# Interfacing with data
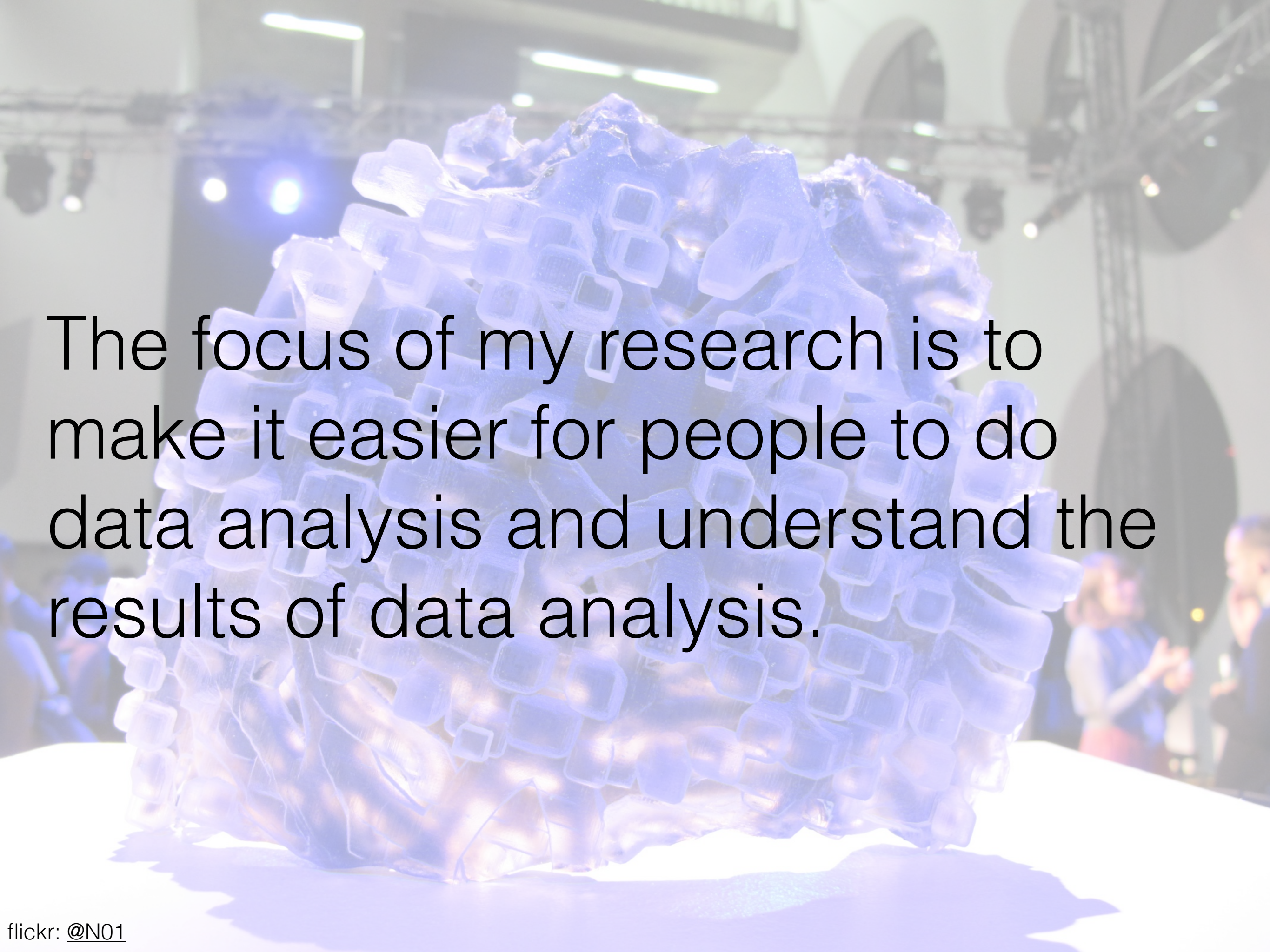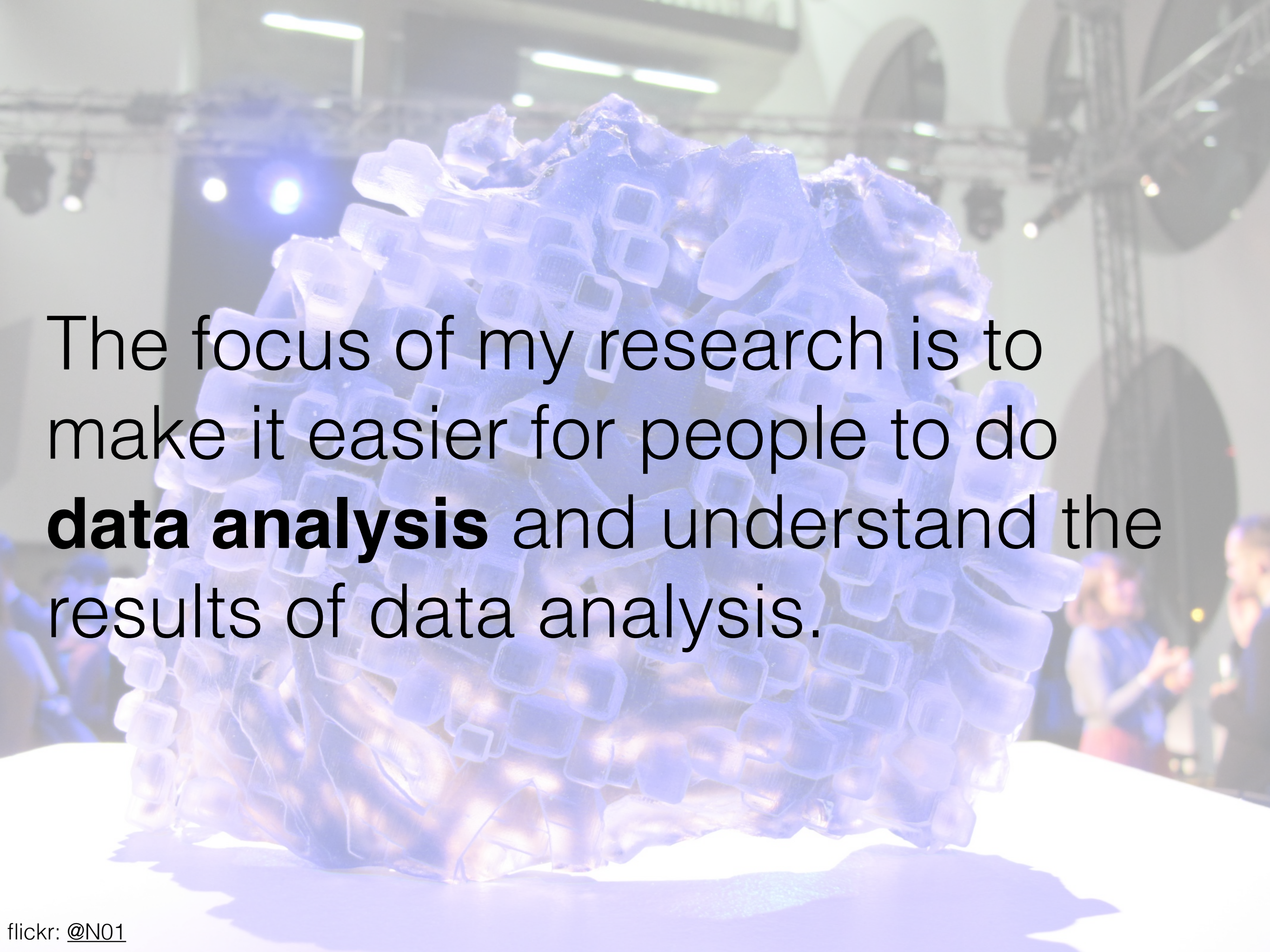
Amelia McNamara (@AmeliaMN)
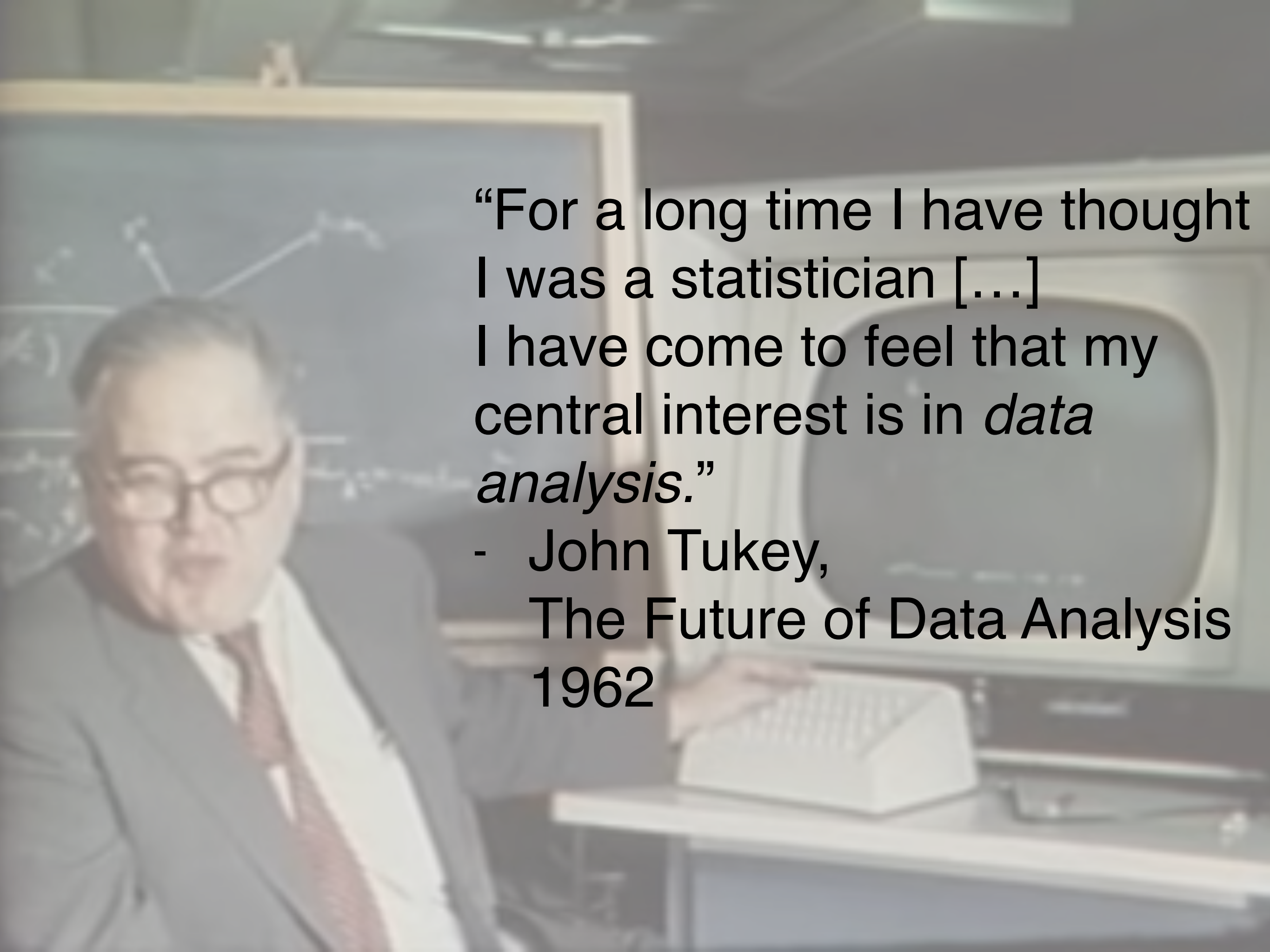Visiting Assistant Professor of Statistical and Data Sciences
Smith College Northampton, MA, USA

flickr: @N01

The focus of my research is to make it easier for people to do data analysis and understand the results of data analysis.

The focus of my research is to make it easier for people to do **data analysis** and understand the results of data analysis.

"For a long time I have thought I was a statistician […]
I have come to feel that my central interest is in *data analysis.*"
- John Tukey,
  The Future of Data Analysis
  1962

# CONTENTS

## I INTRODUCTORY

## II DIAGRAMS

## III DISTRIBUTIONS

## IV TESTS OF GOODNESS OF FIT, INDEPENDENCE AND HOMOGENEITY; WITH TABLE OF $\chi^2$

## V TESTS OF SIGNIFICANCE OF MEANS, DIFFERENCES OF MEANS, AND REGRESSION COEFFICIENTS

## VI THE CORRELATION COEFFICIENT

## VII INTRACLASS CORRELATIONS AND THE ANALYSIS OF VARIANCE

Images via Mark Hansen

# Contents

Images via Mark Hansen

The focus of my research is to make it easier for people to do data analysis and understand the results of data analysis.

The focus of my research is to make it easier for **people** to do data analysis and understand the results of data analysis.

# In particular

High school teachers/students

College students

The focus of my research is to make it easier for people to do data analysis and understand the results of data analysis.

The focus of my research is to make it easier for people to **do** data analysis and **understand** the results of data analysis.

tools

we can push the
boundary with
**tools**
and/or
**curriculum**

curriculum
(and, tools)

"We will be remiss in our duty to our students if we do not see that they learn to use the computer more easily, flexibly, and thoroughly than we ever have; we will be remiss in our duties to ourselves if we do not try to improve and broaden our own uses."

- John Tukey,
The Technical Tools of Statistics
talking about the class of 1970

a modern statistical computing tool

- Accessibility
- Easy entry for novice users
- Data as a first-order persistent object
- Support for a cycle of exploratory and confirmatory analysis
- Flexible plot creation
- Support for randomization throughout
- Interactivity at every level
- Inherent visual documentation
- Simple support for narrative, publishing, and reproducibility
- Flexibility to build extensions

McNamara, Amelia. "Key Attributes of a Modern Statistical Computing Tool." Preprint, http://bit.ly/ModernStatComputing.

tools designed for **learning** statistics are typically:

- graphical
- interactive
- intuitive
- supportive of EDA

but:

- don't support reproducibility
- can't handle real data

tools designed for **doing** statistics are typically:

- powerful
- flexible
- reproducible
- supportive of extensions

but:

- hard to get started using
- not interactive

# Easy entry— TinkerPlots

# Extensible—R

# Extensible—R

# Extensible—R

# Interactivity— Fathom

# Interactivity— Fathom

# R/ggplot2

# R/ggplot2

# Tableau

# Tableau

## Gather your data

A histogram is based on a collection of data about a numeric variable. Our first step is to gather some values for that variable. The initial dataset we will consider consists of fuel consumption (in miles per gallon) from a sample of car models available in 1974 (yes, rather out of date). We can visualize the dataset as a pool of items, with each item identified by its value—which in theory lets us "see" all the items, but makes it hard to get the gestalt of the variable. What are some common values? Is there a lot of variation?

## Sort into an ordered list

A useful first step towards describing the variable's distribution is to sort the items into a list. Now we can see the maximum value and the minimum value. Beyond that, it is hard to say much about the center, shape, and spread of the distribution. Part of the problem is that the list is completely filled; the space between any two items is the same, no matter how dissimilar their values may be. We need a way to see how the items relate to each other. Are they clustered around a few specific values? Is there one lonely item, with a value far removed from all the others?

## Draw the number line

A common convention is to use a number line, on which higher values are displayed to the right and smaller (or negative) values to the left. We can draw a line representing all possible numbers between the minimum and maximum data values.

## Add data to the number line

Now, we map each item to a dot at the appropriate point along the number line. In our visualization we draw the path followed by each item on its way from the list to the line, helping to reveal how adjacent list items end up close or far apart on the number line

gather data items
sort items into list
draw a number line
place items on number line
...(keep scrolling)

(vis scale: 100%)

unit: seconds

dataset: Geyser—272 records of delay (in seconds) between eruptions of Old Faithful

http://tinlizzie.org/histograms/

## Gather your data

A histogram is based on a collection of data about a numeric variable. Our first step is to gather some values for that variable. The initial dataset we will consider consists of fuel consumption (in miles per gallon) from a sample of car models available in 1974 (yes, rather out of date). We can visualize the dataset as a pool of items, with each item identified by its value—which in theory lets us "see" all the items, but makes it hard to get the gestalt of the variable. What are some common values? Is there a lot of variation?

gather data items
sort items into list ⊳
draw a number line ⊳
place items on number line ⊳
...(keep scrolling)

(vis scale: 100%)

unit: seconds

dataset: Geyser—272 records of delay (in seconds) between eruptions of Old Faithful

## Sort into an ordered list

A useful first step towards describing the variable's distribution is to sort the items into a list. Now we can see the maximum value and the minimum value. Beyond that, it is hard to say much about the center, shape, and spread of the distribution. Part of the problem is that the list is completely filled; the space between any two items is the same, no matter how dissimilar their values may be. We need a way to see how the items relate to each other. Are they clustered around a few specific values? Is there one lonely item, with a value far removed from all the others?

## Draw the number line

A common convention is to use a number line, on which higher values are displayed to the right and smaller (or negative) values to the left. We can draw a line representing all possible numbers between the minimum and maximum data values.

## Add data to the number line

Now, we map each item to a dot at the appropriate point along the number line. In our visualization we draw the path followed by each item on its way from the list to the line, helping to reveal how adjacent list items end up close or far apart on the number line

http://tinlizzie.org/histograms/

| View of Data | Perceptual Unit | Data Structure | Student Observation |
|---|---|---|---|
| Pointer | ? | | We said our favorite colors |
| Case Value | ● | | Juan likes red |
| Classifier | Red | | Three like red |
| Aggregate | Red    Not Red | | Half like red |

Konold, C. et al. "Data seen through different lenses." *Educational Studies in Mathematics*, 2014.

# Spatial aggregation toy

# Spatial aggregation toy

# Interactivity in published work

**Choose a Ranking** (choose a weighting or make your own)

| IEEE Spectrum | Trending | Jobs | Open | Custom |
|---|---|---|---|---|

[Edit Ranking] [Add a Comparison]

**Language Types** (click to hide)

| Web | Mobile | Enterprise | Embedded |
|---|---|---|---|

| Language Rank | Types | Spectrum Ranking |
|---|---|---|
| 1. Java | Web Mobile Enterprise | 100.0 |
| 2. C | Mobile Enterprise Embedded | 99.2 |
| 3. C++ | Mobile Enterprise Embedded | 95.5 |
| 4. Python | Web Enterprise | 93.4 |
| 5. C# | Web Mobile Enterprise | 92.2 |
| 6. PHP | Web | 84.6 |
| 7. Javascript | Web Mobile | 84.3 |
| 8. Ruby | Web | 78.6 |
| 9. R | Enterprise | 74.0 |
| 10. MATLAB | Enterprise | 72.6 |

[Show Extended Ranking]

http://bit.ly/IEEE_languagerank

# Interactivity in published work

**Choose a Ranking** (choose a weighting or make your own)

| IEEE Spectrum | Trending | Jobs | Open | Custom |

[Edit Ranking] [Add a Comparison]

**Language Types** (click to hide)

| 🌐 Web | 📱 Mobile | 🖥 Enterprise | ▦ Embedded |

The ranking is calculated using 12 weighted data sources. Click a data source to toggle its inclusion in the ranking and drag its slider to reweight it.

| Google (search) | ——O—— 50 | Google (trends) | ———O 50 |
| Github (active) | ——O—— 50 | Github (created) | —O——— 30 |
| Stack Overflow (?s) | —O——— 30 | Stack Overflow (views) | —O——— 30 |
| Reddit | —O——— 20 | Hacker News | —O——— 20 |
| Career Builder | O———— 5 | Dice | O———— 5 |
| Topsy | —O——— 20 | IEEE Xplore | ———O 100 |

[Cancel] [Save as Custom]

http://bit.ly/IEEE_languagerank

# CODAP

## (Common Online Data Analysis Platform)



http://concord.org/projects/codap

# Curriculum

www.mobilizingcs.org/

**IDS**

**Introduction to Data Science**

Robert Gould

Suyen Moncada-Machado

Terri Anna Johnson

James Molyneux

- Year-long course

- Validates Algebra II requirement

- "Data science"

- Taught in R within RStudio server

- Participatory sensing

- Content includes:

  - Exploratory data analysis

  - Randomization, simulation, bootstrapping

  - Simple linear regression, multiple regression

  - Decision trees, clustering, k-means

www.mobilizingcs.org/

# IDS

## Introduction to Data Science

**Robert Gould**

**Suyen Moncada-Machado**

**Terri Anna Johnson**

**James Molyneux**

---

## RStudio Lab Codes and Functions

### Contents

### Loading, saving and viewing data

**data()**: Loads and displays a pre-loaded data file from RStudio.

*Example*:

```
data(cdc)
```

---

**read.csv()**: Imports data from a *.csv* formatted file in R.

*Example*:

```
read.csv("Time Use.csv")
```

---

**View()**: Displays the data as a spreadsheet in a new tab.

*Example*:

```
View()
```

---

**head()**: Prints the first 6 values or rows of data in the console.

*Examples*:

```
# Observations of a dataset
head(cdc)

# Observations of a variable
head(-gender, data = cdc)
```

---

**tail()**: Prints the last 6 values or rows of data in the console.

*Examples*:

```
# Observations of a dataset
tail(cdc)
```

# R Syntax

```
xyplot(mpg ~ wt | as.factor(cyl), data = mtcars)
```

vs.

```
par(mfrow = c(1,3))
plot(mtcars$wt[mtcars$cyl == 4], mtcars$mpg[mtcars$cyl == 4])
plot(mtcars$wt[mtcars$cyl == 6], mtcars$mpg[mtcars$cyl == 6])
plot(mtcars$wt[mtcars$cyl == 8], mtcars$mpg[mtcars$cyl == 8])
```

# R Syntax

# Formula-based syntax

- `lattice` graphics

- `mosaic` package for statistics

- `mobilizr` additional functions

# Want to check out the labs?

```r
library(devtools)
install_github("mobilizingcs/mobilizr")
library(mobilizr)
load_labs()
```

IDS Teachers are engaging students in exploratory data analysis and statistical thinking to deepen their understanding of STEM concepts.

http://stemforall2016.videohall.com/presentations/790

# Data science at Smith



flickr: leslee

# Introductory course

Introductory Statistics with
Randomization and
Simulation.
David M Diez,  Christopher
D Barr, Mine Çetinkaya-
Rundel
www.openintro.org



Introductory
Statistics with
Randomization
and Simulation

**First Edition**

**OpenIntro**

David M Diez
Christopher D Barr
Mine Çetinkaya-Rundel

# R syntax comparison
## Cheat Sheet

**SMITH COLLEGE**

## Syntax

**Syntax** is the set of rules that govern what code works and doesn't work. Most programming languages offer one standardized syntax, but R has many.

**Most people use some combination of all the syntaxes available to them.**

1. Dollar sign syntax uses the dollar sign to locate a variable within a dataset. It is expected by most **base** R functions.
2. Formula syntax uses the **data=** argument at the end of a list of function arguments. The formula syntax is used by modeling functions like **lm()**, **lattice** graphics like **xyplot()**, and **mosaic** summary statistics like **mean()**.
3. Tidyverse syntax uses data as the first argument to function calls. It is used by the packages **dplyr** and **tidyr**, among others. The associated graphics library is **ggplot2**.

## Dollar sign syntax

`goal(data$x, data$y)`

**Summary statistics:**
one continuous variable:
mean(mtcars$mpg)

one categorical variable:
table(mtcars$cyl)

two categorical variables:
table(mtcars$cyl, mtcars$am)

one continuous, one categorical:
mean(mtcars$mpg[mtcars$cyl==4])
mean(mtcars$mpg[mtcars$cyl==6])
mean(mtcars$mpg[mtcars$cyl==8])

**Plotting:**
one continuous variable:
hist(mtcars$disp)

boxplot(mtcars$disp)

one categorical variable:
barplot(table(mtcars$cyl))

two continuous variables:
plot(mtcars$disp, mtcars$mpg)

two categorical variables:
mosaicplot(table(mtcars$am, mtcars$cyl))

one continuous, one categorical:
histogram(mtcars$disp[mtcars$cyl==4])
histogram(mtcars$disp[mtcars$cyl==6])
histogram(mtcars$disp[mtcars$cyl==8])

boxplot(mtcars$disp[mtcars$cyl==4])
boxplot(mtcars$disp[mtcars$cyl==6])
boxplot(mtcars$disp[mtcars$cyl==8])

**Wrangling:**
subsetting:
mtcars[mtcars$mpg>30, ]

making a new variable:
mtcars$efficient[mtcars$mpg>30] <- TRUE
mtcars$efficient[mtcars$mpg<30] <- FALSE

## Formula syntax

`goal(y~x|z, data=data, group=w)`

**Summary statistics:**
one continuous variable:
mosaic:: mean(~mpg, data=mtcars)

one categorical variable:
mosaic:: tally(~cyl, data=mtcars)

two categorical variables:
mosaic:: tally(cyl~am, data=mtcars)

one continuous, one categorical:
mosaic:: mean(mpg~cyl, data=mtcars)

*tilde*

**Plotting:**
one continuous variable:
lattice::histogram(~disp, data=mtcars)

lattice::bwplot(~disp, data=mtcars)

one categorical variable:
mosaic::bargraph(~cyl, data=mtcars)

two continuous variables:
lattice::xyplot(mpg~disp, data=mtcars)

two categorical variables:
mosaic::bargraph(~am, data=mtcars, group=cyl)

one continuous, one categorical:
lattice::histogram(~disp|cyl, data=mtcars)

lattice::bwplot(cyl~disp, data=mtcars)

## Tidyverse syntax

`data %>% goal(x)`

**Summary statistics:**
one continuous variable:
mtcars %>% dplyr:: summarize(mean(mpg))

one categorical variable:
mtcars %>% dplyr:: group_by(cyl) %>%
dplyr:: summarize(n())

*the pipe*

two categorical variables:
mtcars %>% dplyr:: group_by(cyl, am) %>%
dplyr:: summarize(n())

one continuous, one categorical:
mtcars %>% dplyr::group_by(cyl) %>%
dplyr:: summarize(mean(mpg))

**Plotting:**
one continuous variable:
ggplot2::qplot(x=mpg, data=mtcars, geom = "histogram")

ggplot2::qplot(y=disp, x=1, data=mtcars, geom="boxplot")

one categorical variable:
ggplot2::qplot(x=cyl, data=mtcars, geom="bar")

two continuous variables:
ggplot2::qplot(x=disp, y=mpg, data=mtcars, geom="point")

two categorical variables:
ggplot2::qplot(x=factor(cyl), data=mtcars, geom="bar") + facet_grid(.~am)

one continuous, one categorical:
ggplot2::qplot(y=disp, x=factor(cyl), data=mtcars, geom="boxplot")

ggplot2::qplot(x=disp, data=mtcars, geom = "histogram") + facet_grid(.~cyl)

**Wrangling:**
subsetting:
mtcars %>%
dplyr:: filter(mpg>30)

making a new variable:
mtcars <- mtcars %>%
dplyr:: mutate(efficient = if_else(mpg>30, TRUE, FALSE))

- Graphical literacy skills should be taught in the context of science and social science.
- Translating between representations may be beneficial
- Explicitly focus on the links between visual features and meaning.
- Make graph reading metacognitive

Shah, P. and Hoeffner, J. "Review of Graph Comprehension Research: Implications for Instruction." *Educational Psychology Review*, 2002.

New faculty



- SDS 136: Communicating with Data
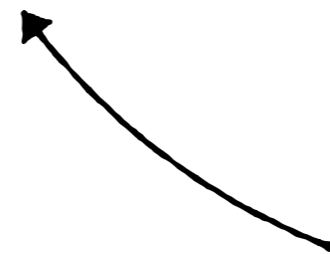
- SDS 192: Introduction to Data Science

- SDS 235: Visual Analytics

- SDS 236: Data Journalism

- SDS 293: Machine Learning

New courses

- Started at UCLA
- Now a national and international event
- Groups of undergraduate students compete to find insight in data
- Past data sponsors: LAPD, Kiva.com, eHarmony, Travelocity

# Data Science as a Science

*Figure 3.* Prestructural representation from Group 5. (The caption at the top reads "This is a graph of a group of kids who showed their ages".)

Chick, H. and Watson, J. "Data representation and interpretation by primary school students working in groups." *Mathematics Education Research Journal,* 2001.

# Thank you

Amelia McNamara (@AmeliaMN)