

PanelWhiz: Efficient Data Extraction using the German SOEP

John P. Haisken-DeNew (University of Bochum, Germany) and
Markus Hahn (Melbourne Institute, Australia)

Abstract: This paper outlines a panel data retrieval program written for Stata/SE or better, which allows easier accessing of the German SOEP household panel data set. Using a drop-down menu and mouse click system, the researcher selects variables from any and all available years of the panel. The data is automatically retrieved and merged to form a “long file”, which can be directly used by the Stata panel estimators. The system implements modular data cleaning programs called “plugins”. Yearly updates to the data retrievals can be made automatically. Projects can be stored in libraries allowing modular administration and appending.

Keywords: Panel data storage and retrieval
JEL: C81, C87, C23

1. Introduction

Applied social scientists have forever been faced with different data interfaces for different data sets. In most cases, an interface is not even available, forcing the researcher to address data files by name, and extract the information required by hand. However, the specific structure of a panel data can vary dramatically as described in Haisken-DeNew (2001), such that some data sets provide many files per year, differing by their population, or level of aggregation etc., creating many obstacles for researchers.

PanelWhiz is a collection of subroutines that allows researchers to use an intuitive “common” graphical interface for accessing many datasets directly within the statistical package Stata/SE, (<http://www.stata.com>) whereby the researcher does not select individual variables, but rather vectors of variables (items) with one mouse click. This allows for an efficient method of selecting information for a data set retrieval, especially if the data set contains many waves (years) of information. With one mouse-click, data can be automatically retrieved, with merging and matching done automatically.

With the PanelWhiz system, the user can open data files by clicking on a browse page. The idea behind the tool is that because of the intrinsically longitudinal nature of the data, one is typically not interested in retrieving a variable in a particular wave, but rather in retrieving the variable for several waves, i.e. an item-correspondence. For all data sets, a variable renaming algorithm (where necessary) is used to ensure time consistent variable names (See Haisken-DeNew, 2001 for more information on this). Thus, if one opens a data file and one finds a variable of interest, one clicks on the variable and information for the entire item correspondence is also collected and added to a "PanelWhiz project". Straightforwardly, the object is to collect items and save them into the data project, allowing an automatic data retrieval. The data are extracted in “long” format allowing easy further data cleaning or direct estimation using Stata’s panel “xt” commands.

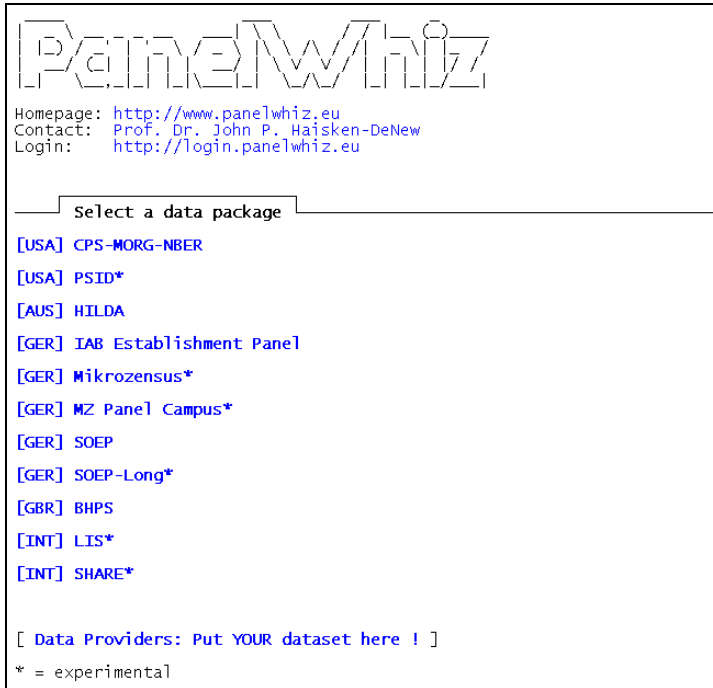
PanelWhiz, since appearing in 2005, now has several hundred registered users, using the common interface to access many different datasets other than the SOEP, such as the British BHPS, the German IAB Establishment Panel, the Australian HILDA, the American CPS, etc. Recently, support has been extended to the American PSID. This paper describes use of PanelWhiz for the German SOEP. However due to the generalized nature of PanelWhiz, the interface is almost completely identical for all other supported datasets.

2. Overview of PanelWhiz

2.1 Starting Up Stata with PanelWhiz

PanelWhiz is installed as a Stata Add-On and loads every time Stata is started. For example in following Screen Shot, one can select by mouse click the desired data set to be supported.

Screen Shot 1: Select Data Set



In this example, the SOEP has been selected. One can select an already existing project, or create a new one from scratch. In this example, we will examine an existing project **zufu.soep**. Because it has been already saved, PanelWhiz keep a note of the last 10 saved projects and allows easily loading by simply clicking on the blue link indicating the project name.

Screen Shot 2: Open a Project



Here we have indeed opened the PanelWhiz project **zufr.soep**, and have a heads-up display indicating the contents of the project and the possible project commands. The top area displays the possible project commands and the bottom area the contents of the project. Here the project already contains 4 items (vectors of variables). The item labels are blue (and therefore clickable, linked to a keyword thesaurus).

Screen Shot 3: Project Page

Start Page > Project Page <Project Page>

SOEP Project Browse Page
PanelWhiz SOEP (4.0) Sep 2010 <john@panelwhiz.eu>

Project: [zufr] from <10 Sep 2010> at <17:11:19>
<Untitled>
<John Häisken-DeNew> at <jhaiskendenew@rwi-essen.de>

Project Prefs Retrieved Data Library Help

Build Project: [Add More Items | ?]
Old Projects: [Update English or German | ?]
My Own Projects: [New | Open | Save, As | Save | Append | ?]
Recover Projects: [Acquire | Restore | ?]
Item Undo/Redo: [Undo | Redo | Erase History | ?]
Item Auto-Refresh: [On | OFF | Refresh NOW | ?]
Plugins <alt> : [Refresh | Remove ALL | Include ALL | ?]
Retrieval <alt> : [Do Retrieval NOW | ?]

Items [4] Items+Vars Items+Vars+Files

All Categories

[X P p w i]	IK622	EN: General Satisfaction with Life Now
[X _ _ w _]	IK2743	EN: Current Monthly HH Net Income
[X _ _ w i]	IK3111	EN: Federal State/Province
[X _ _ w i]	IK5600	EN: NUTS1 Region Code (Federal State)

The variables underlying the 4 items in the project are listed below. Each variable listed under the item displays the associated wave/year. In the SOEP, the “a” wave is 1984, the “b” wave is 1985 and so on. The “z” wave is 2009. Thus for the item *IK622*, in 1984 the underlying variable is **AP6801** and in 2009 it is **ZP15701**.

Screen Shot 4: Start Page

Items [4] Items+Vars Items+Vars+Files

All Categories

[X P p w i]	IK622	EN: General Satisfaction with Life Now	a-ap6801	b-bp9301	c-cp9601	d-dp9801	e-ep89	f-fp108
			g-gp109	G-gp6401e	h-hp10901	H--	i-ip10901	j-jp10901
			k-kp10401	l-lp10401	m-mp11001	n-np11701	o-op12301	p-pp13501
			q-qp14301	r-rp13501	s-sp13501	t-tp14201	u-up14501	v-vp154
			w-wp142	x-xp149	y-yp15501	z-zp15701		
[X _ _ w _]	IK2743	EN: Current Monthly HH Net Income	a-ah46	b-bh39	c-ch51	d-dh51	e-eh42	f-fh42
			g-gh42	G-gh36e	h-hh48	H--	i-ih49	j-jh49
			k-kh49	l-lh50	m-mh50	n-nh50	o-oh50	p-ph50
			q-qh54	r-rh49	s-sh4901	t-th4801	u-uh4801	v-vh5101
			w-wh5101	x-xh5101	y-yh5201	z-zh5201		
[X _ _ w i]	IK3111	EN: Federal State/Province	a-abu1a	b-bbu1a	c-cbu1a	d-dbu1a	e-ebu1a	f-fbu1a
			g-gbu1a	G--	h-hbu1a	H--	i-ibu1a	j-jbu1a
			k-kbu1a	l-lbu1a	m-mbu1a	n-nbu1a	o-obu1a	p-pbu1a
			q-qbu1a	r-rbu1a	s-sbu1a	t-tbu1a	u-ubu1a	v-vbu1a
			w-wbu1a	x-xbu1a	y-ybu1a	z-zbu1a		
[X _ _ w i]	IK5600	EN: NUTS1 Region Code (Federal State)	a-nuts184	b-nuts185	c-nuts186	d-nuts187	e-nuts188	f-nuts189
			g-nuts190	G--	h-nuts191	H--	i-nuts192	j-nuts193
			k-nuts194	l-nuts195	m-nuts196	n-nuts197	o-nuts198	p-nuts199
			q-nuts100	r-nuts101	s-nuts102	t-nuts103	u-nuts104	v-nuts199
			w-nuts106	x-nuts107	y-nuts108	z-nuts109		

One can update the project using the automatic “Update” function on the project page. When a new data distribution becomes available, for each item, the newest variable is added automatically to the relevant item. The retrieval can be run again, having now the most recent information.

2.2 Items and Specials

Assuming that one would like to add any items to the project, one can choose between two types of concepts: “items” or “specials”. Items are vectors of variables that have a standard time dimension associated with them, i.e. one variable for each year. SOEP examples of these files would be **AP**, **APGEN**, **AH**, **AHGEN** etc. Specials have a non-standard time dimension, i.e. they may have one observation per person and be time invariant, or may already be in long format, with person-year observations as the unit of analysis. Here we will first examine the page associated with items.

Clicking on a year like [**a - 1984**], will load a browse page allowing one to click on all variables/items associated with the year 1984. Alternatively, one can click on a special file, such as [**bioimmig**] where the data is already in long format. In contrast, the information from the special file [**bioparen**] is time invariant and contains only one entry per person. PanelWhiz knows how to extract and merge the information from all of these kinds of files.

Screen Shot 5: Items and Specials

The image displays two screenshots of the PanelWhiz web interface. Both screenshots show the 'My PanelWhiz SOEP Project: zufr' and navigation tabs for 'Projects', 'Waves', 'Specials', and 'Help'.

The left screenshot shows the 'Select a wave' section with a grid of year-based items:

[a - 1984]	[b - 1985]	[c - 1986]	[d - 1987]
[e - 1988]	[f - 1989]	[g - 1990]	[h - 1991]
[i - 1992]	[j - 1993]	[k - 1994]	[l - 1995]
[m - 1996]	[n - 1997]	[o - 1998]	[p - 1999]
[q - 2000]	[r - 2001]	[s - 2002]	[t - 2003]
[u - 2004]	[v - 2005]	[w - 2006]	[x - 2007]
[y - 2008]	[z - 2009]		

Below this is the 'Waves keyword index' with a grid of letters: [*****] [a] [b] [c] [d] [e] [f] [g] [h] [i] [j] [k] [l] [m] [n] [o] [p] [q] [r] [s] [t] [u] [v] [w] [x] [y] [z].

The right screenshot shows the 'Select a file' section with a list of special files:

- [bioage01] Bio Age under 1
- [bioage03] Bio Age under 3
- [bioage06] Bio Age under 6
- [bioage17] Bio Age under 17
- [biobirth] Bio Birth Mother
- [biobirthm] Bio Birth Father
- [bioimmig] Bio Immigration
- [biojob] Bio First Job
- [bioparen] Bio Parents
- [bioresid] Bio Residence
- [biosoc] Bio Societal Background
- [biotwin] Bio Twin/Multiples
- [cognit06] Cognitive Abilities
- [exit] Sample Exit
- [pwealth] Wealth Person
- [hwealth] Wealth Household
- [health] Health
- [gripstr] Grip Strength

Below this is the 'Specials keyword index' with a grid of letters: [*****] [a] [b] [c] [d] [e] [f] [g] [h] [i] [j] [k] [l] [m] [n] [o] [p] [q] [r] [s] [t] [u] [v] [w] [x] [y] [z].

For both items and specials, all associated items have been scanned and the contents of the item labels have been catalogued into a thesaurus of keywords. Thus, if one were interested in all items or specials regarding the topic of “occupation”, one would click on the [**o**] of the keyword index, to examine all keywords starting with the letter “o”.

Technically speaking, PanelWhiz works because the item correspondence information (for each item, that vector of variables over all 26 years) is injected into each relevant variables as a Stata variable characteristics. PanelWhiz reads this information from a variable in one particular file/wave and automatically knows where to find corresponding information in all other files/waves.

2.3 Item Browse Page

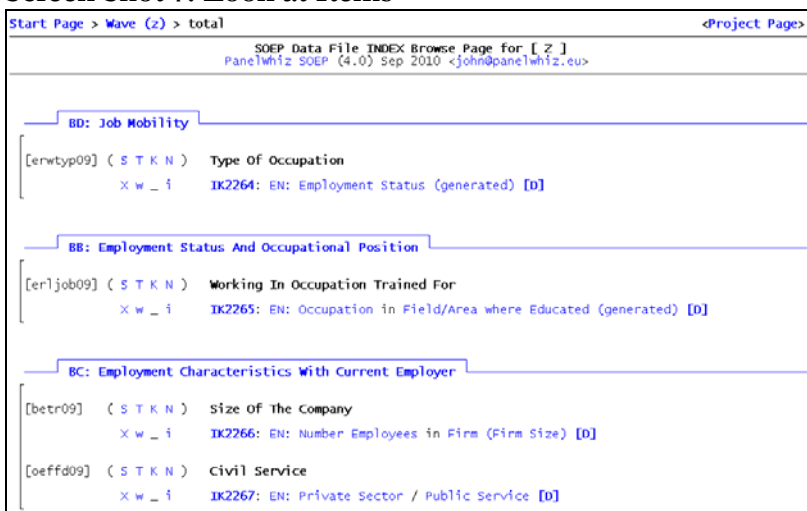
To find an item found say in the year 2009, we click on the year [**z** – **2009**], and receive the following browse page. It contains a browse page containing all variables/items from all files in the year 2009. Alternatively, one can jump to only specific variables/items in a particular file. Further, the variables are sorted using the same hierarchical scheme as in SOEPinfo. See Haisken-DeNew and Frick (2005) for more information on SOEPinfo.

Screen Shot 6: Top Level Item Browse Page



In this example, we select the blue clickable button [**ALL FILES**] from wave Z and get the following browse page. All variables are listed in the order they naturally exist in the respective physical data files. The example shows that for the SOEP variable **erwtyp09**, there is a PanelWhiz item **IK2264** associated with it. By actually clicking on **IK2264**, one would select the entire item (potentially all underlying variables from wave “a” through “z”).

Screen Shot 7: Look at Items



One can also examine the changing nature of the item over time. The variable **erwtyp09** contains value labels. Thus we can click on blue “i” to the left of the item name and label. There is a ready-made HTML page showing all labelled values for all variables of the entire item. This will be especially useful information for data cleaning requirements. Just because a variable has been coded one way in one year/wave, it does not mean it will remain so over all time. The

following Screen Shot illustrates this example. Jumping from wave 1984 to 1985, there have been some additional outcome values added. These changes are colour coded in grey.

Screen Shot 8: Start Page

(2264) EN: Employment Status (generated)

[~ := no data available] [x := no label/value available]

VALUE	1984 (a)	1985 (b)	1986 (c)	1987 (d)
1	x	[1] Not Employed, Green	[1] Not Employed, Green	[1] Not Employed, Green
2	[2] Not Employed (First Surveyed) Not Applicable Since 94	[2] Not Employed (First Surveyed) Not Applicable Since 94	[2] Not Employed (First Surveyed) Not Applicable Since 94	[2] Not Employed (First Surveyed) Not Applicable Since 94
3	[3] Employed (First Surveyed) Not Applicable Since 94	[3] Employed (First Surveyed) Not Applicable Since 94	[3] Employed (First Surveyed) Not Applicable Since 94	[3] Employed (First Surveyed) Not Applicable Since 94
4	x	[4] Empl. Exc Change	[4] Empl. Exc Change	[4] Empl. Exc Change
5	x	[5] Empl. No Info If Change	[5] Empl. No Info If Change	[5] Empl. No Info If Change
6	x	[6] Empl. With Change, Also First Time Employment	[6] Empl. With Change, Also First Time Employment	[6] Empl. With Change, Also First Time Employment
7	x	x	x	x

By clicking on “N” to the left of the variable label, one can view the item “notes”, giving an indication of the variable names of variables belonging to the item over all years.

Screen Shot 9: Item Notes

erwtyp09	
Itemname	2264
Itemlabel	EN: Employment Status (generated) DE: generierter Erwerbstatus (Erwerbstyp)
Category	BD: Job Market And Occupation: Job Mobility BD: Arbeitsmarkt und Beschaeftigung: Berufliche Mobilitaet
JEL	J21
Itemvector	erwtyp84 erwtyp85 erwtyp86 erwtyp87 erwtyp88 erwtyp89 erwtyp90 ----- erwtyp91 ----- erwtyp92 erwtyp93 erwtyp94 erwtyp95 erwtyp96 erwtyp97 erwtyp98 erwtyp99 erwtyp00 erwtyp01 erwtyp02 erwtyp03 erwtyp04 erwtyp05 erwtyp06 erwtyp07 erwtyp08 erwtyp09

All words appearing in an item label have been added to a keyword thesaurus. Each keyword is linked to all items in the entire dataset containing the keyword.

Screen Shot 10: Start Page

Start Page > Entries for 0 <Project Page>

SOEP Keyword Browse Page for [0]
PanelWhiz SOEP (4.0) Sep 2010 <john@panelwhiz.eu>

[*****] [a] [b] [c] [d] [e] [f] [g] [h] [i] [j] [k] [l]
[m] [n] [o] [p] [q] [r] [s] [t] [u] [v] [w] [x] [y] [z]

occ	occupationally	occasional	occupied	occupation	occupiers	occupational
october	office	offs	off	offer	offer	offered
okt	one	one's	openly	oneself	operation	oil
onw	optimism	oral	openly	operation	opinion	once
organisation	original	organization	order	order	orderd	only
orphans	out	orthopaedist	organizational	organizations	origin	opinion
out	over	outgoing	orp	orpha	orphan	orderd
over	own	overall	other	other	others	origin
own		owner	outside	outside	outwork	orphan
			overnight	overnight	overtime	others
			ownership	ownership		outwork
						overtime

[occ] TOP
 X N J - IK3796 EN: K.A. Item Nonresponse (Occ Status)

[occasional] TOP
 X N J - IK299 EN: Occasional Work for Money

[occupation] TOP
 X N J - IK2265 EN: Occupation in Field/Area where Educated (generated)
 X N J - IK2279 EN: ISCO88 Occupation Code
 X N J - IK229 EN: Occupation in Field/Area where Educated
 X N J - IK255 EN: Starting Fresh in Different Occupation
 X N J - IK3478 EN: Occupation of Individual
 X N J - IK3846 EN: Nat.Stat.Office-Occupation (Infratest)
 X N J - IK428 EN: Occupation Specific Insurance: Retirement Pension
 X N J - IK438 EN: Occupation Specific Insurance: EU/BU widow(er) Orphan
 X N J - IK439 EN: Occupation Specific Insurance: widow(er) Orphan EU/BU

2.4 Plugins

The variables underlying an item may vary of course from year to year. Using the PanelWhiz plugin system, the user can automatically create small scripts to clean time inconsistent data.

Screen Shot 11: Plugin Wizard

Start Page > Wave (z) > total_zpgen > Plugin Wizard <Project Page>

Panelwhiz SOEP - Plugin Wizard
Panelwhiz SOEP (4.0) Sep 2010 <john@panelwhiz.eu>

de_IK2264.ado Based On: NEW Plugin

[Create] [View] [Erase] [Include Now] [New] [Open (Original)] [Open (Personal)]

Plugin Information [Default] [Edit]

Plugin : [IK2264] EN: Employment Status (generated)
 Author : John Haisken-DeNew
 Email : jhaisken@panelwhiz.eu
 Web : http://www.panelwhiz.eu
 Copyright : Copyright 2006. All rights reserved. Unauthorized copying prohibited.
 Version : 1.0
 Years : a b c d e f g h i j k l m n o p q r s t u v w x y z
 Possible : a b c d e f g h i j k l m n o p q r s t u v w x y z
 Changes : - - - - -

Plugin Initialization [Default] [Edit]

Update year : 2009
 Depends on : <nothing>

Plugin Action List [Clear List]

[+|-] [up|down]

Depending on the dataset, various actions can be selected. For example, using plugins, nominal money values can automatically be deflated using an integrated CPI function.

Screen Shot 12: Plugin Wizard and Defining Functions

Plugin Wizard: Plugin Action Manager

Select Plugin Action: dm2euro

- dm2euro
- east
- label
- newitem
- real
- realshare
- recode
- replace
- labellang
- comlong
- recodelong
- replacelong
- renamelong
- verbalmlong

Converting DM to EURO

This action is for converting an amount of money from Deutsch Mark (DM) to EURO (€). The conversion is done automatically for the appropriate waves.

OK Cancel Submit

2.5 Retrievals

Once a project has been created, PanelWhiz has enough information to retrieve the actual data from the SOEP dataset. PanelWhiz dynamically creates a command file and executes it on-the-fly. PanelWhiz opens the files that the user specifically addressed and pulls out the variables specifically selected. It stores these variables in a temporary file and then moves on to the next data file. It does this many times until all data chunks have been extracted, then merges all data chunks together in the manner prescribed by the user.

The extracted data are automatically in Stata long format, ready to be processed using Stata's rich panel "xt" commands. To document exactly the retrieval that was run, PanelWhiz creates an executable DO file. The following Screen Shot shows an excerpt of a generated DO file.

Screen Shot 13: Excerpt of a generated DO file

```

/* -----( create master )----- */;
save      "$tmp/master", replace;
/* -----( pull: hp / 1991 Person )----- */;
use       hhnr      persnr
          hp10901
using     "$soep/hp";
label    lang DE;
pwtclone hp10901  IK622;
drop     hp10901;
sort     persnr;
save     "$tmp/hp", replace;

/* -----( pull: ip / 1992 Person )----- */;
use       hhnr      persnr
          ip10901
using     "$soep/ip";
label    lang DE;
pwtclone ip10901  IK622;
drop     ip10901;
sort     persnr;
save     "$tmp/ip", replace;

```

3. Summary

PanelWhiz is a data retrieval tool that eases data extractions from the many large scale data sets such as the German SOEP. PanelWhiz is directly combined into Stata SE, allowing a seamless interaction between the micro data and the statistics package. Vectors of variables, called item-correspondences can be selected all at once. Special cleaning programs written in Stata called "plugins" can clean a particular item-correspondence and make it time and/or content consistent. Groups of item-correspondences can be stored as projects. Groups of projects can be stored as libraries. This method of organizing the projects and plugins allows for a modular administration, facilitating knowledge transfer and group work. Data can be retrieved by mouse-click, providing rectangularized data in long format. As new releases of the micro data become available, the user can "automatically" update his projects to include the latest wave of information. All programs used are available in source Stata code which allows complete transparency of content. All commands used in the generated retrieval are documented in a full functional retrieval DO file, capable of recreating the identical retrieval at any time.

References

Haisken-DeNew, John P. (2001) "A Hitchhiker's Guide to the World's Household Panel Data Sets." *The Australian Economic Review*. 34(3), 356-366.

Haisken-DeNew, John P. and Joachim R. Frick (2005) "The Desktop Companion to the German Socio-Economic Panel Study", DIW Berlin, Germany