

# Advanced data visualizations with Stata

#StataViz ++

Asjad Naqvi

International Institute for Applied Systems Analysis (IIASA)  
Wirtschaftsuniversität Wien (WU)

Stata UK Conference  
10 September, 2021

## Background

- Ph.D. Economics (2007-12), New School for Social Research, New York, USA.
- Started using Stata during my M.Sc. Economics degree (2003-04)
- Worked/still working on a ton of research projects mostly using Stata
  - Overseeing several large data projects
  - Going through more code than I like
- Why do all the dataviz stuff?
  - Online COVID-19 dataviz explosion + work-from-home → curiosity
  - Huge positive response from the online community

## Why we need more #StataViz?

- Few official releases for graphs
- Still lots of options available in the default Stata structure
- Lots of great development in terms of dataviz packages:
  - [colorpalette](#), [colrspace](#), [heatmap](#) (Ben Jann), [spmap](#) (Maurizio Pisati), various packages by Nick Cox, [nwcommands](#) (Thomas Grund)

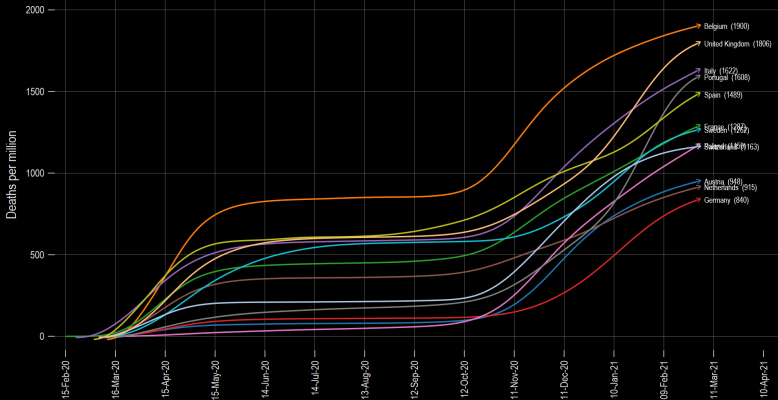


## Why we need more #StataViz?

- Few official releases for graphs
- Still lots of options available in the default Stata structure
- Lots of great development in terms of dataviz packages:
  - `colorpalette`, `colrspace`, `heatmap` (Ben Jann), `spmap` (Maurizio Pisati), various packages by Nick Cox, `nwcommands` (Thomas Grund)
- What I will discuss here:
  - What has been developed
  - What is still possible
  - What core elements Stata should still add

# colorpalette

## COVID-19 cumulative deaths per million population ( 1 Mar 21)

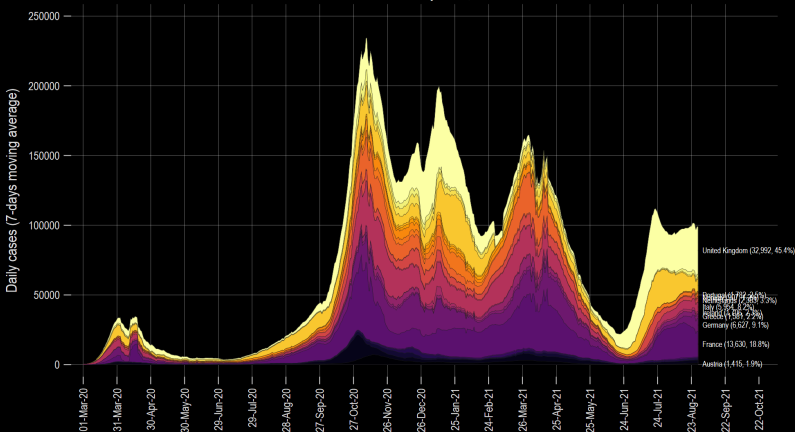


Data source: Our World in Data COVID-19 tracker. Total deaths for the last reported date given in brackets.

# colorpalette + stacked-area graphs

## New cases - European countries

Total for today = 72,670



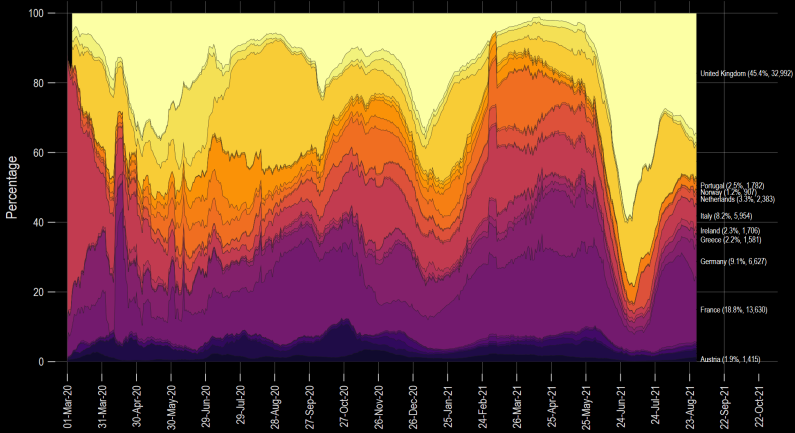
Data sources: Our World in Data, JHU, ECDC. World Bank classifications used for country groups. Top 10 countries are labeled.



# colorpalette + stacked-area graphs

## Share of new cases - European countries

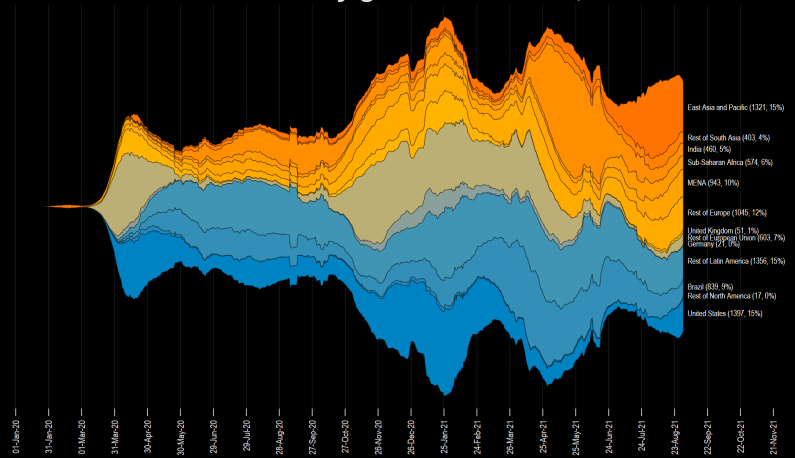
Total for today = 72,670



Data sources: Our World in Data, World Bank classifications used for country groups. Top 10 countries are labeled.

# colorpalette + stream plot

## COVID-19 daily global deaths: 9,130

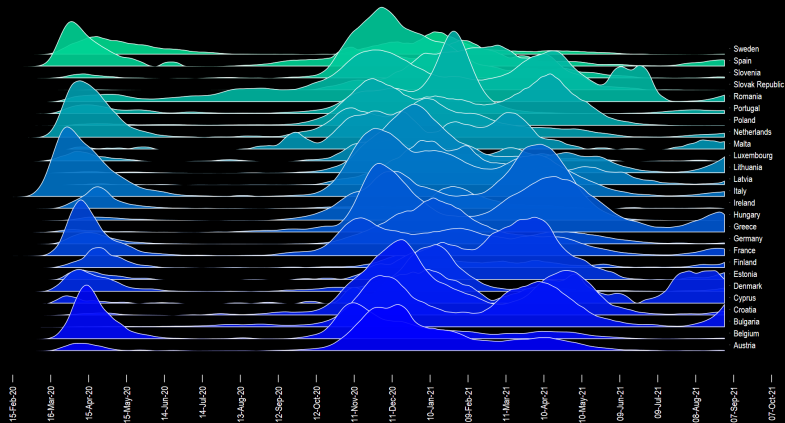


Data sources: Our World in Data, World Bank 2020 classifications used for country groups.



# Ridgeline (Joy) plots

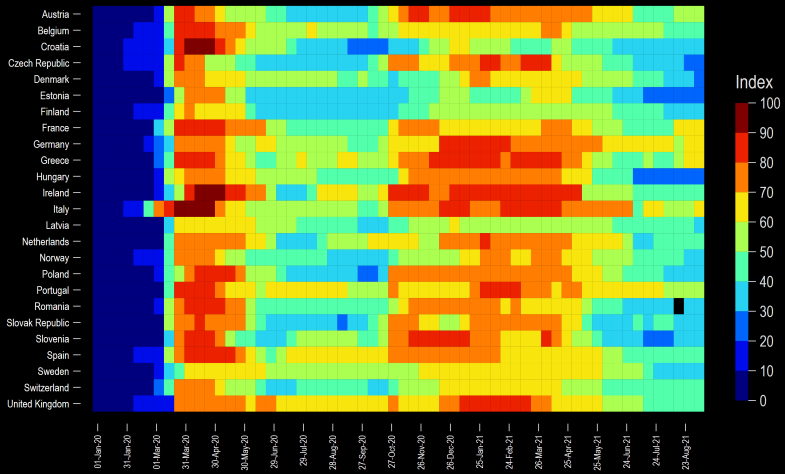
## COVID-19 daily deaths in Europe



Data sources: Our World in Data, World Bank classifications used for country groups. Each country plot is normalized by its maximum value.

# colorpalette + heatmap

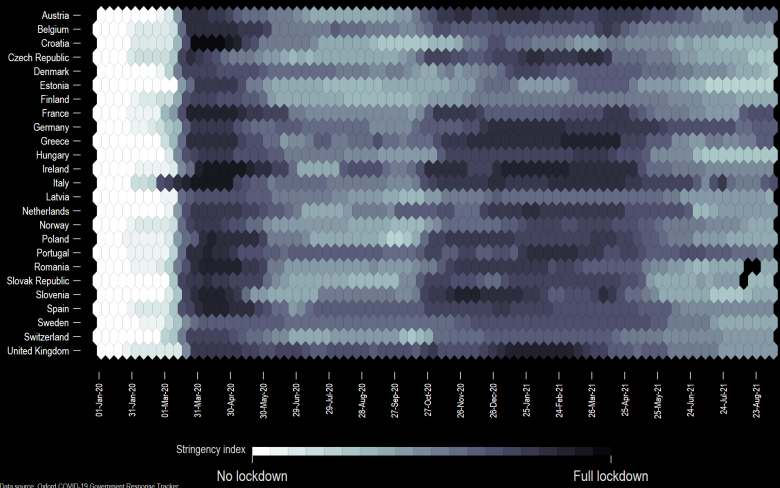
## COVID-19 Policy Stringency Index



Data source: Oxford COVID-19 Government Response Tracker

# colorpalette + heatmap

## COVID-19 Policy Stringency Index

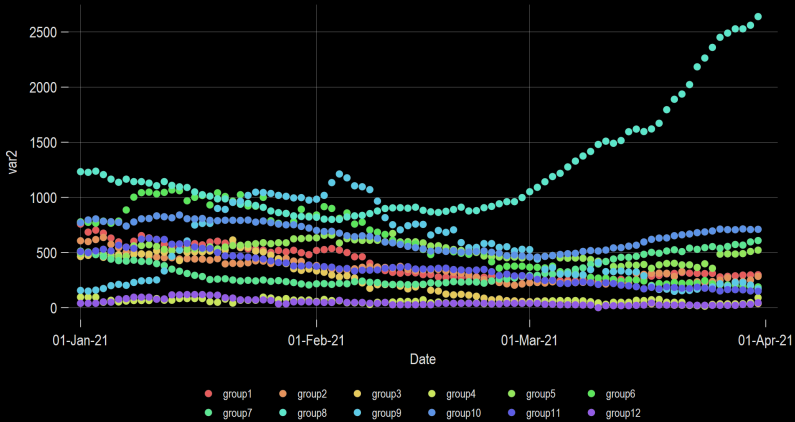


Data source: Oxford COVID-19 Government Response Tracker



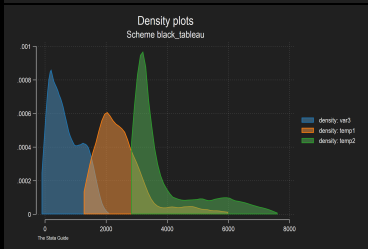
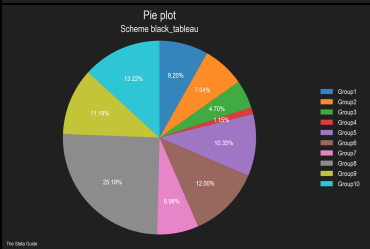
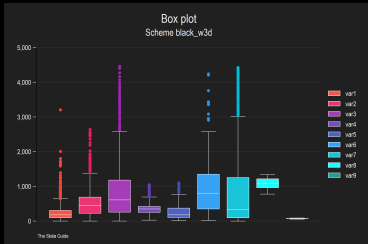
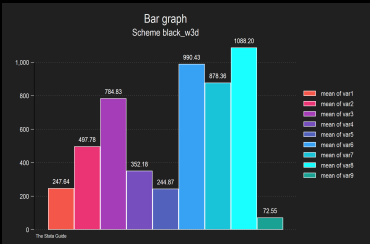
# Schemes + colorpalette

## Scatter plot Scheme neon

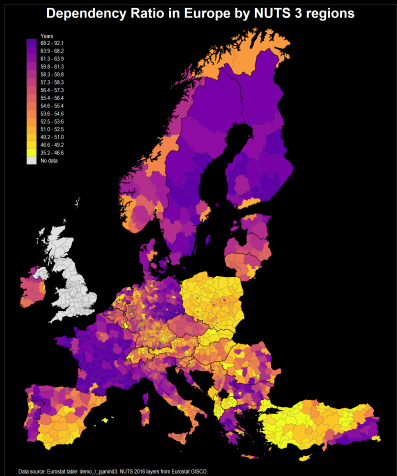
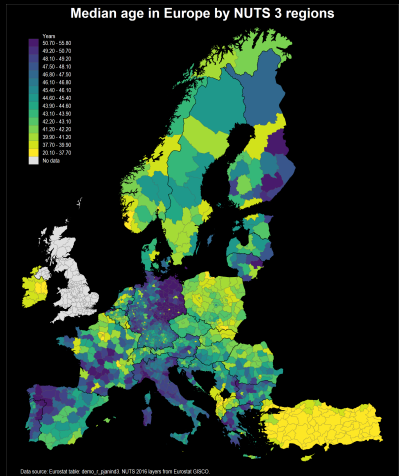


The Stata Guide

# Schemes + colorpalette



# spmap + colorpalette



## OSM + QGIS to Stata

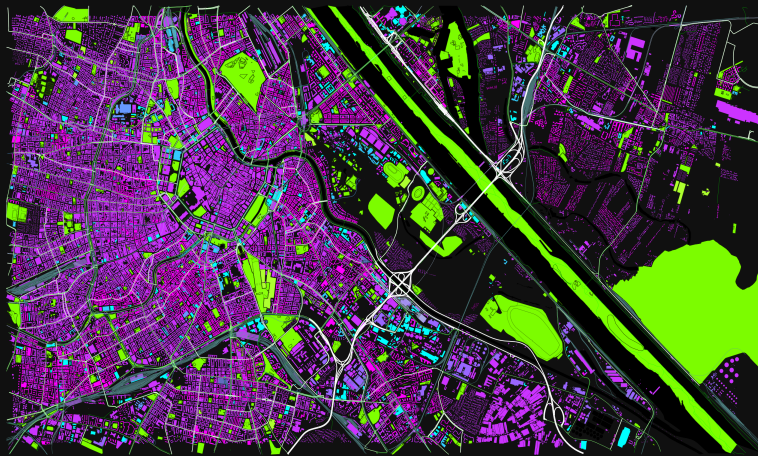
- GIS data from OpenStreetMaps (OSM)
- QGIS can be used to extract multiple layers (points, lines, polygons)
- Why do this?
  - Fully customized maps
  - Spatial layer clipping + control over bounding box
  - Allow for very rich visualizations

## Vienna





# Vienna

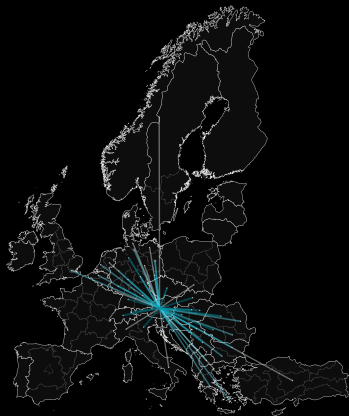


# New York



# Maps + Networks

Social Connectivity Index - Austria (70+ pctile)

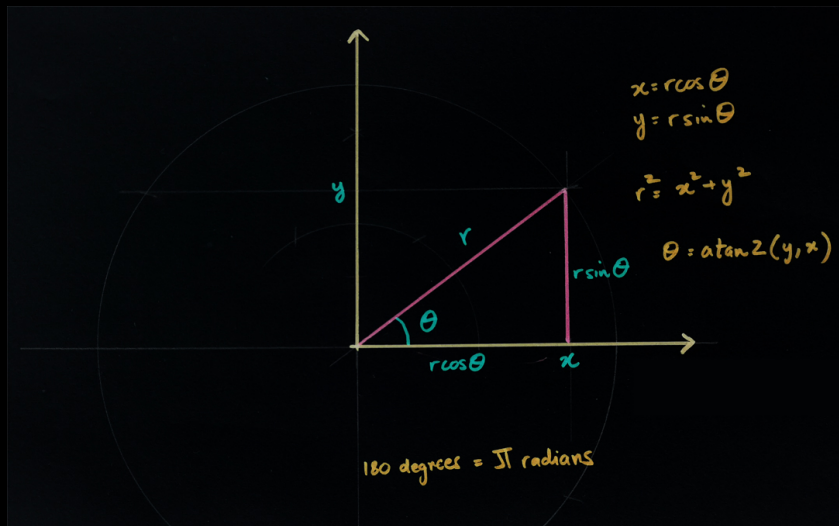


Map layers: GISCO Eurostat, SCI Facebook. Within country ties not included. Darker shade is higher SCI.

Social Connectivity Index - UK (70+ pctile)



Map layers: GISCO Eurostat, SCI Facebook. Within country ties not included. Darker shade is higher SCI.



## Spider plots

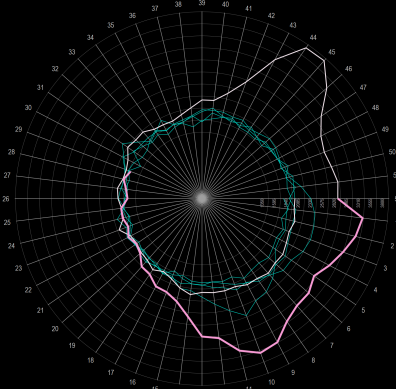
## COVID-19 policies: 29 Aug 2021



Data source: Oxford COVID-19 Government Response Tracker.

## Polar plots

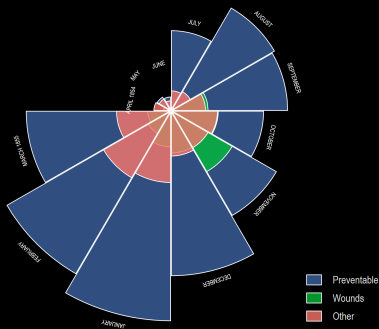
## Excess weekly deaths for 65+ - CZ



Source: Eurostat table demo\_r\_mmswk3. 2020 = light pink, 2021 = dark pink, and 2016-2019 = green shades.

## Coxcomb plots (by Florence Nightingale)

DIAGRAM OF THE CAUSES OF MORTALITY  
IN THE ARMY IN THE EAST

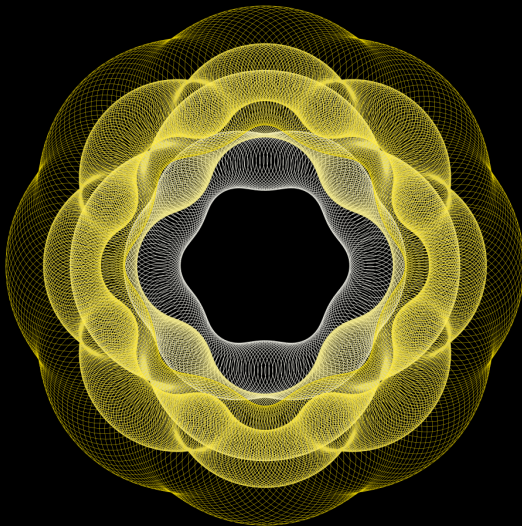


I believe in observation, measurement, and reasoning, confirmed by evidence. I believe in independent observers. I'll believe anything, however, the flimsier and more solid the evidence will have to be. Isaac Asimov, The Robot Mind (1983)

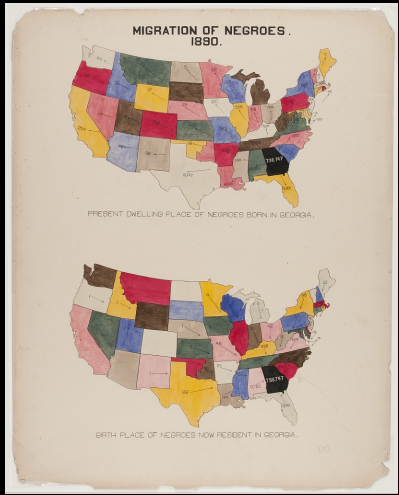




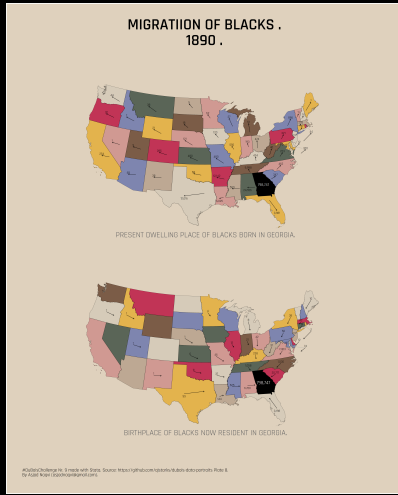
# Hypotrochoids (spiralographs), Guilloché patterns



# The Du Bois Challenge

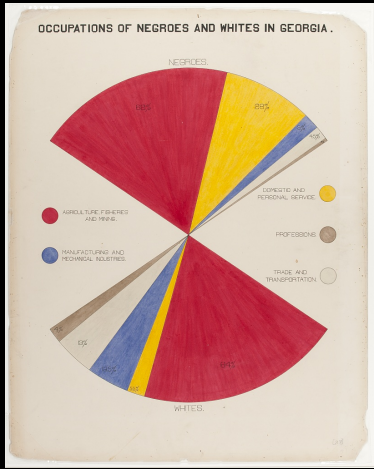


Original

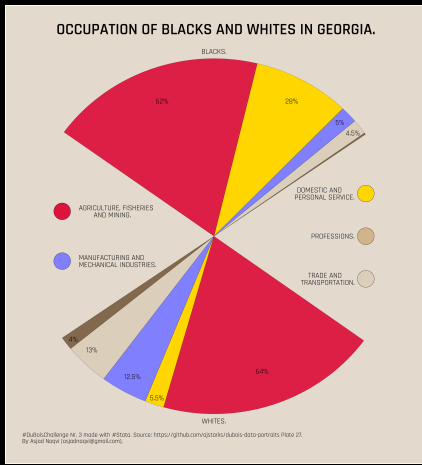


Replication

# The Du Bois Challenge

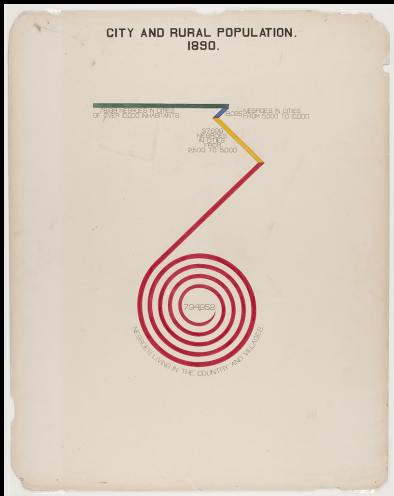


Original

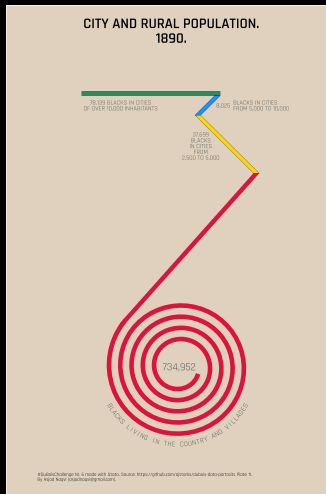


Replication

# The Du Bois Challenge



Original



Replication

## Relational and hierarchical datasets

- Both are in a network structure
- *Relational*: flows from one node to another for all layers
- *Hierarchical*: Flows are partitioned and ordered by layers

## Dummy relational data

---

From	To	Value
App	Blog	325
App	Homepage	255
App	Medium	220
App	Twitter	180
Blog	App	144
Blog	Homepage	184
Blog	Twitter	320
Blog	Website	50
Direct	App	124
Direct	Blog	175
Direct	Homepage	119
Facebook	App	124
Facebook	Blog	167
Facebook	Homepage	185

---

---

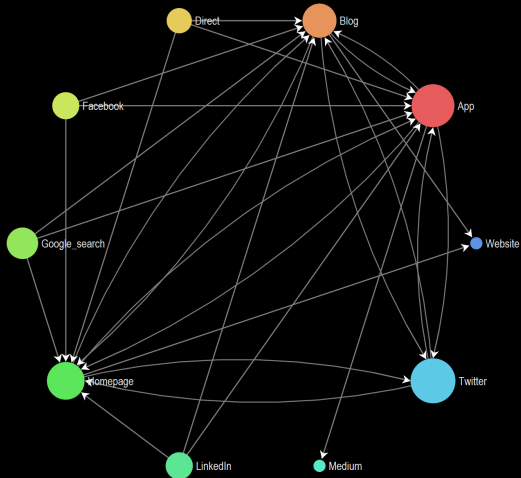
From	To	Value
Google search	App	233
Google search	Blog	252
Google search	Homepage	131
Homepage	App	375
Homepage	Blog	287
Homepage	Twitter	50
Homepage	Website	100
LinkedIn	App	185
LinkedIn	Blog	173
LinkedIn	Homepage	123
Twitter	App	485
Twitter	Blog	341
Twitter	Homepage	219

---

# A matrix representation

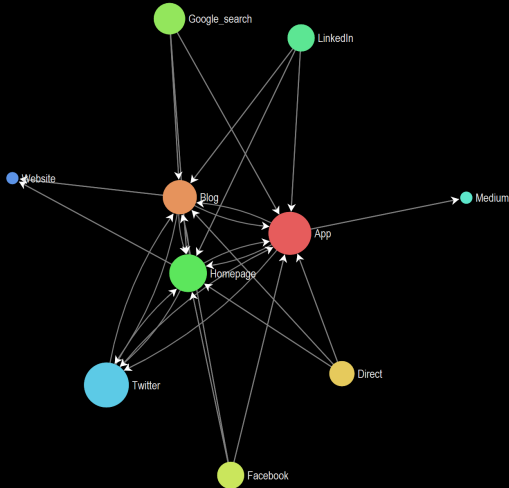


# A network representation (circle layout)





# A network representation (MDS layout)

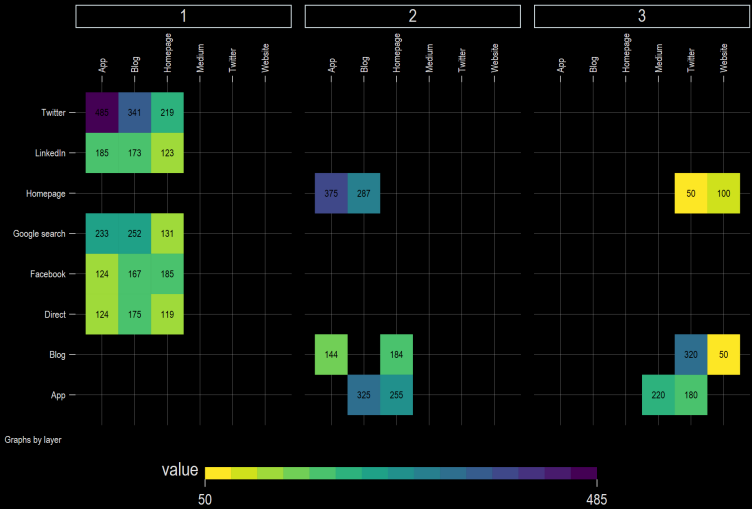


## Sample data - Hierarchical

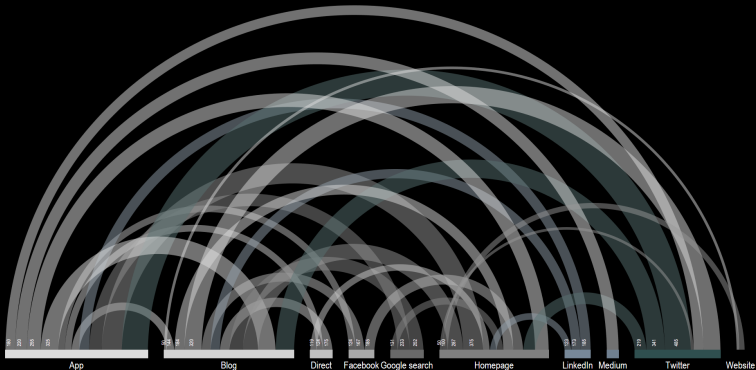
Layer	From	To	Value
1	Twitter	Blog	341
1	Twitter	Homepage	219
1	Twitter	App	485
1	Google search	Blog	252
1	Google search	Homepage	131
1	Google search	App	233
1	Facebook	Blog	167
1	Facebook	Homepage	185
1	Facebook	App	124
1	LinkedIn	Blog	173
1	LinkedIn	Homepage	123
1	LinkedIn	App	185
1	Direct	Blog	175
1	Direct	Homepage	119
1	Direct	App	124

Layer	From	To	Value
2	Blog	Homepage	184
2	Blog	App	144
2	Homepage	Blog	287
2	Homepage	App	375
2	App	Homepage	255
2	App	Blog	325
3	Homepage	Website	100
3	Homepage	Twitter	50
3	App	Twitter	180
3	App	Medium	220
3	Blog	Twitter	320
3	Blog	Website	50

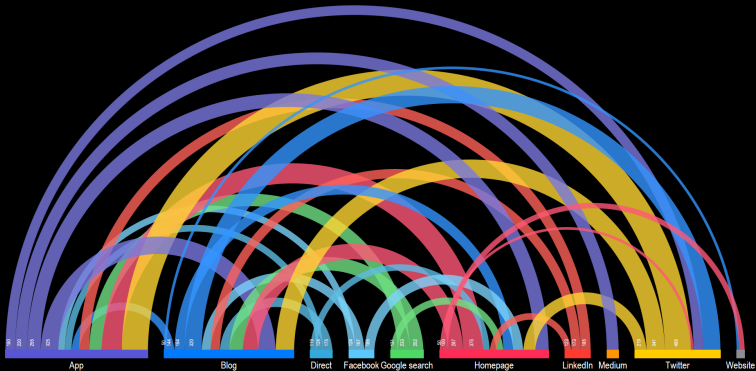
# A matrix representation



# Arc



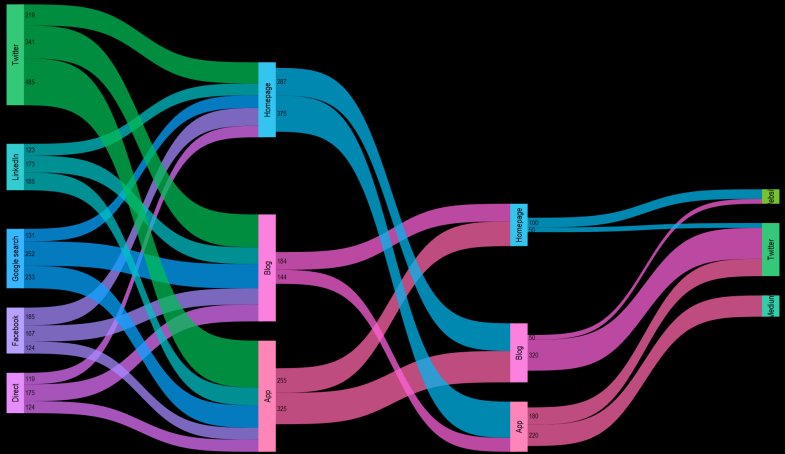
# Arc



# Chord



# Sankey



## Why do all of this in Stata?

- One can pass the DataViz stuff to more established languages:
  - R using `rcall` (Haghish 2019) or Python using `PyStata`
  - But these should be viewed more as stop-gap options
  - Develop and internalize the `dataviz` option within Stata
  - Why?



## Why do all of this in Stata?

- One can pass the DataViz stuff to more established languages:
  - R using `rcall` (Haghish 2019) or Python using `PyStata`
  - But these should be viewed more as stop-gap options
  - Develop and internalize the `dataviz` option within Stata
  - Why?
- Easy to integrate with the Stata code workflow
- Reduce the language burden on end-users
- Time investment? One-time (high) sunk costs *or*
- Some of the `#StataViz` shown above can be easily packaged for end-users

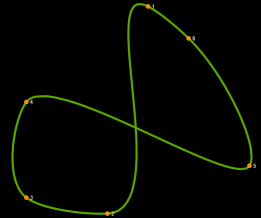
## Core features that can take Stata to the next level (wish list)

- Ability to control marker size scaling
- Line weights (like marker weights)
- Allow angles/sizes/colors to be read from variables (like mlabels)
- Color/alpha scaling for markers/lines
- Color gradients for lines/areas
- Ability to add custom markers (or increase marker pool)
- Ability to added colored text in graphs
- Ability to add images in figures
- Ability to read images (pixel data/colors)
- Ability to draw/add/subtract areas between two or more functions
- Line curvature

Thank you!

For more #Stataviz:

- [The Stata Guide](#) on Medium
- GitHub: [asjadnaqvi](#)
- Twitter: [@AsjadNaqvi](#)



Catmull-Rom spline

A lot more is in development for next year