup> Parameters that may be adjusted within a loop that are used several times should be defined as locals, for example, colors, text sizes, and spacing adjustments in a loop producing a figure composed of multiple `-twoway-` commands.

## 6 USC-CARE's computing environment and tools

### 6.1 Project-related storage

Project data are stored on OneDrive servers and are accessible through a browser or through a locally mapped drive, which allows access to files through your local machine. To sync a SharePoint folder to your local machine, follow the steps below:

1. Open SharePoint and go to the main Projects folder under Research and Evaluation.
2. Click the Sync button.
3. This will open up a dialogue about setting up OneDrive, sign in and work through this process.
  - create that folder if it does not exist
4. This should create a linkage whose path will depend on the user's operating system. For instance, in MacOS: `/Users/(USERNAME)/Library/CloudStorage/OneDrive-SharedLibraries-UniversityofSouthernCalifornia/CESR-Education - (folder)`.

The main SharePoint folder in which project- and non-project-related data are stored is the Data folder. Within the Data folder, projects (including the UAS) should have a common, predictable structure. Below is a directory tree for the One8 survey project which illustrates several of the principles advocated for in this document. I discuss each in turn.

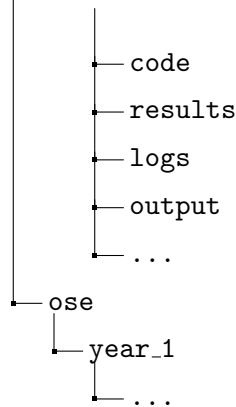
```
Data
├── One8
│   ├── pltw
│   │   ├── year_1
│   │   │   ├── clean
│   │   │   └── raw
```

---

<sup>15</sup>Standard repetitive procedures should be turned into an `.ado` file (in Stata, or package library in R)

<sup>16</sup>Stata 16 introduced a new suite of commands for managing variable lists called `v1`. A good summary of this functionality is available [here](#). These commands are useful for compiling lists of variables that are related in some way (e.g., demographic characteristics or test scores), and sorting them into globals based on their level of measurement (e.g., binary, continuous). These globals can then be used in subsequent operations, and can be manipulated using Stata's factor variables functions.





This project provides a useful illustration of how to design a project folder with an eye towards extensibility, but also project growth. Within the `Data/One8/` folder, there are two upper-level directories for each sub-project commissioned by the funder: `pltw` and `ose`. This demonstrates the idea that, for a given *client*, we may have multiple *projects*, and nesting these within a single client folder makes it easier to find projects associated with the same funder. This is particularly important since projects funded by the same client tend to have similarities which are conducive to sharing resources (e.g., data, code, reporting templates), so navigating to (and pointing file path constants to) is simplified with this organization principle in mind. This same principle informs the organization of the UAS folder, where waves are nested within the UAS directory.

## 7 Researcher written Stata programs

This section lists and provides a brief description of the programs contributed by current and former colleagues at Gibson Consulting Group, as well as current researchers at USC-CARE. The programs were developed to simplify reoccurring actions that we encounter in our projects. These are monolingual, and are all developed for Stata: we encourage polyglots to port these to their preferred programming language.

### 7.1 `preamble.ado`

#### 7.1.1 *Author: Marshall Garland*

#### 7.1.2 *Description*

`preamble.ado` incorporates the typical information that's included in a preamble into a single command. It has several options that simplifies this routine step.

#### 7.1.3 *Usage and options*

See the GitHub ReadMe for installation instructions.

## 7.2 `uas_preamble.ado`

7.2.1 *Author: Marshall Garland*

7.2.2 *Description*

`preamble.ado` is a convenience program for simplifying and standardizing the initial loading, cleaning, and variable derivation for each UAS education survey wave data file. See the package GitHub page for installation and usage instructions.

## 7.3 `varcount.ado`

7.3.1 *Author: Marshall Garland*

7.3.2 *Description*

`varcount.ado` is a command that counts the number of variables in a variable list or, if a variable list isn't provided, the number of variables in the dataset.

7.3.3 *Usage and options*

See the GitHub ReadMe for installation and usage instructions.

## 7.4 `graphsout.ado`

7.4.1 *Author: Marshall Garland*

7.4.2 *Description*

`graphsout.ado` is a wrapper for `graph export` that allows the user to export a graph from Stata to multiple file formats.

7.4.3 *Usage and options*

1. The command requires a file name.
2. `replace`: Overwrite the file name on disk, if it exists.
3. `type`: File format from the list of allowable file formats for `graph export`. Type `help graph export` to see the list of available formats.
4. `font`: Graph text font.

```
graphsout /Volumes/Desktop/, ///  
replace type(emf pdf svg) font("Calibri")
```

## 7.5 qualtrics.ado

7.5.1 *Author: Danial Hoepfner*

7.5.2 *Description*

qualtrics.ado interacts with the qualtrics API to list, download, and clean survey data collected on Qualtrics.

7.5.3 *Usage and options*

The command has three sub-commands, `qualtrics set`, `qualtrics list`, and `qualtrics get`. Unlike most of our .ados it has a help file.

1. `qualtrics set`: Set a password and optional user name to use instead of the token and data center each call.
  - token: Qualtrics token for our organization. **This is sensitive! Someone could delete all of our survey data with it**, required.
  - center: Qualtrics data center for our organization, “col” is our data center, required.
  - password: a password you specify, required.
  - user: a user name you specify, only needed if using multiple Qualtrics accounts or sharing an .ado folder with another.
2. `qualtrics list`: List surveys associated with your account, filter results.
  - match: Only list surveys whose name matches the included regular expression.
  - active: Only list surveys that are active.
  - modrange(MDY:MDY): Only list surveys that were modified in a date range.
  - createrange(MDY:MDY): Only list surveys that were created in a date range.
  - user: The username, optional.
  - password: The password, required unless specifying token and data center.
  - token: Qualtrics token for our organization, not needed if using password/user.
  - center: Qualtrics data center for our organization, not needed if using password/user.
3. `qualtrics get`: Download, convert, and clean a survey.
  - id: Only list surveys whose name matches the included regular expression.
  - csv: Folder to save the .csv file (name is maintained from Qualtricsm required).
  - dta: Folder and file name to save the .dta file, required if using clean option.
  - valuelabs: Request value labels rather than numeric codes, incompatible with clean option.

- clean: Clean dataset, apply labels, value labels and characteristics from survey meta data.
- relab: Display code to variable adjust labels.
- revallab: Display code to value adjust labels.
- preserve: Restore current dataset, rather than loading downloaded file.
- user: The username, optional.
- password: The password, required unless specifying token and data center.
- token: Qualtrics token for our organization, not needed if using password/user
- center: Qualtrics data center for our organization, not needed if using password/user

```

qualtrics set, t("QualtricsAPITokenHere") c("az1") ///
    password(ApassWOrd) u(danial)

qualtrics list, password(ApassWOrd) m(Fall 20[1-2][890])

qualtrics get, id("SV_XYXYXYXYX") password(ApassWOrd) ///
    csv("C:/Users/dhoepfner/Desktop/") ///
    dta("C:\Users\dhoepfner\Desktop\testdata.dta") clean relab reval

```

## 7.6 gtab.ado

7.6.1 *Author: Danial Hoepfner*

7.6.2 *Description*

gtab.ado is a tool to build custom tables easily. It works by creating a set of string vectors which you then fill with subsequent gtab commands. You can preview, import into Stata or export to excel

7.6.3 *Usage and options*

The command has six sub-commands, gtab init, gtab [column\_name], and gtab fill, gtab preview, gtab import, and gtab export. You can also build multiple tables simultaneously by inserting a number between gtab and the subcommand.

1. gtab [number] initialize: Initialize a table (or tables using a number) with the desired columns names and order. gtab will then print commands to add to each column.

2. `gtab [number] [column_name]`: After initializing the table, add row(s) to that column. Information can be added to columns in any order, though cells meant to be empty for a particular row must be filled (either by having a blank `gtab [column_name]` command or using `-gtab fill-`).
3. `gtab [number] fill`: Since `gtab` works by using a series of string vectors, in empty cells aren't specified, the table can become mis-aligned. `gtab fill` equalizes the vector lengths so that all columns have the same number of rows, by adding blank rows to columns shorter than the longest column.
4. `gtab [number] preview`: Preview table in the results window.
5. `gtab [number] import`: Import table into current data frame.
6. `gtab [number] export 'file_path_and_name'`: Export table to an excel file.
  - sheet: Workbook sheet to export to.
  - replace: Replace excel file if it exists.
  - sheetreplace: Replace sheet if it exists.
  - sheetmodify: Modify sheet if it exists.

Basic usage with a number, allowing multiple tables simultaneously

```
sysuse auto, clear
gtab 1 init var mean min max

gtab 1 var Variable
gtab 1 mean Mean
gtab 1 min Min
gtab 1 max Max

ds *
foreach v in `r(varlist)' {
    if !regexp("`v'", "type 'v'", "str")) {
        gtab var `:variable label `v'`
        sum `v'
        gtab 1 mean `:di %3.1f `r(mean)''
        gtab 1 min `:di %3.0f `r(min)''
        gtab 1 max `:di %3.0f `r(max)''
    }
}
gtab 1 mean Legend: Data from auto.dta
```

```

gtab 1 fill
gtab 1 pre
gtab 1 export "C:/users/dhoepfner/results/example.xlsx", ///
replace sheet("descriptives")

```

### Advanced Usage

```

sysuse auto, clear
/*Here I am going to use the levels of foreign to help define the columns,
this will make tables extensible to additional levels added to a variable.
I'll define and add to a local called gtab that defines the columns.
When filling the table, you can use of the same loops for the table.*/

local gtab var
levelsof foreign
foreach l in `r(levels)' {
    local gtab `gtab' m`l'
}
*Now adding a column for the difference
local gtab `gtab' d
*Now adding columns for the Ns
levelsof foreign
foreach l in `r(levels)' {
    local gtab `gtab' m`l'
}
*Now initializing gtab
gtab init `gtab'
*Now filling the title row
gtab var Variable
levelsof foreign
foreach l in `r(levels)' {
    gtab m`l' `:label (foreign) `l'' Mean
    gtab n`l' `:label (foreign) `l'' N
}
gtab d Difference

/*Now we'll loop over each variable and level of foreign
and extract and export the desired statistics.

```

Notice that we can use stored results to define stars for significance, or whatever else we want.\*/

```

local vct=0
ds foreign, not
foreach v in `r(varlist)' {
    if !regexm("`':type `v''", "str") {
        local ++vct
        gtab var `:variable label `v''
        levelsof foreign
        foreach l in `r(levels)' {
            sum `v' if foreign == `l'
            local m`l' `r(mean)'
            local n`l' `r(N)'
            gtab m`l' `:di %7.1fc `m`l''
            gtab n`l' `:di %6.0fc `n`l''
        }
        local stars
        ranksum `v', by(foreign)
        if `r(p_exact)' < .05 local stars `"'*"'
        if `r(p_exact)' < .01 local stars `"'**"'
        if `r(p_exact)' < .001 local stars `"'***"'
        if `vct' !=5 gtab d `:di %3.2f `='m1'-`m0''''`stars'
        if `vct' ==5 gtab d `:di %3.2f `='m1'-`m0''''Custom Tables!
    }
}
*Import the table into the current dataset/frame.
gtab import

```

## 7.7 fmtexl.ado

7.7.1 *Author: Danial Hoepfner*

7.7.2 *Description*

fmtexl.ado is a wrapper for putexcel which formats excel tables nicely. By default it loops over each sheet in an excel file. Notes: You need to use -fmtexl using- syntax to specify excel file. Default settings are not optimal unless formatting tables like JS likes, the bw option is closer to journal style tables.

**7.7.3** *Usage and options*

1. headfill: Header fill color, default is Gibson blue
2. headtext: Header text color, default is white
3. lfill: Light row fill color, default is white
4. dfill: Dark row fill color, default is light blue
5. rowcol: Table body rows alternate colors
6. doublerowcol: Table body rows alternate colors in blocks of two, for example, when including standard error in row below estimate.
7. headrows: More than 1 header row, specify number
8. labcols: More than 1 label column, specify number
9. specrange: Specify table cell range rather than detect from import excel, allows formatting of multiple tables on the same sheet. Formatted as [A-Z]\*[1-9]\*:[A-Z]\*[1-9]\*
10. bw: Black and white journal style format. Overwrites all formatting options except for labcols and headrows

```
fmtexl using "C/users/dhoepfner/results/example.xlsx", ///
bw labc(2) headr(2)
```

**7.8** **apchx.ado****7.8.1** *Author: Danial Hoepfner***7.8.2** *Description*

apchx.ado loops over all variables in current dataset and looks for type differences in variables in a specified dataset. If you attempt to append a dataset in Stata and a variable is byte in one dataset, but string in another, you will get an error. Without this, Stata errors after the first mismatch and so this issue can repeatedly bite. It also has options to compare value ranges and labels in the two datasets.

**7.8.3** *Usage and options*

1. range: Check variable ranges are consistent across datasets.
2. label: Check value labels are consistent across datasets.

```
apchx using "C/users/dhoepfner/cleaned/example.dta", label range
```



## 7.9 conwithin.ado

7.9.1 *Author: Danial Hoepfner*

7.9.2 *Description*

conwithin.ado checks whether variable list A is constant within variable list B. Very useful for panel or nested data.

7.9.3 *Usage and options*

1. within: Variable list B, typically ID variables in which variable list A should be constant.
2. tag: Tags observations where a variable in list A is not constant within list B.

```
conwithin race gender disability, wi(student_number) t(t)
sort student_number
list student_number race if race_t == 1
```

## 7.10 edrop.ado

7.10.1 *Author: Eric Booth*

7.10.2 *Description*

edrop.ado: -cap drop- won't drop all variables in a list if one does not exist. -edrop- drops all variables in the list, and reports those that were not present in the first place.

7.10.3 *Usage and options*

1. keeP: Keep that list instead of dropping
2. temp: Also drop temporary variables

```
edrop race1 race2 race3
```

## 7.11 latest.ado

7.11.1 *Author: Danial Hoepfner*

**7.11.2** *Description*

latest.ado: Best practice is to save analytic datasets with a date appended, so the point errors were introduced can be tracked down. However, this means you need to specify the date of the most recent file for each -use- command. This returns 'r(file)' which is the most recently saved file, optionally matching a regular expression.

**7.11.3** *Usage and options*

1. `match`: Return most recent file matching regular expression.

```
latest `"{converted}"', m("student.+\.dta")
use `"{converted}'r(file)'"', clear
```

**7.12** **char\_finder.ado****7.12.1** *Author: Danial Hoepfner***7.12.2** *Description*

char\_finder.ado: -charlist- on SSC by Nick Cox works great to identify all of the unique characters in a string variable. However, sometimes order matters, and so char\_finder returns the text of a string, with the ASCII codes displayed below. Works well for finding Microsoft Word style quotes and other symbols (which Stata reads as a series of ASCII codes, which can then be replaced with subinstr() and '=char(x)').

**7.12.3** *Usage and options*

```
char_finder `";asdklfjsd ;lfjd(*)*)(@"'
```

**7.13** **fixtex.ado****7.13.1** *Author: Danial Hoepfner***7.13.2** *Description*

fixtex.ado converts symbols in text that will cause a  $\text{\LaTeX}$  script to choke into their proper codes.

**7.13.3** *Usage and options*

```
fixtex A sentence_with $ some probl&em symbols  
tex 'r(output)'
```

## 7.14 texpdf.ado

7.14.1 *Author: Danial Hoepfner*

7.14.2 *Description*

fixtex.ado compiles a .tex document.

7.14.3 *Usage and options*

1. copy: Copy resulting PDF to a new location.
2. open: Open PDF once compiled.

```
texpdf '$\{tex\}region_10_report.tex', copy('$\{to_distribute\}')
```

## 7.15 syntax\_saver.ado

7.15.1 *Author: Danial Hoepfner*

7.15.2 *Description*

syntax\_saver.ado attempts to recover syntax when Stata crashes. Copies temporary syntax files to a restore folder, or another specified folder. Typically, Stata only stores the most recent run in an instance, so if you only highlighted one line, that's likely all you'll get.

7.15.3 *Usage and options*

1. restore: Specify a folder to save restored files to. Default is C:/Users/'c(username)'/Desktop/restore/
2. Today: Only restore files from today.

```
syntax_saver, today
```

## 7.16 latout.ado

7.16.1 *Author: Danial Hoepfner*

### 7.16.2 *Description*

latout.ado Loops over a list of variables, or all variables in a dataset, and produces a PDF to check for issues. Can select a -by- variable to present figures/tables by another variable. Particularly useful for multi-year data submissions.

### 7.16.3 *Usage and options*

1. title: Title of PDF
2. path: Where to save PDF, default is working directory
3. footer: Text to place in footer which links back to TOC, default is "Return to Top"
4. byvars: Variables to split output by, max =2
5. holes: Where to have a hole in panelled figures, useful when byvars levels is an odd number
6. columns: Number of columns for by variable panelled figures
7. cutoff: Cutoff between number of categories before variables are treated as continuous, default is 5
8. strings: How to display strings with more than CUT categories, options are -list-, -sample-, and -skip-, default is sample
9. varnames: Include variable names (not just labels) in output.

```
latout , title("AVID Submission 2017") ///  
pa'("${output}")' bvy(institution year) varn
```

(2023) All rights reserved

### **Authors**

Marshall Garland,  
[garlandm@usc.edu](mailto:garlandm@usc.edu)

Eric Booth,

Jill Carle,

Danial Hoepfner,  
[dhoepfner@gibsonconsult.com](mailto:dhoepfner@gibsonconsult.com)

David Osman,  
[dosman@gibsonconsult.com](mailto:dosman@gibsonconsult.com)

Mitchell Kilborn,

Gracie Petty,  
[gpetty@gibsonconsult.com](mailto:gpetty@gibsonconsult.com)

**For more information, please visit:**

<https://cesr.usc.edu/care>