

# HEATMAPS FOR ECONOMIC ANALYSIS

Tom Cui, Eric Zwick  
(DRAFT)

October 5, 2016

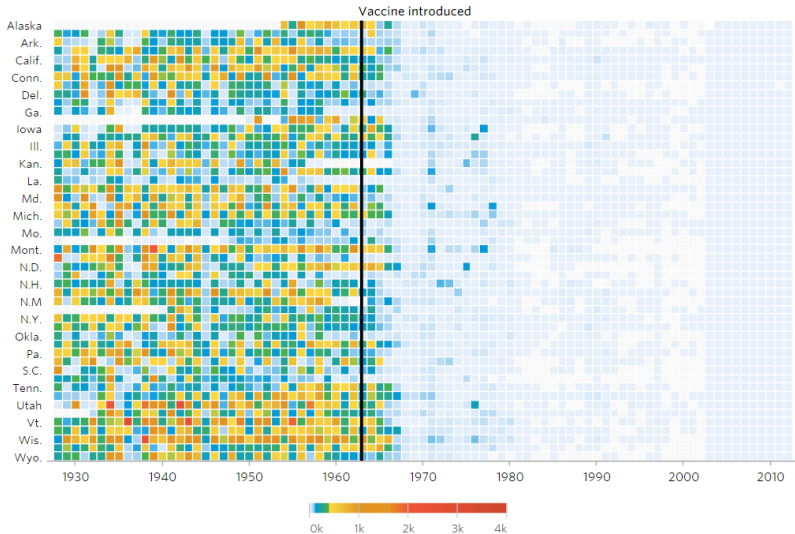
# WHAT IS A HEATMAP?

- ▶ **A two-dimensional visualization of data using colour to represent magnitude**
- ▶ Broad definition, which could be divided into
- ▶ **Embedded** heatmaps that overlay colour on an actual map or image (not covered here)
- ▶ **Matrix** heatmaps that presents a grid of values where colours differ by cell

# WHAT IS A HEATMAP?

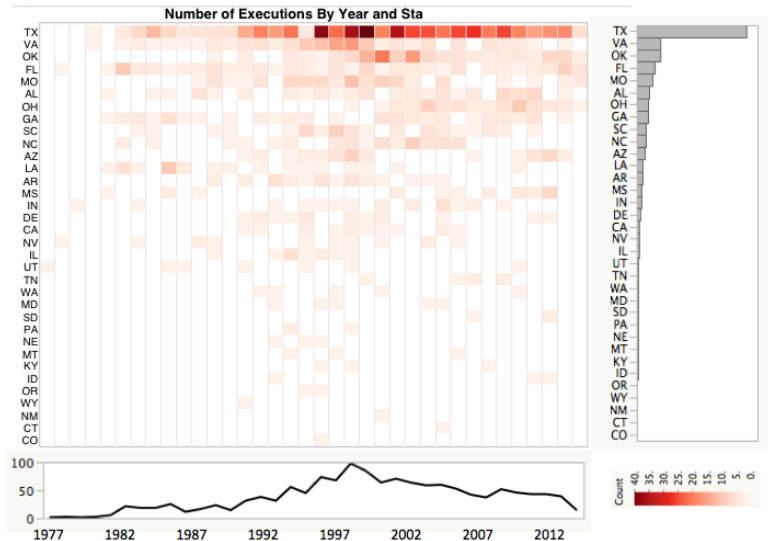
**Example:** The WSJ vaccine visualization (DeBold, Friedman 2015)

## Measles



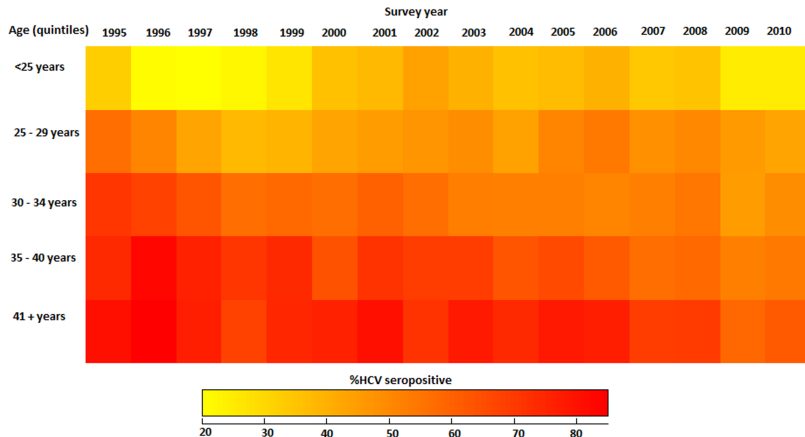
# WHAT IS A HEATMAP?

**Example:** Kaiser Fung's executions data



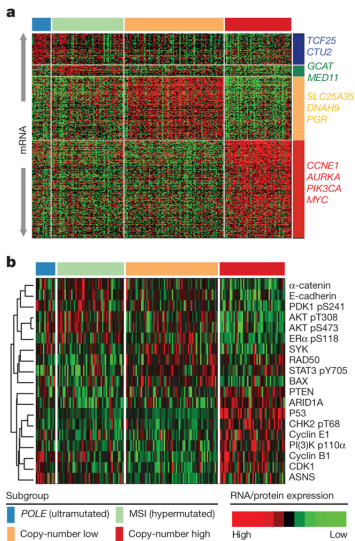
# WHAT IS A HEATMAP?

**Example (Bad):** A “quilt plot” of Hep C prevalence (Wand et al)



# WHAT IS A HEATMAP?

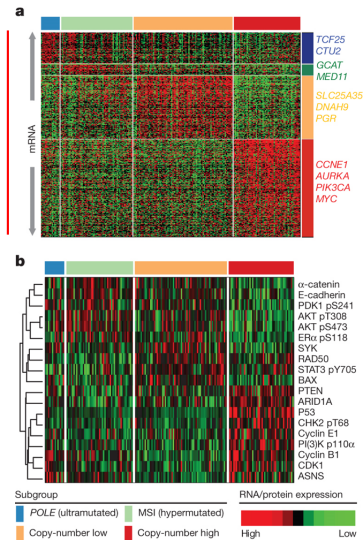
**Example:** Plotting gene expression data over samples (TCGN 2013)



# WHAT IS A HEATMAP?

**Example:** Plotting gene expression data over samples (TCGN 2013)

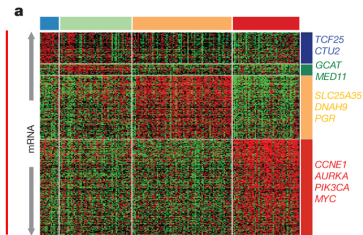
Each row ( $\sim 1500$ )  
is one gene



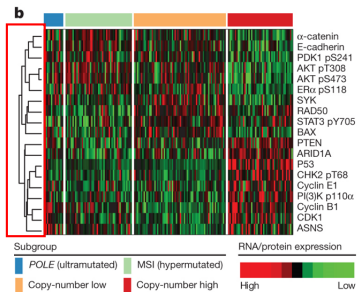
# WHAT IS A HEATMAP?

**Example:** Plotting gene expression data over samples (TCGN 2013)

Each row ( $\sim 1500$ )  
is one gene



**Dendrogram**



Each row is  
a protein



# WHAT IS A HEATMAP?

Some takeaways from these examples:

- ▶ The axes change the interpretation  
(1) - (3) use time as the X and factors as the Y, (4) uses factors for both

# WHAT IS A HEATMAP?

Some takeaways from these examples:

- ▶ The axes change the interpretation  
(1) - (3) use time as the X and factors as the Y, (4) uses factors for both
- ▶ Good representation of high-dimensional data  
(4) is an extreme example of this, but common in bioinformatics

# WHAT IS A HEATMAP?

Some takeaways from these examples:

- ▶ The axes change the interpretation  
(1) - (3) use time as the X and factors as the Y, (4) uses factors for both
- ▶ Good representation of high-dimensional data  
(4) is an extreme example of this, but common in bioinformatics
- ▶ Permuting axis order improves interpretation  
(2) sorts Y by total count over the sampling period, (4) uses cluster analysis (recall dendrogram)

# SETTING UP A HEATMAP FOR ECONOMICS

- ▶ In an ideal world, we could derive causal effects in a model  $Y = g(W)$  using exogeneous assignment of  $W$  and observing the entire support of  $W$

# SETTING UP A HEATMAP FOR ECONOMICS

- ▶ In an ideal world, we could derive causal effects in a model  $Y = g(W)$  using exogeneous assignment of  $W$  and observing the entire support of  $W$
- ▶ Big data makes the latter easier. Former still hard!
- ▶ Hence research designs that exploit a policy introduction or kink are popular

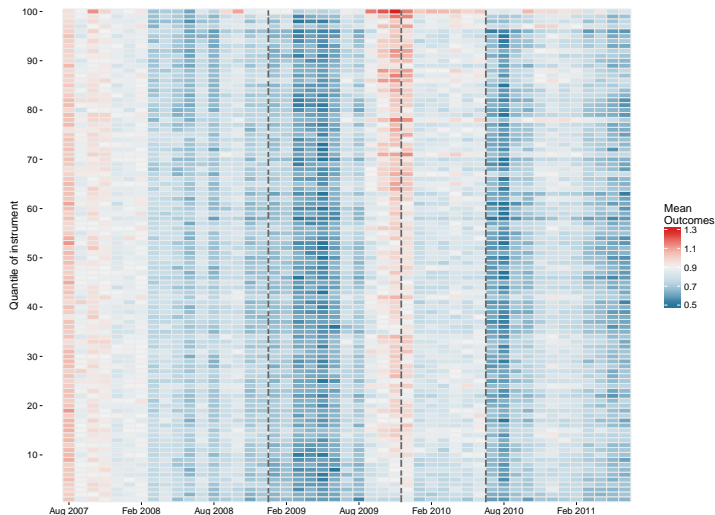
# SETTING UP A HEATMAP FOR ECONOMICS

- ▶ In an ideal world, we could derive causal effects in a model  $Y = g(W)$  using exogeneous assignment of  $W$  and observing the entire support of  $W$
- ▶ Big data makes the latter easier. Former still hard!
- ▶ Hence research designs that exploit a policy introduction or kink are popular

Now consider a heatmap where time is on the X axis (**showing the policy introduction**) and where  $W$ , a variable of interest or one related to a latent factor is binned on the Y axis (**showing the support of  $W$** )

# SETTING UP A HEATMAP FOR ECONOMICS

**Example:** Scaled house sales in a heatmap sorted by FTHB exposure, from Berger, Turner, Zwick (2016)



# SETTING UP A HEATMAP FOR ECONOMICS

Using earlier takeaways:

- ▶ The axes change the interpretation  
Placing time on  $X$  and an instrument of  $W$  on  $Y$  implies this heatmap is a visualization of nonparametric regression



# SETTING UP A HEATMAP FOR ECONOMICS

Using earlier takeaways:

- ▶ The axes change the interpretation  
Placing time on  $X$  and an instrument of  $W$  on  $Y$  implies this heatmap is a visualization of nonparametric regression
- ▶ Good representation of high-dimensional data  
Around 8600 ZIPs binned into 100 percentiles

# SETTING UP A HEATMAP FOR ECONOMICS

Using earlier takeaways:

- ▶ The axes change the interpretation  
Placing time on X and an instrument of W on Y implies this heatmap is a visualization of nonparametric regression
- ▶ Good representation of high-dimensional data  
Around 8600 ZIPs binned into 100 percentiles
- ▶ Permuting axis order improves interpretation  
Y axis sorted to be increasing in W's instrument, and figure tells us the effect of W on Y is positive in a linear model

# SETTING UP A HEATMAP FOR ECONOMICS

Extensions:

- ▶ Quantiles of instrument on  $X$ , other variables on  $Y$ , plotting means  
= **Covariate balance check**

# SETTING UP A HEATMAP FOR ECONOMICS

Extensions:

- ▶ Quantiles of instrument on X, other variables on Y, plotting means  
= **Covariate balance check**
- ▶ Time on X, portfolios on Y, plotting market-adjusted returns  
= **Financial event study**

# SETTING UP A HEATMAP FOR ECONOMICS

## Extensions:

- ▶ Quantiles of instrument on X, other variables on Y, plotting means  
= **Covariate balance check**
- ▶ Time on X, portfolios on Y, plotting market-adjusted returns  
= **Financial event study**
- ▶ Time on X, generation on Y, plotting average of a simulated policy function  
= **OLG model dynamics**

# SETTING UP A HEATMAP FOR ECONOMICS

## Extensions:

- ▶ Quantiles of instrument on X, other variables on Y, plotting means  
= **Covariate balance check**
- ▶ Time on X, portfolios on Y, plotting market-adjusted returns  
= **Financial event study**
- ▶ Time on X, generation on Y, plotting average of a simulated policy function  
= **OLG model dynamics**
- ▶ Index determining policy entry on X, quantiles of dependent variable on Y, plotting obs. counts in bin  
= **Fuzzy RDD**

# SETTING UP A HEATMAP FOR ECONOMICS

Extensions:

- ▶ Quantiles of instrument on X, other variables on Y, plotting means  
= **Covariate balance check**
- ▶ Time on X, portfolios on Y, plotting market-adjusted returns  
= **Financial event study**
- ▶ Time on X, generation on Y, plotting average of a simulated policy function  
= **OLG model dynamics**
- ▶ Index determining policy entry on X, quantiles of dependent variable on Y, plotting obs. counts in bin  
= **Fuzzy RDD**

and so on.

# The heatmapEco package



# THE HEATMAPECO PACKAGE

- ▶ **Many** programs for creating heatmaps exist

So why another package?

# THE HEATMAP<sub>ECO</sub> PACKAGE

- ▶ **Many** programs for creating heatmaps exist
  - ▶ Stata `twoway contour`, `hmap`
  - ▶ R base, `gplots`, `ggplot2`, `d3heatmap` ...
  - ▶ Matlab and Python `matplotlib`

So why another package?

# THE HEATMAPECO PACKAGE

- ▶ **Many** programs for creating heatmaps exist
  - ▶ Stata `twoway contour`, `hmap`
  - ▶ R `base`, `gplots`, `ggplot2`, `d3heatmap` ...
  - ▶ Matlab and Python `matplotlib`

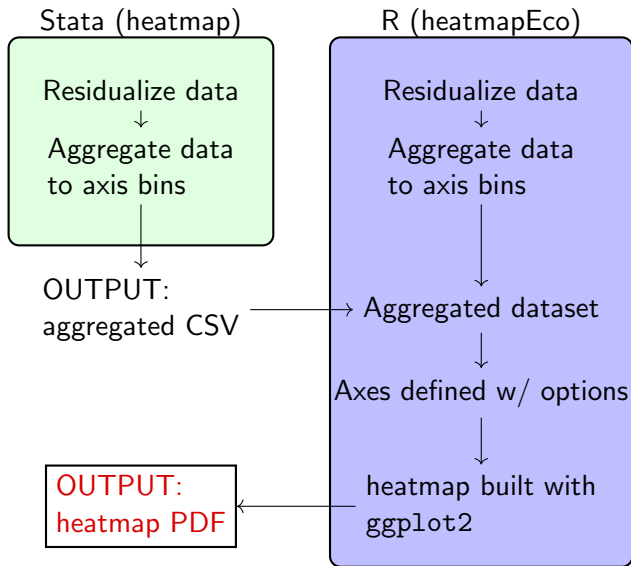
So why another package?

- ▶ `heatmapEco` makes it easy building informative heatmaps by
  - ▶ **Focusing on axis setup as a design framework;**
  - ▶ **Computing relevant axis permutations;**
  - ▶ **Executing prerequisite data cleaning.**

# THE HEATMAPECO PACKAGE

- ▶ Complicated heatmaps like TCGN's are also quite uncomplicated; they are literally a projection of some tabular data
- ▶ In other words, the data loaded in is a 373x1500 matrix. The values are then standardized, variables are clustered and given a colour
- ▶ But instead data may need to be aggregated, reshaped; axes relabelled; colour palettes adjusted to show significant results
- ▶ heatmapEco combines R packages to simplify these changes and adds design features of its own

# THE HEATMAP<sub>E</sub>CO PACKAGE



# HEATMAPECO AXES

- ▶ Currently, X axis can be set up as:
  - ▶ An **index axis** over numeric values (income, policy thresholds)
  - ▶ A **time axis** where time strings are converted into valid axis values by the package

# HEATMAP ECO AXES

- ▶ Currently, X axis can be set up as:
  - ▶ An **index axis** over numeric values (income, policy thresholds)
  - ▶ A **time axis** where time strings are converted into valid axis values by the package
- ▶ Currently, Y axis can be set up as:
  - ▶ A **factor axis** where each entry is some (aggregated) grouping
  - ▶ A **quantile axis** where a continuous instrument is split into N quantiles

# HEATMAP ECO AXES

- ▶ Currently, X axis can be set up as:
  - ▶ An **index axis** over numeric values (income, policy thresholds)
  - ▶ A **time axis** where time strings are converted into valid axis values by the package
- ▶ Currently, Y axis can be set up as:
  - ▶ A **factor axis** where each entry is some (aggregated) grouping
  - ▶ A **quantile axis** where a continuous instrument is split into N quantiles

Currently output is in landscape letter format, but ultimately axis placement should be arbitrary and portrait format heatmaps possible



# HEATMAP<sub>ECO</sub> AGGREGATION

In R the aggregation process is inputted using a pseudo-formula

$$Z \sim \text{CrS}(Y, \text{ID}, w) : X(t)$$

where

- ▶ Z is the dependent variable, or the fill variable
- ▶ Y is the factor independent variable or a continuous instrument to be binned
- ▶ X is the index or time axis
- ▶ t allows time varying Y to be sorted on its values at a time t, (**use caution**)
- ▶ ID is the individual identifier, either unique or unique with t
- ▶ w are quantile weights

# HEATMAP<sub>ECO</sub> AGGREGATION

In R the aggregation process is inputted using a pseudo-formula

$$Z \sim \text{CrS}(Y, \text{ID}, w) : X(t)$$

where

- ▶ Z is the dependent variable, or the fill variable
- ▶ Y is the factor independent variable or a continuous instrument to be binned
- ▶ X is the index or time axis
- ▶ t allows time varying Y to be sorted on its values at a time t, (**use caution**)
- ▶ ID is the individual identifier, either unique or unique with t
- ▶ w are quantile weights

In Stata the syntax is

```
heatmap Z Y X [weights], id(varname) [t_sort(string)]
```

## HEATMAPECO AGGREGATION

- ▶ Note that, in R, an anonymous function could be passed as an argument
- ▶ This means the aggregation function argument `grp.func` can take many forms, so long as a summary function is involved

## HEATMAP ECO AGGREGATION

- ▶ Note that, in R, an anonymous function could be passed as an argument
- ▶ This means the aggregation function argument `grp.func` can take many forms, so long as a summary function is involved
- ▶ E.g. take the median of a quantile-month bin. Or take the log transform of that median
- ▶ Or add control flow; if data censored, first remove censored data and output log median of what remains
- ▶ Stata's aggregation features are much less rich: every collapse function could be inputted into `grpfunc`

# HEATMAP ECO RESIDUALIZATION

Both dependent and independent variables (fill and Y axis) can be first residualized according to a model

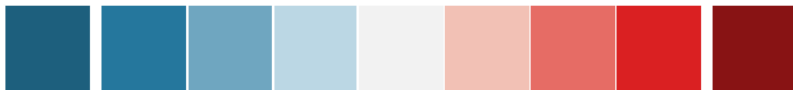
$$Y = \beta W + D\theta + F\psi + X\gamma + \varepsilon$$

Where D, F are fixed effects and X are controls.

Stata implementation uses `base areg`. R implementation uses `plm` or `lfe` (TODO)

# COLOUR PALETTES

Standard divergent color palette



Semi-sequential palette for count data



- ▶ On standard palette, far two shades reserved for outlier detection: binned values above the  $1.5 + \text{IQR}$  range are considerably darker
- ▶ Standard colors are not equally spaced: distribution below median take longer to get to dark blue hues. This is to emphasize “Ashenfelter dips”
- ▶ Count data palette is ColorBrewer YlOrBr, with high outliers and a muted hue to deemphasize data censored by 0 (by default)

## heatmapEco Examples

# WSJ REPLICATION

Download data from Project Tycho. The cleaning in R:

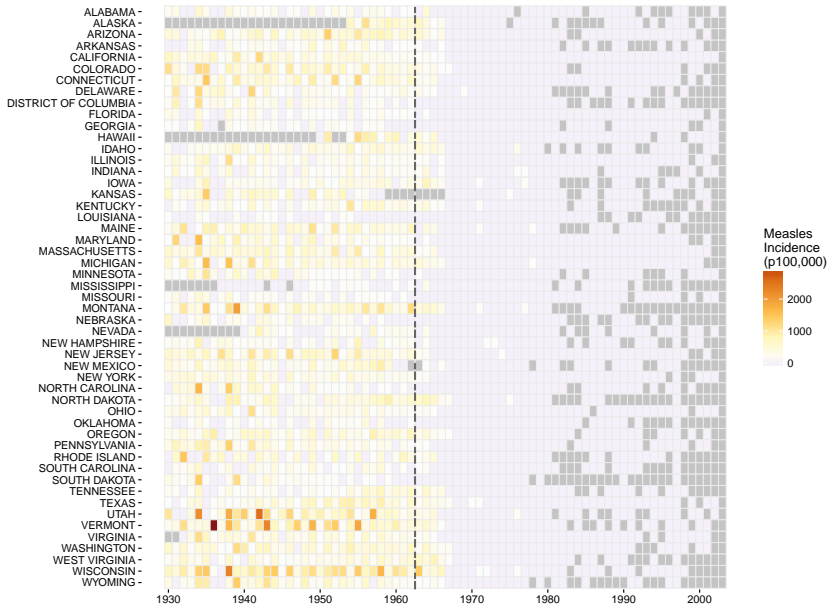
```
library(data.table)
obj <- melt(fread("MEASLES_Incidence_1930-2003.csv"),
            c("YEAR", "WEEK"))
obj[, value := as.numeric(value)]
```

Calling heatmapEco:

```
nasum <- function(...)
  if (all(is.na(...))) NA else sum(..., na.rm=TRUE)
heatmapEco(value ~ CrS(variable,variable):YEAR, obj,
t.fmt="%Y", t.per="year", pol.break=c("Jan 1963"),
grp.func=nasum, count=T, factor.ax=T, outliers=T, split.x=10,
zlab="Measles Incidence (p100,000)", save="measlesRep.pdf")
```



# WSJ REPLICATION



# WSJ REPLICATION

Line by line:

# WSJ REPLICATION

Line by line:

- ▶ `heatmapEco(value ~ CrS(variable,variable):YEAR,obj,`  
Inputs formula for aggregation and dataset
- ▶ `t.fmt="%Y", t.per="year", pol.break=c("Jan 1963"),`  
Data object, time is in pure "year" format, policy line date

# WSJ REPLICATION

Line by line:

- ▶ `heatmapEco(value ~ CrS(variable,variable):YEAR,obj,`  
Inputs formula for aggregation and dataset
- ▶ `t.fmt="%Y", t.per="year", pol.break=c("Jan 1963"),`  
Data object, time is in pure "year" format, policy line date
- ▶ `grp.func=naum [naum <- function(...)`  
`if (all(is.na(...))) NA else sum(..., na.rm=TRUE)]`  
Grouping function is summation, excluding NAs (a year with NAs is inputted as NA, grayed out)

# WSJ REPLICATION

Line by line:

- ▶ `heatmapEco(value ~ CrS(variable,variable):YEAR,obj,`  
Inputs formula for aggregation and dataset
- ▶ `t.fmt="%Y", t.per="year", pol.break=c("Jan 1963"),`  
Data object, time is in pure "year" format, policy line date
- ▶ `grp.func=naum [naum <- function(...)`  
`if (all(is.na(...))) NA else sum(..., na.rm=TRUE)]`  
Grouping function is summation, excluding NAs (a year with NAs is inputted as NA, grayed out)
- ▶ `count=T, factor.ax=T, outliers=T, split.x=10,`  
Use the count colour palette; the Y-axis are state factors; turn on outlier perception; X tick every ten units

# WSJ REPLICATION

Line by line:

- ▶ `heatmapEco(value ~ CrS(variable,variable):YEAR,obj,`  
Inputs formula for aggregation and dataset
- ▶ `t.fmt="%Y", t.per="year", pol.break=c("Jan 1963"),`  
Data object, time is in pure "year" format, policy line date
- ▶ `grp.func=NASUM [NASUM <- function(...)`  
`if (all(is.na(...))) NA else sum(..., na.rm=TRUE)]`  
Grouping function is summation, excluding NAs (a year with NAs is inputted as NA, grayed out)
- ▶ `count=T, factor.ax=T, outliers=T, split.x=10,`  
Use the count colour palette; the Y-axis are state factors; turn on outlier perception; X tick every ten units
- ▶ `zlab="Measles Incidence (p100,000)",save="measlesRep.pdf")`  
Policy line, labels, output location.

# WSJ REPLICATION

Line by line:

- ▶ `heatmapEco(value ~ CrS(variable,variable):YEAR,obj,`  
Inputs formula for aggregation and dataset
- ▶ `t.fmt="%Y", t.per="year", pol.break=c("Jan 1963"),`  
Data object, time is in pure "year" format, policy line date
- ▶ `grp.func=naSUM [naSUM <- function(...)`  
`if (all(is.na(...))) NA else sum(..., na.rm=TRUE)]`  
Grouping function is summation, excluding NAs (a year with NAs is inputted as NA, grayed out)
- ▶ `count=T, factor.ax=T, outliers=T, split.x=10,`  
Use the count colour palette; the Y-axis are state factors; turn on outlier perception; X tick every ten units
- ▶ `zlab="Measles Incidence (p100,000)",save="measlesRep.pdf")`  
Policy line, labels, output location.

Overall: **9 lines of code w/ data.table**

- ▶ **9 lines fewer** than base w/ `heatmap.2`
- ▶ **25 lines fewer** than pure `ggplot2`

# THE BERGER, TURNER, ZWICK HEATMAP

Let's call the program from Stata this time

```
heatmap y3_trim fthomebuyers_filingunits_2000 mdate ///  
        [aw=totalhsales_base], n(100) id(zip) tperiod(yearmon) ///  
        ylabel(10) polbreak(Jan 2009, Dec 2009, Jul 2010) ///  
        save(BTZRep.pdf)
```

- Default group function is mean, but the quantiles are weighted



# THE BERGER, TURNER, ZWICK HEATMAP

Let's call the program from Stata this time

```
heatmap y3_trim fthomebuyers_filingunits_2000 mdate ///  
        [aw=totalhsales_base], n(100) id(zip) tperiod(yearmon) ///  
        ylabel(10) polbreak(Jan 2009, Dec 2009, Jul 2010) ///  
        save(BTZRep.pdf)
```

- ▶ Default group function is mean, but the quantiles are weighted
- ▶ Each column is a month, labelled appropriately

# THE BERGER, TURNER, ZWICK HEATMAP

Let's call the program from Stata this time

```
heatmap y3_trim fthomebuyers_filingunits_2000 mdate ///  
        [aw=totalhsales_base], n(100) id(zip) tperiod(yearmon) ///  
        ylabel(10) polbreak(Jan 2009, Dec 2009, Jul 2010) ///  
        save(BTZRep.pdf)
```

- ▶ Default group function is mean, but the quantiles are weighted
- ▶ Each column is a month, labelled appropriately
- ▶ `polbreak()` interprets time strings and adds policy lines accordingly

# THE BERGER, TURNER, ZWICK HEATMAP

Let's call the program from Stata this time

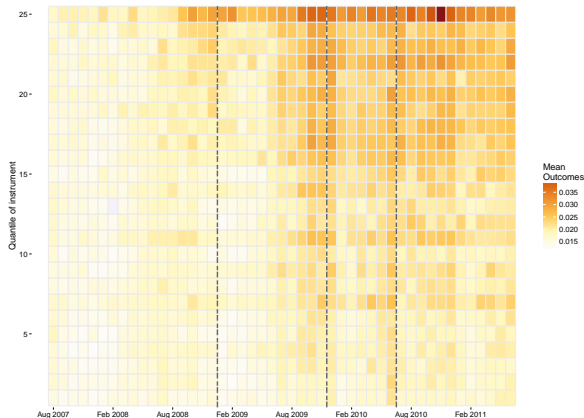
```
heatmap y3_trim fthomebuyers_filingunits_2000 mdate ///  
        [aw=totalhsales_base], n(100) id(zip) tperiod(yearmon) ///  
        ylabel(10) polbreak(Jan 2009, Dec 2009, Jul 2010) ///  
        save(BTZRep.pdf)
```

- ▶ Default group function is mean, but the quantiles are weighted
- ▶ Each column is a month, labelled appropriately
- ▶ `polbreak()` interprets time strings and adds policy lines accordingly
- ▶ `ylabel(n)` divides y-axis labels into `n` even intervals

# THE BERGER, TURNER, ZWICK HEATMAP

Another perspective: check the standard errors on the mean estimates over a coarser partition

```
heatmap y3_trim fthomebuyers_filingunits_2000 mdate ///  
[aw=totalhsales_base], n(25) id(zip) tperiod(yearmon) ///  
grpfunc(sem) ylabel(5) count out ///  
polbreak(Jan 2009, Dec 2009, Jul 2010) save(BTZRep_se.pdf)
```



# Conclusions

# WHEN NOT TO USE HEATMAPS

- ▶ Heatmaps are not a panacea: there is a tradeoff between
  - ▶ Higher density of effectively presented data;
  - ▶ Information lost in using colours, instead of geometric shapes, to represent change
- ▶ It is also unclear how heatmaps can display uncertainty of statistics plotted in each bin, e.g. confidence intervals
- ▶ A good argument for a package that simplifies heatmap creation — the less time spent making a visualization, the less likely one gets overattached to one when a better solution exists

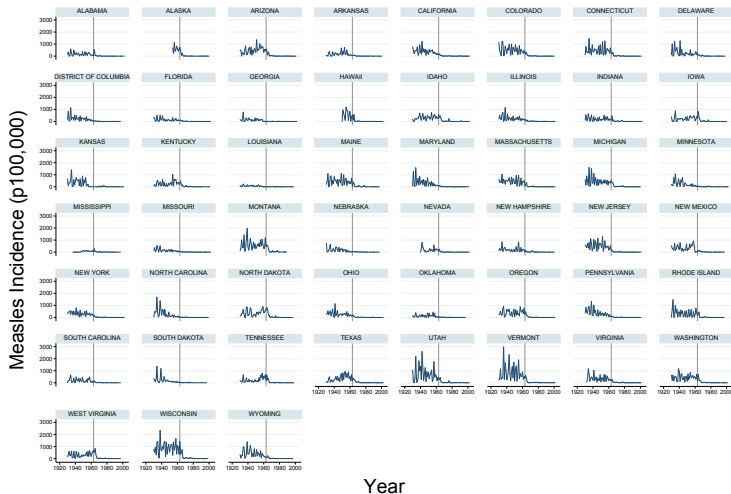
# WHEN NOT TO USE HEATMAPS

A good heuristic (define  $Z$  as the variable plotted with colour):

- ▶ Plotting quantiles on the  $Y$  axis: How much clarity is gained relative to overlapping line graphs split by  $Y$ ? What information is lost?
- ▶ Plotting a factor variable on the  $Y$  axis: How much clarity is gained relative to a small multiple plot split by  $Y$ ? What information is lost?

# WHEN NOT TO USE HEATMAPS

## Example: Measles vaccine revisited





# WHEN NOT TO USE HEATMAPS

**Example:** visualizing positive assortative matching

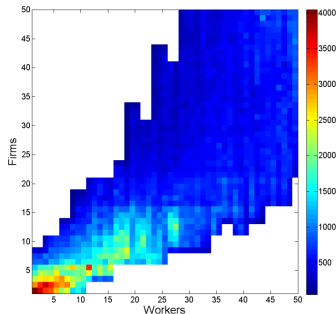
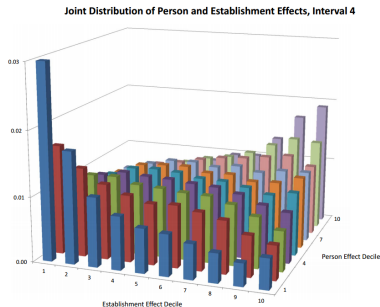


Figure 6: Estimated Match Density.

(L: Card, Heining & Kline (2012); R: Hagedorn, Law & Manovskii (2016))  
2016 How would the interpretation change if the visualization was instead overlaying many marginals over each other? Small multiples of marginals?

# FUTURE UPDATES

- ▶ Easy addition of side plots to the heatmap (a histogram on both axes, time series, bar plot of differences over two periods. . . )
- ▶ Syntax revisions
- ▶ Let either axis support variables belonging in one of four types (time, factor, quantile, index)
- ▶ Variable dimensions for heatmap cells (for uneven discretizations of a continuous variable)
- ▶ ???

# REFERENCES I

- Berger, David, Nicholas Turner, and Eric Zwick.** 2016. "Stimulating Housing Markets." *Working Paper*.
- Card, David, Jörg Heining, and Patrick Kline.** 2012. "Workplace heterogeneity and the rise of West German wage inequality." National Bureau of Economic Research.
- DeBold, Tynan, and Dov Friedman.** 2015. "Battling Infectious Diseases in the 20th Century: The Impact of Vaccines." *The Wall Street Journal*, , (11).
- Eisen, Michael B, Paul T Spellman, Patrick O Brown, and David Botstein.** 1998. "Cluster analysis and display of genome-wide expression patterns." *Proceedings of the National Academy of Sciences*, 95(25): 14863–14868.
- Fung, Kaiser.** n.d.. "Advocacy graphics." [http://junkcharts.typepad.com/junk\\_charts/2014/04/advocacy-graphics.html](http://junkcharts.typepad.com/junk_charts/2014/04/advocacy-graphics.html), Accessed: 2016-03-14.
- Hagedorn, Marcus, Tzuo Hann Law, and Iouri Manovskii.** 2016. "Identifying equilibrium models of labor market sorting."
- Network, Cancer Genome Atlas Research, et al.** 2013. "Integrated genomic characterization of endometrial carcinoma." *Nature*, 497(7447): 67–73.
- Wand, Handan et. al.** 2014. "Quilt Plots: A Simple Tool for the Visualisation of Large Epidemiological Data." *PLOS One*, , (11).

Thanks!