

# Converting html tables into a Stata dataset

Jan Ditzen

Free University of Bozen-Bolzano, Bozen, Italy

[www.jan.ditzen.net](http://www.jan.ditzen.net), [jan.ditzen@unibz.it](mailto:jan.ditzen@unibz.it)

<https://janditzen.github.io/htmltab2stata/>

May 9, 2024

# Motivation

- Data often available online in form of html Tables.
- Copying content of tables to Stata slow and manual.
- `htmltab2stata` automatises downloading html Tables directly into a Stata dataset.
- Restricted to HTML “<table>” environment.

# htmltab2stata

## Syntax

```
htmltab2stata , url(string) [tablenumber(string) firstrow href]
```

- `url(url)` the url of the html website to be processed. The url has to be a downloadable html website. url can be a web address or a local html file.
- `tablenumber(integer)` number of table within the html document. Default is 1, i.e. the first table is processed.
- `firstrow` Use firstrow of table as variable names.
- `href` Links enclosed in `<a href=....></a>` are added to the content transferred to Stata.

# Example

## Tables on My Github Page



A hillwalking economist working on panel time series with a love for coding econometric methods.

Edit profile

AK 92 followers · 23 following

Free University of Bozano-Bozen  
www.jan.ditzen.net  
@janditzen  
https://bsky.app/profile/janditzen.bsky.social  
in/jan-ditzen-54594477

Achievements



Send feedback

janitzen / README.md

### Welcome

I am an Assistant Professor (RFO-A) at the Free University of Bozen-Bolzano, Bolzano, Italy.

My research interests are in the field of Applied Econometrics with a focus on Panel-Time Series and Spatial Econometrics and Growth Empirics. I have also an interest in Growth Theory, Simulation Studies and the implementation of econometric methods into statistical software.

See also my [Personal webpage](#).

### Stata packages - Econometric Methods

Package	Version	Updated	Description	Article/Sides
<a href="#">xtfce2</a>	Release: <a href="#">v4.3</a>	<a href="#">Release date: 31 March 2024</a>	Estimation of Common Correlated Effects (CCE) Estimator. Includes <code>xtfce2</code> and <code>xtfce2</code> to test for and estimate exponent of cross-section dependence	<a href="#">Stata Journal 18:3</a> , <a href="#">Stata Journal 21:3</a> , <a href="#">Slides</a> .
<a href="#">xtbreak</a>	Release: <a href="#">v1.4</a>	<a href="#">Release date: 27 Feb 2024</a>	Estimating and testing for many known and unknown structural breaks in time series and panel data. Work with <a href="#">Yannis Karanas</a> and <a href="#">Joachim Westerlund</a> .	<a href="#">arXiv 2110.14550</a> , <a href="#">arXiv 2211.06707</a> (Main paper), <a href="#">Slides</a> .
<a href="#">nwetregress</a>	Release: <a href="#">v0.2</a>	<a href="#">Release date: 20 Feb 2024</a>	Network Regressions in Stata with unbalanced panel data and time varying network structures or spatial weight matrices. Work with <a href="#">William Gitterer</a> and <a href="#">Mona Zekhnin</a> .	<a href="#">Slides</a>
<a href="#">xtst</a>	Release: <a href="#">v1.0</a>	<a href="#">Release date: 20 Feb 2024</a>	Testing for slope homogeneity in Stata. Work with <a href="#">Tore Bersvendsen</a> .	<a href="#">Stata Journal 21:1</a>
<a href="#">xtgranger</a>	Release: <a href="#">v1.0</a>	<a href="#">Release date: 20 Feb 2024</a>	Improved tests for Granger noncausality in panel data. Work with <a href="#">Jiaqi Xiao</a> , <a href="#">Yannis Karanas</a> , <a href="#">Arturas Juoda</a> and <a href="#">Yasits Sarafidis</a> .	<a href="#">Stata Journal 23:1</a> , <a href="#">Slides</a>
<a href="#">xtnumfac</a>	Release: <a href="#">v1.0</a>	<a href="#">Release date: 20 Feb 2024</a>	A battery of estimators for the number of common factors in time series and panel-data models. Work with <a href="#">Simon Rensy</a> .	<a href="#">Stata Journal 23:2</a> .

### Stata packages - Data Processing

Package	Version	Updated	Description
<a href="#">simulat2</a>	Release: <a href="#">v1.0</a>	<a href="#">Release date: 20 Feb 2024</a>	Enhanced and parallelised simulations in Stata.
<a href="#">mata2Tex</a>	Release: <a href="#">v0.7</a>	<a href="#">Release date: 20 Feb 2024</a>	Export Mata Matrix to LaTeX Tables.
<a href="#">htmltab2stata</a>	Release: <a href="#">v0.2</a>	<a href="#">Release date: 19 Feb 2024</a>	Converting html tables into a Stata dataset.
<a href="#">stataid</a>	Release: <a href="#">v1.1</a>	<a href="#">Release date: 19 Feb 2024</a>	Obtaining and displaying information about running Stata instances and closing Stata instances in Microsoft Windows.
<a href="#">multiStata</a>	Release: <a href="#">v1.5</a>	<a href="#">Release date: 17 Feb 2024</a>	Parallel loops in Stata (discontinued, no further bug fixing/development).

# Example

## Tables on My Github Page

- There are two tables on My GitHub page. The pages are defined using the HTML “<table>” environment.
- Let's say we want to load the first table.

```
htmltab2stata , url(https://github.com/JanDitzen) firstrow href
```

- Downloads the website, parses it for tables and loads the first table into a new Stata dataset:

The screenshot shows the Stata Data Editor interface with a table of packages. The table has columns for Package, Version, Updated, Description, and Article\_Slides. The first row is highlighted.

	Package	Version	Updated	Description	Article_Slides
1	xtcsc2			Estimation of Common Correlated Effects (CCE) Estimator. Includes xtcsc2 and xtcsc2to test for and estimate exponent of cross-section dependence	Stata Journal 18.3, Stata Journal 21.3, Slides
2	xtbreak			Estimating and testing for many known and unknown structural breaks in time series and panel data. Work with Yannis Karavias and Joakim Westerlund.	arXiv 2110.14550, arXiv 2211.06707 (Main paper), Slides
3	netxtregress			Network Regressions in Stata with unbalanced panel data and time varying network structures or spatial weight matrices. Work with William Greiser and Mored Zekhnini.	Slides
4	xtsthat			Testing for slope homogeneity in Stata. Work with Tore Bersvendsen.	Stata Journal 21.1
5	strangetest			Improved tests for Granger noncausality in panel data. Work with Jiaojiao Xiao, Yannis Karavias, Atanas Jurdas and Vasilis Sarafidis.	Stata Journal 23.1, Slides
6	strumfac			A battery of estimators for the number of common factors in time series and panel-data models. Work with Simon Reese.	Stata Journal 23.2

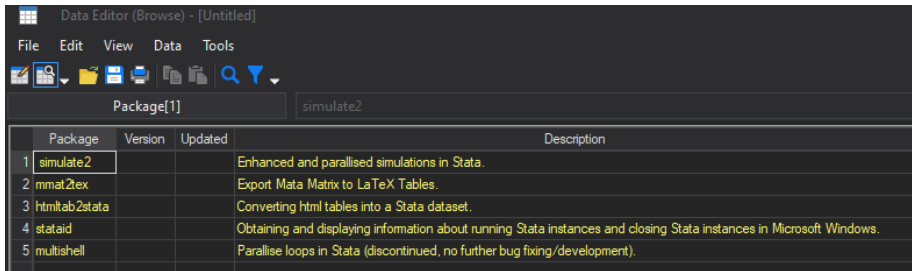
# Example

## Tables on My Github Page

- Let's say we want to load the second table.

```
htmltab2stata , url(https://github.com/JanDitzen) firstrow href tablenumber(2)
```

- Leads to:



Data Editor (Browse) - [Untitled]

File Edit View Data Tools

Package[1] simulate2

	Package	Version	Updated	Description
1	simulate2			Enhanced and paralised simulations in Stata.
2	mmat2tex			Export Mata Matrix to LaTeX Tables.
3	htmltab2stata			Converting html tables into a Stata dataset.
4	stataid			Obtaining and displaying information about running Stata instances and closing Stata instances in Microsoft Windows.
5	multishell			Parallise loops in Stata (discontinued, no further bug fixing/development).

## Conclusion

- `htmltab2stata` is a simple tool to load html Tables into Stata datasets.
- Can load websites from URLs or files.
- Can processes multiple Tables per page and load directly links.
- Limited to static pages with HTML “<table>” environment.
- How to install?

```
net install htmltab2stata , from("https://janditzen.github.io/htmltab2stata/")
```

- More info:



[jan.ditzen.net](https://jan.ditzen.net)



[GitHub](https://github.com/janditzen/htmltab2stata)