# Big Data, Machine Learning & the Social Sciences

## Hanna Wallach

Microsoft Research & University of Massachusetts Amherst

hanna@dirichlet.net

# DATA

# QUESTIONS

# MODELS

# FINDINGS

# DATA

QUESTIONS

MODELS

FINDINGS

# What is Big Data?

> "[B]ig data […] refers to large, diverse, complex, longitudinal, and/or distributed data sets generated from instruments, sensors, Internet transactions, email, video, click streams, and/or all other digital sources available.

— NSF & NIH

# What is Big Data?

> " Big data is high volume, high velocity, and/or high variety information assets that require new forms of processing to enable enhanced decision making [and] insight discovery.

— Gartner, Inc., 2012

# Social Data

> " Big data is the amassing of huge amounts of statistical information on **social and economic trends and human behavior**.
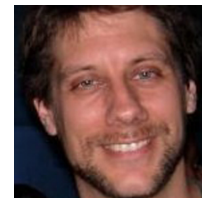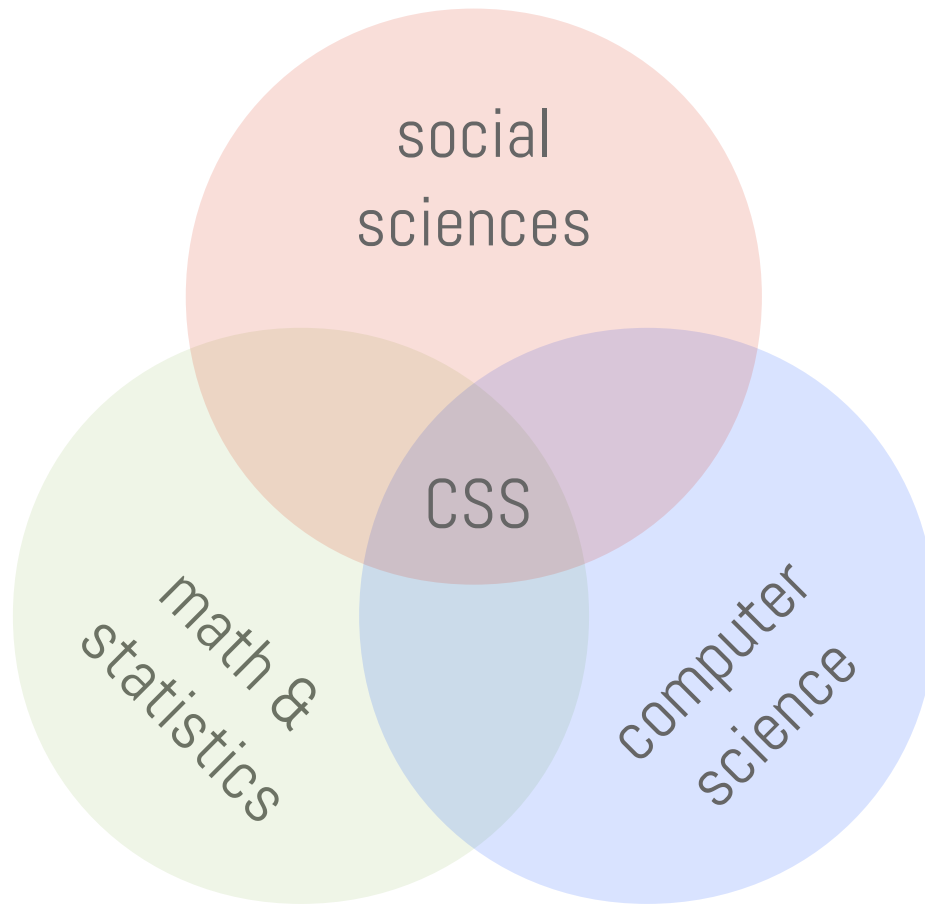
— Michelle Chen, 2014

# Granular Data

> " The issue is not just size – we've always had big data sets; the issue is **granularity**.

— Michael Jordan, 2014

# Computational Social Science
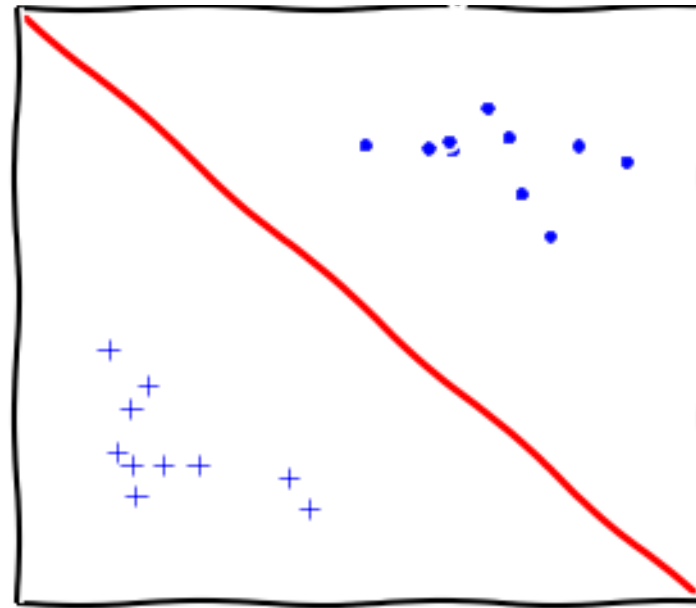
# On Average... Good?



majority          minority

DATA

# DATA

DATA information data
DaTa DatA DAta daTA
INFORMATION dAtA

data → data → data →
data → data → data → ...

*Data*

# DATA

QUESTIONS

MODELS

FINDINGS

DATA

QUESTIONS

MODELS

FINDINGS

# Small Patterns

> " The real, epistemological problem with big data is small patterns. [... But] small patterns may be significant only if properly aggregated. [...]
> So what we need is a better understanding of which data are worth preserving.

— Luciano Floridi, 2013

# Big Questions

" **And this is a matter of grasping which questions are or will be interesting. [… U]ltimately, the game will be won by those who 'know how to ask and answer relevant questions.'**

— Luciano Floridi, 2013

# Data-first?

# Method-first?

$$= \int d\phi \prod_{v=1}^{V} \phi_v^{N_v} \frac{\Gamma(\beta)}{\prod_{v=1}^{V} \Gamma(\beta n_v)} \prod_{v=1}^{V} \phi_v^{\beta n_v - 1} \delta\left(\sum_{v=1}^{V} \phi_v - 1\right)$$

$$= \frac{\Gamma(\beta)}{\prod_{v=1}^{V} \Gamma(\beta n_v)} \int d\phi \prod_{v=1}^{V} \phi_v^{N_v + \beta n_v - 1} \delta\left(\sum_{v=1}^{V} \phi_v - 1\right)$$

$$= \frac{\Gamma(\beta)}{\prod_{v=1}^{V} \Gamma(\beta n_v)} \frac{\prod_{v=1}^{V} \Gamma(N_v + \beta n_v)}{\Gamma(N + \beta)},$$

```
    $line .= <CASEBOOKS>;
    redo unless eof(CASEBOOKS);
}

$line =~ s/\\\t/xyzdrptmpxyz/g;
@columns = split("\t", $line);
$columns[3] = uc $columns[3];
$line = join("\t", @columns);
$line =~ s/xyzdrptmpxyz/\\\t/g;
```

~~Data-first~~

~~Method-first~~

# QUESTION-FIRST

# Gender in Local Government

**?** Is there gender homophily? Do women occupy disadvantaged positions in the network? Are women less central in the dominant coalition?

# "Push" Transparency

# "Pull" Transparency

## Mayor and City Council Contact Information

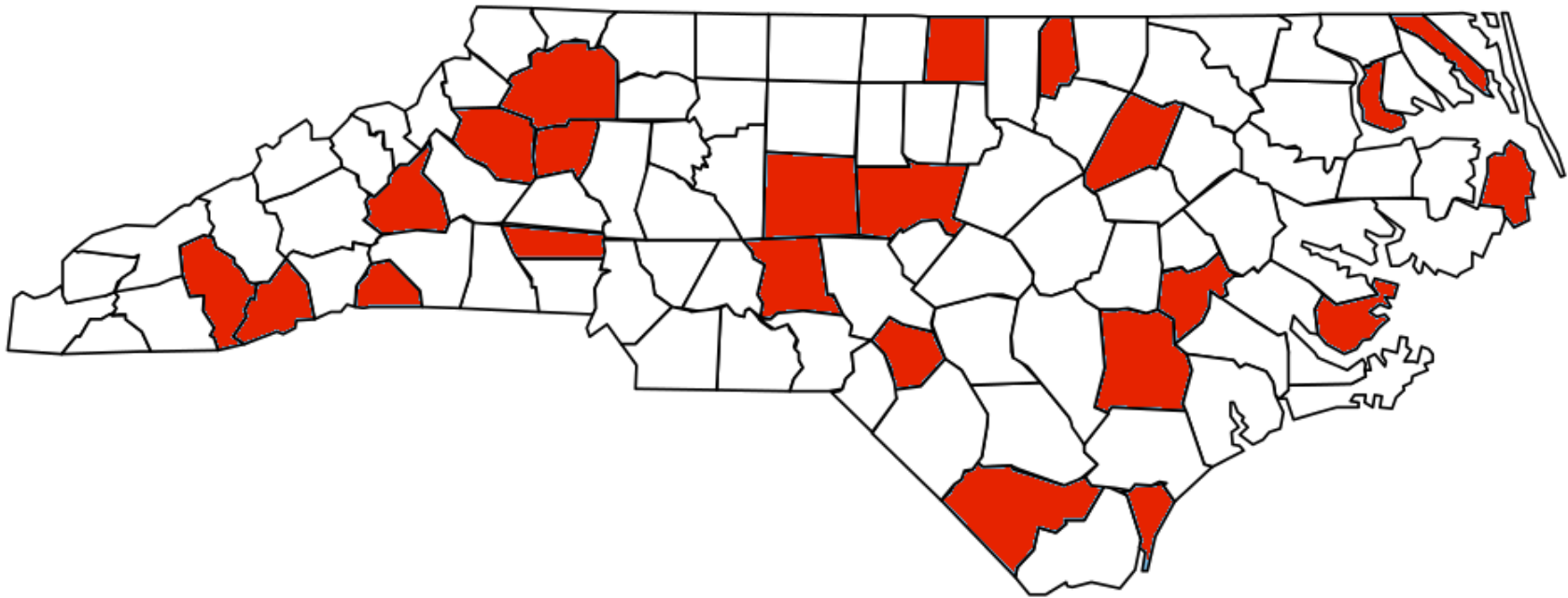Agenda, Minutes and Meeting Info | Contact the Council | Council Mission & Goals

You can contact the entire council at once via Email , or e-mail each council member individually via their addresses below.

Please note that any correspondence, such as e-mail or letters, sent to City officials or staff becomes public record and is available for public/media review.
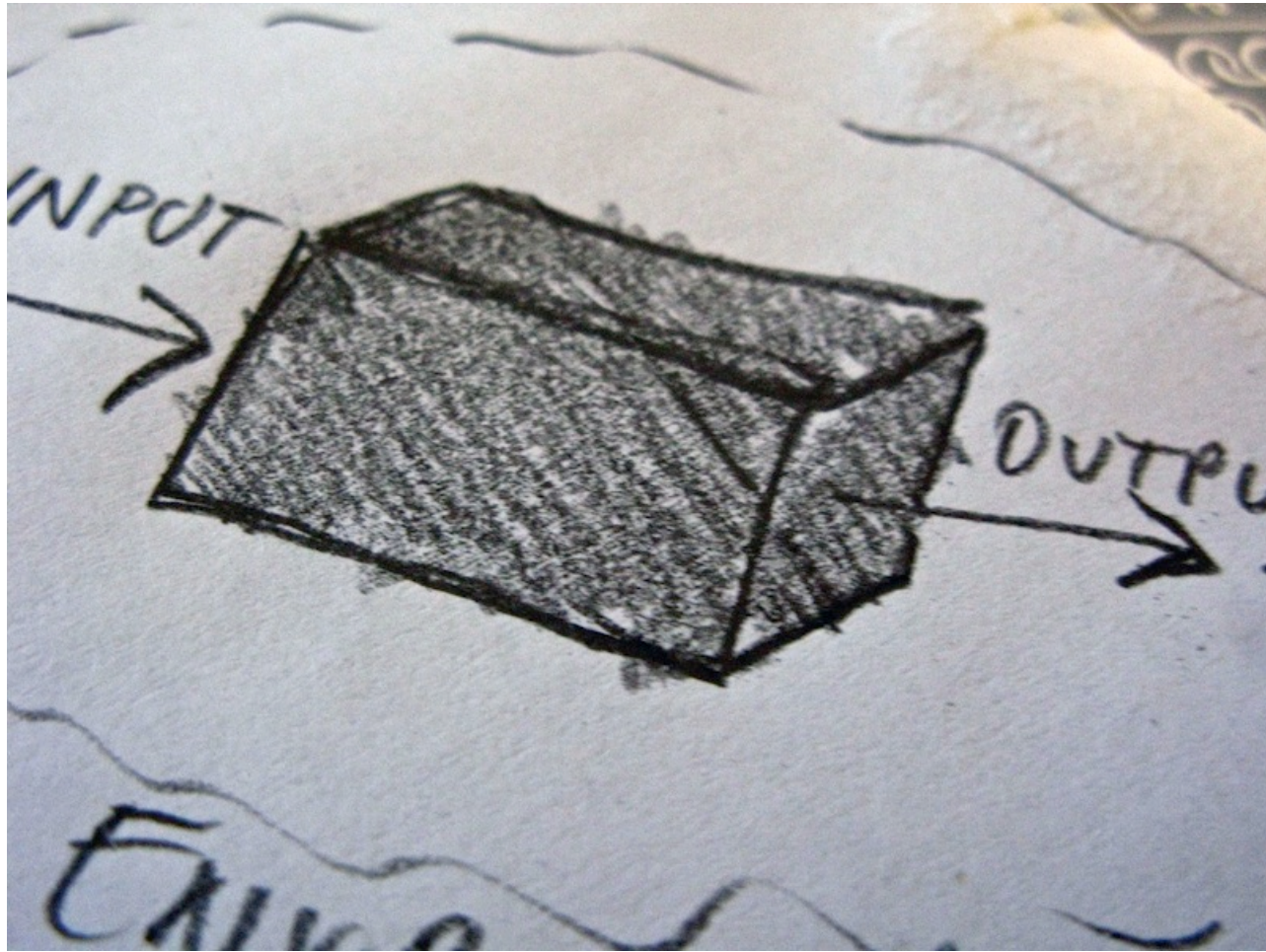
public records:
new data availability

# 23 Counties, 500,000+ Emails

# Algorithmic Accountability

DATA

QUESTIONS

MODELS
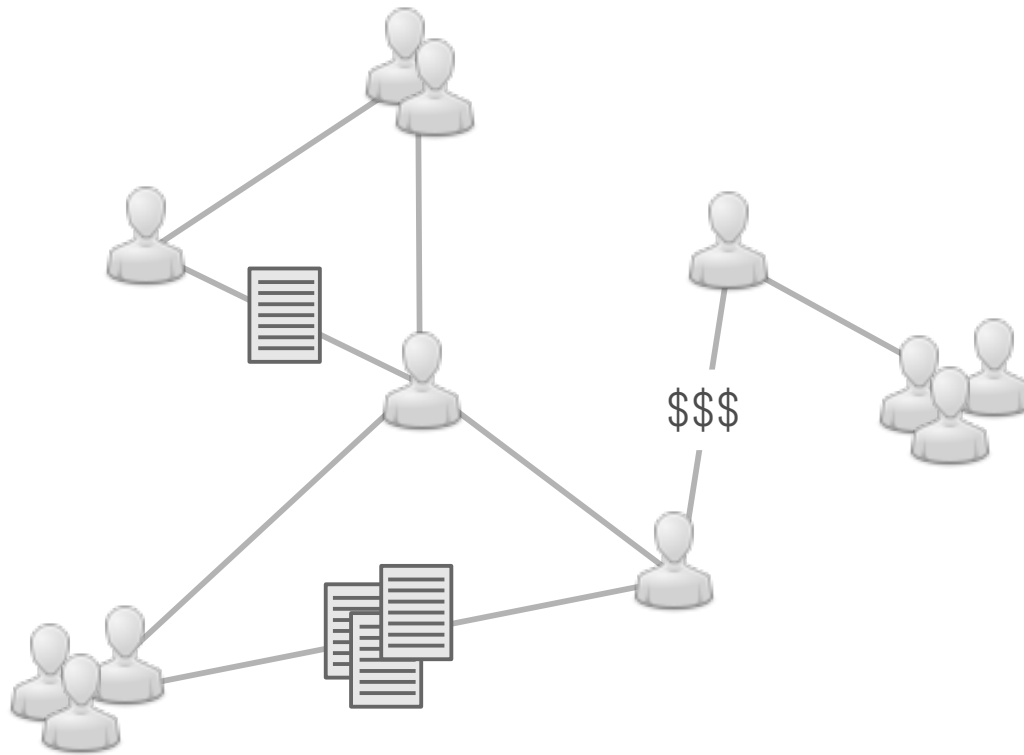
FINDINGS

DATA

QUESTIONS

MODELS

FINDINGS

# Research Goals

"

[C]omputer scientists may be interested in finding the needle in the haystack—such as […] the right web page to display from a search— but social scientists are more commonly interested in characterizing the haystack.
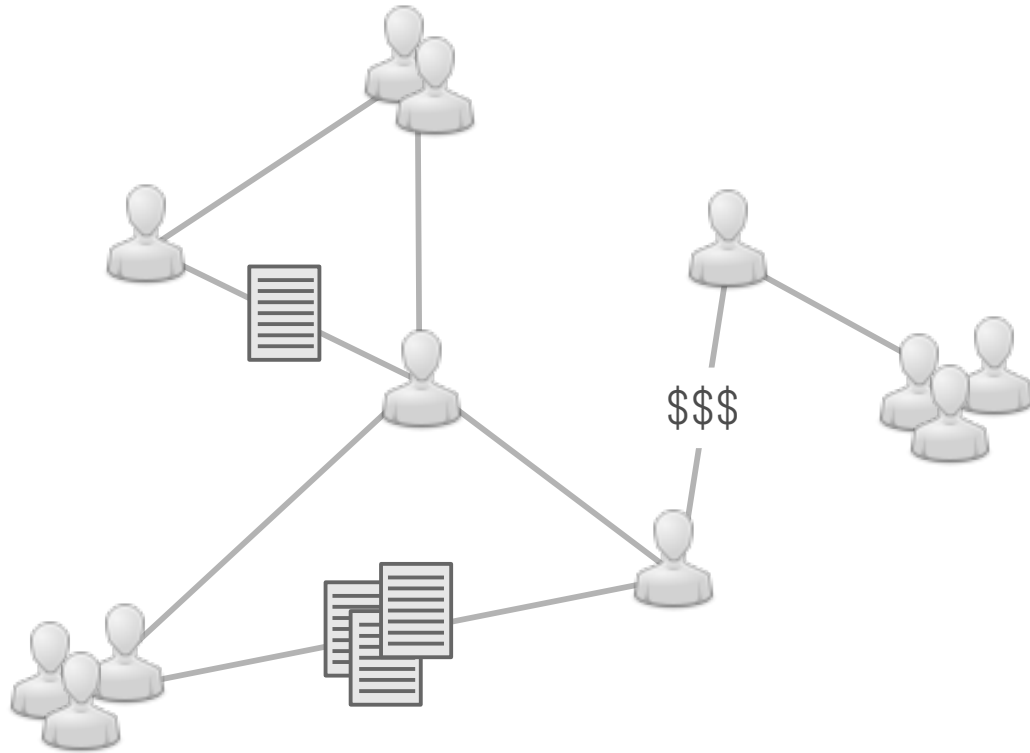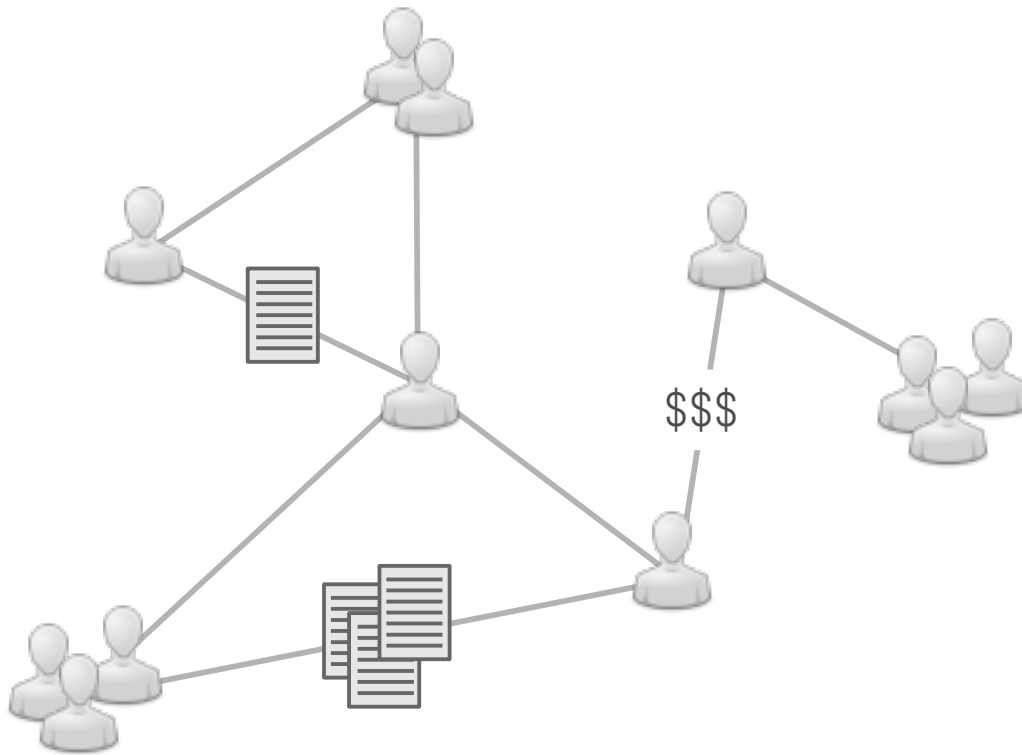
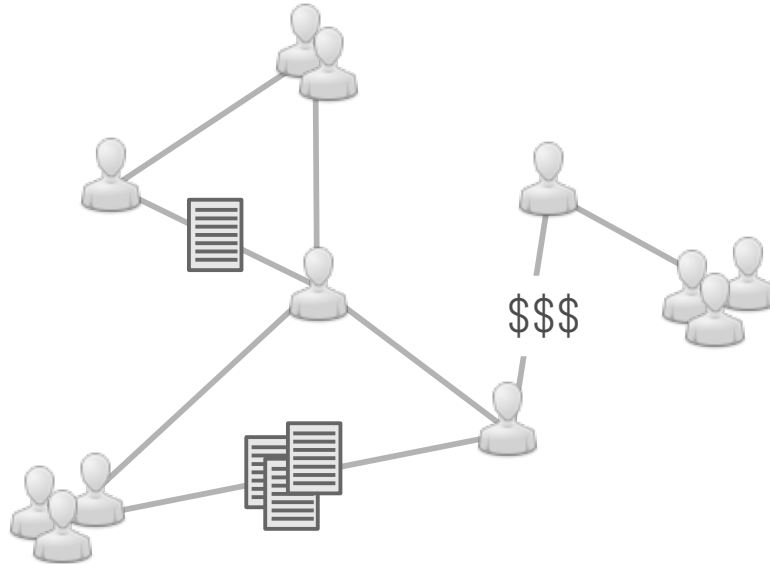— King & Hopkins, 2010
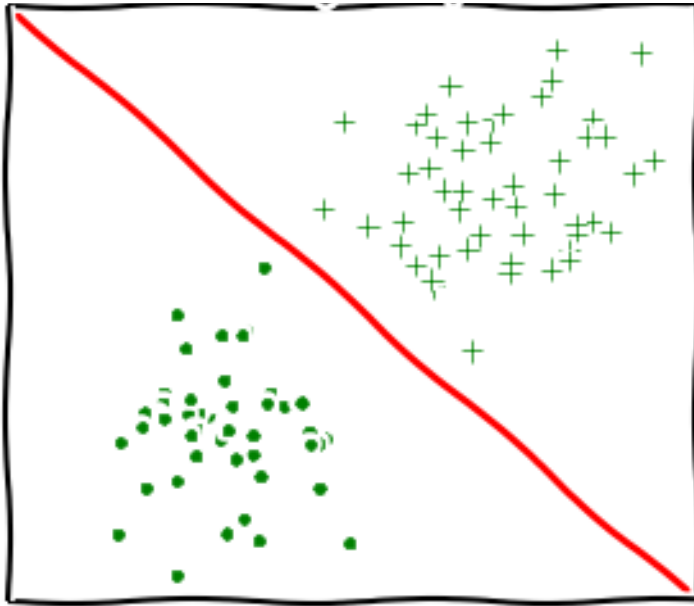
# Prediction

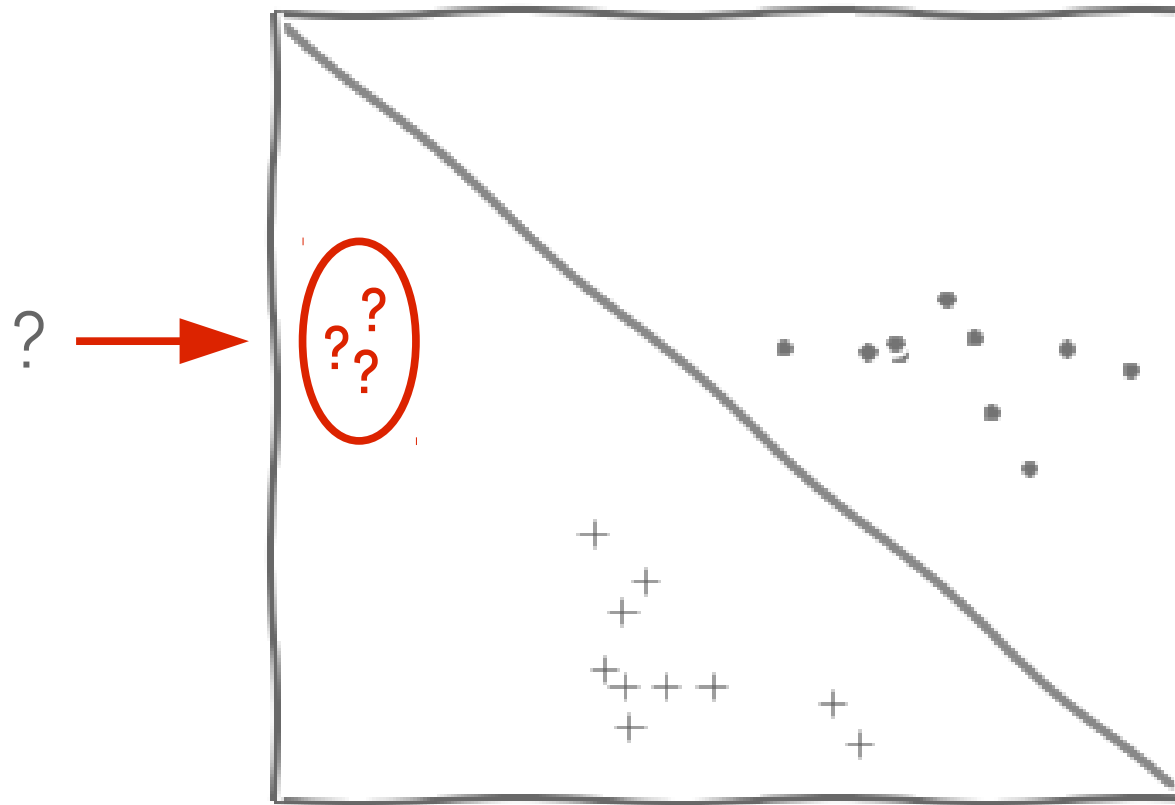# Explanation

???

# Exploration

# Bias & Unfairness

# Error Analysis
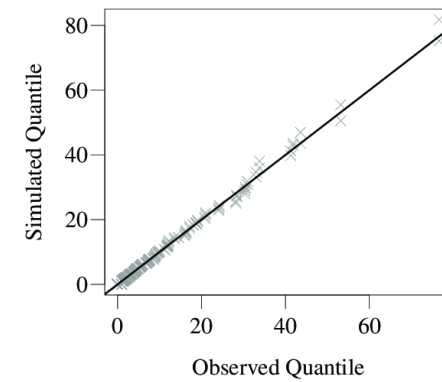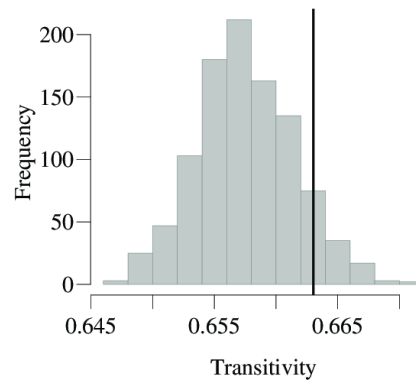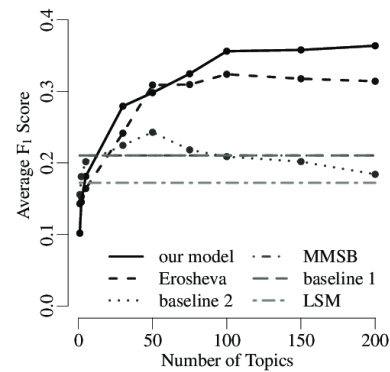
# Uncertainty

# Exploratory Models

DATA

QUESTIONS

# MODELS
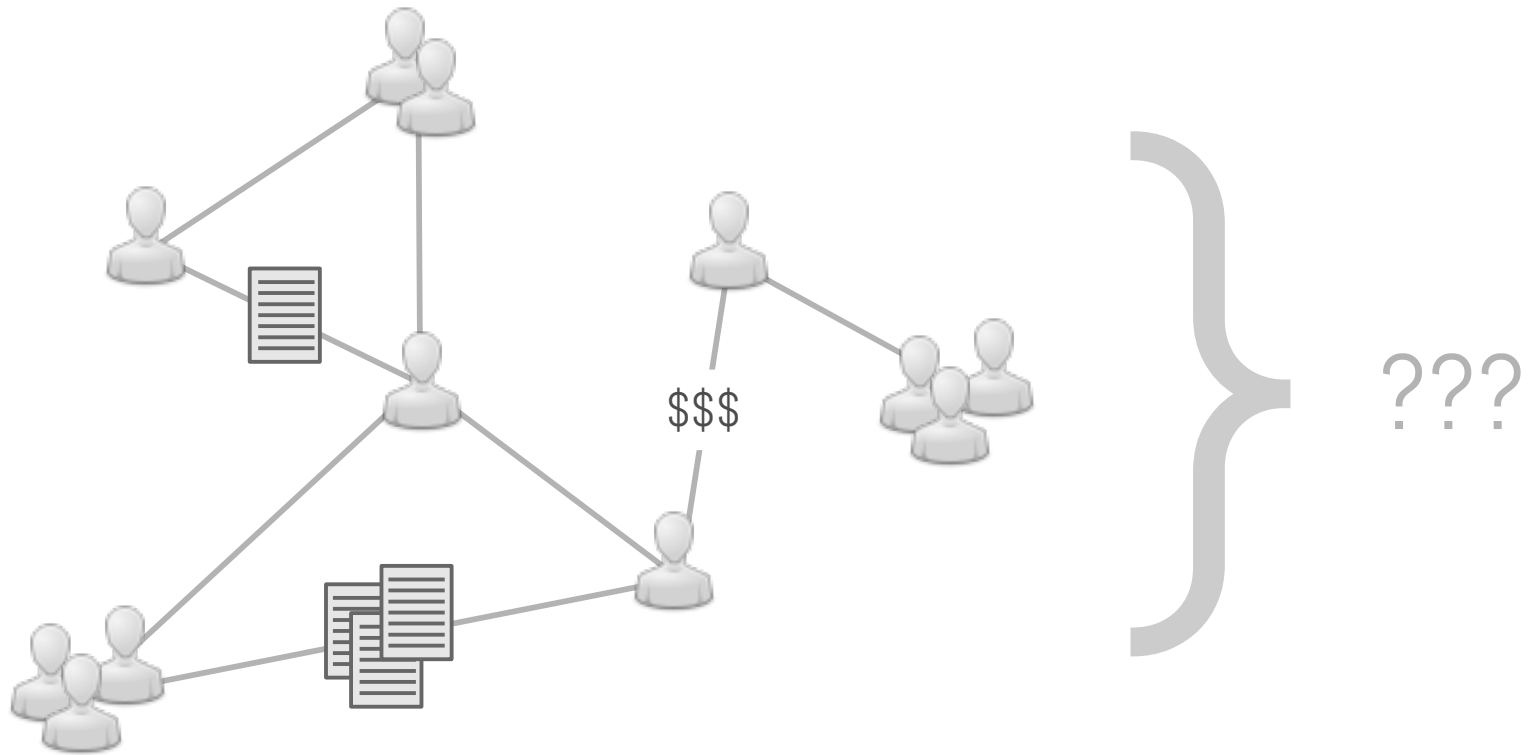
FINDINGS
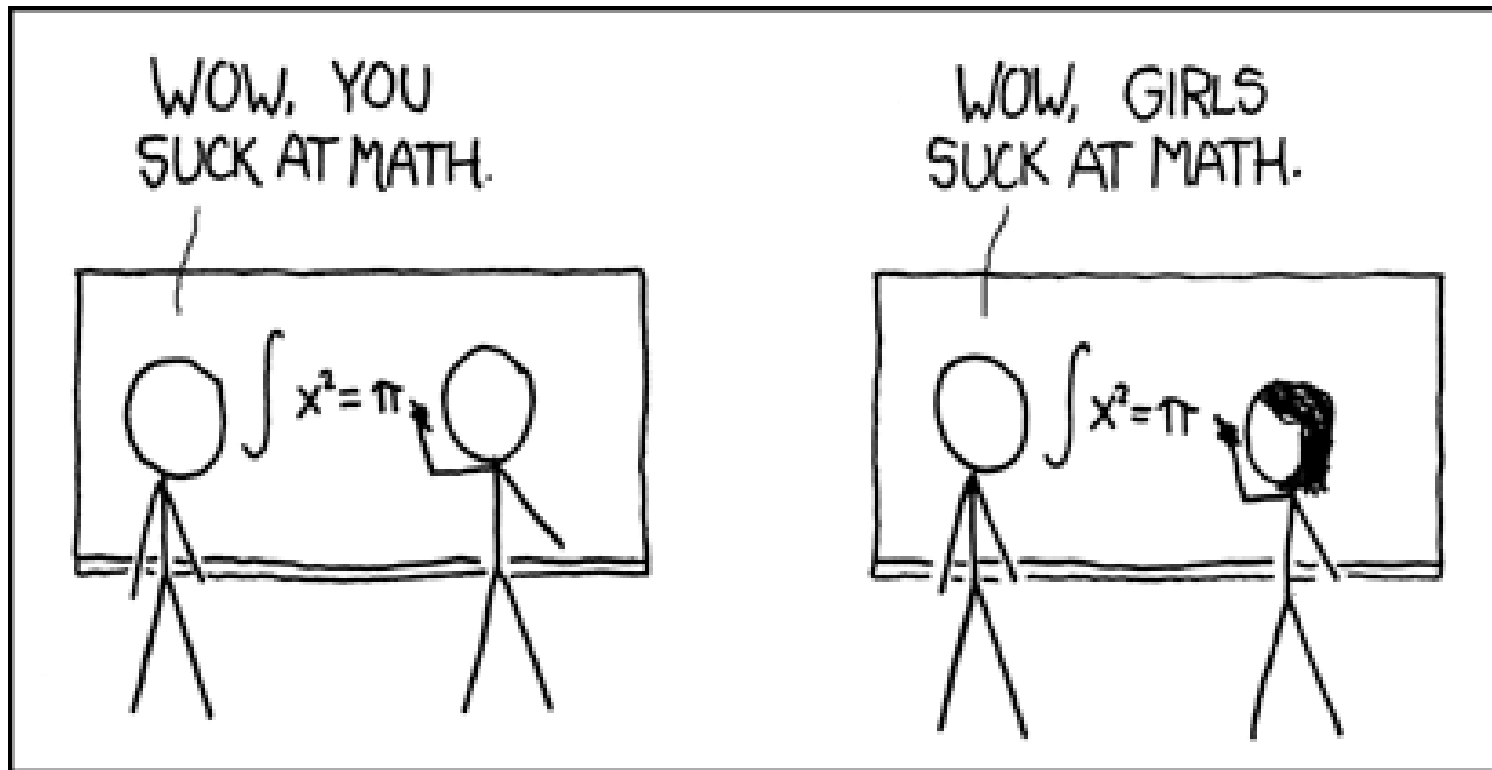
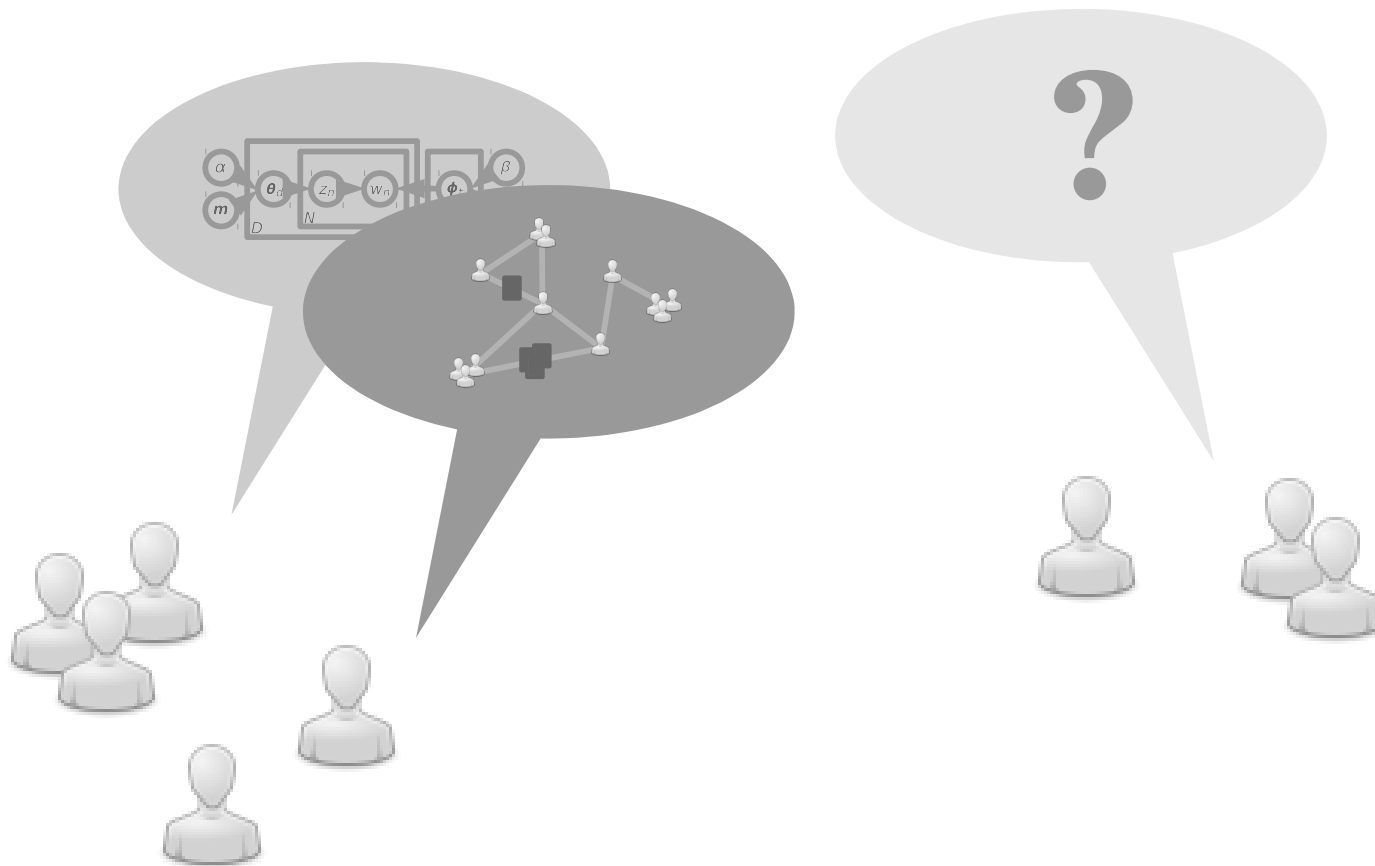DATA

QUESTIONS

MODELS

FINDINGS

# Intuition

# Confirmation Bias

# e.g., Stereotype Threat

# Scientific Communication

DATA

QUESTIONS

MODELS

# FINDINGS

# Thanks!

@hannawallach | http://dirichlet.net