

# Random projections, reweighting and half-sampling for high-dimensional statistical inference

Art B. Owen

Stanford University

based on joint works with:

Dean Eckles      Facebook Inc.

Sarah Emerson      Oregon State University

# About these slides

These are the slides I presented on February 15 at MCQMC 2012 in Sydney Australia.

I have corrected some typos and extended the presentation of the challenging integral over the Stiefel manifold.

A few of these slides were skipped over in order to allow time for questions.

This talk covers two projects. The bootstrap work with Dean Eckles has now been accepted by the Annals of Applied statistics. The projection work with Sarah Emerson is still in progress.

# Monte Carlo methods for statistics

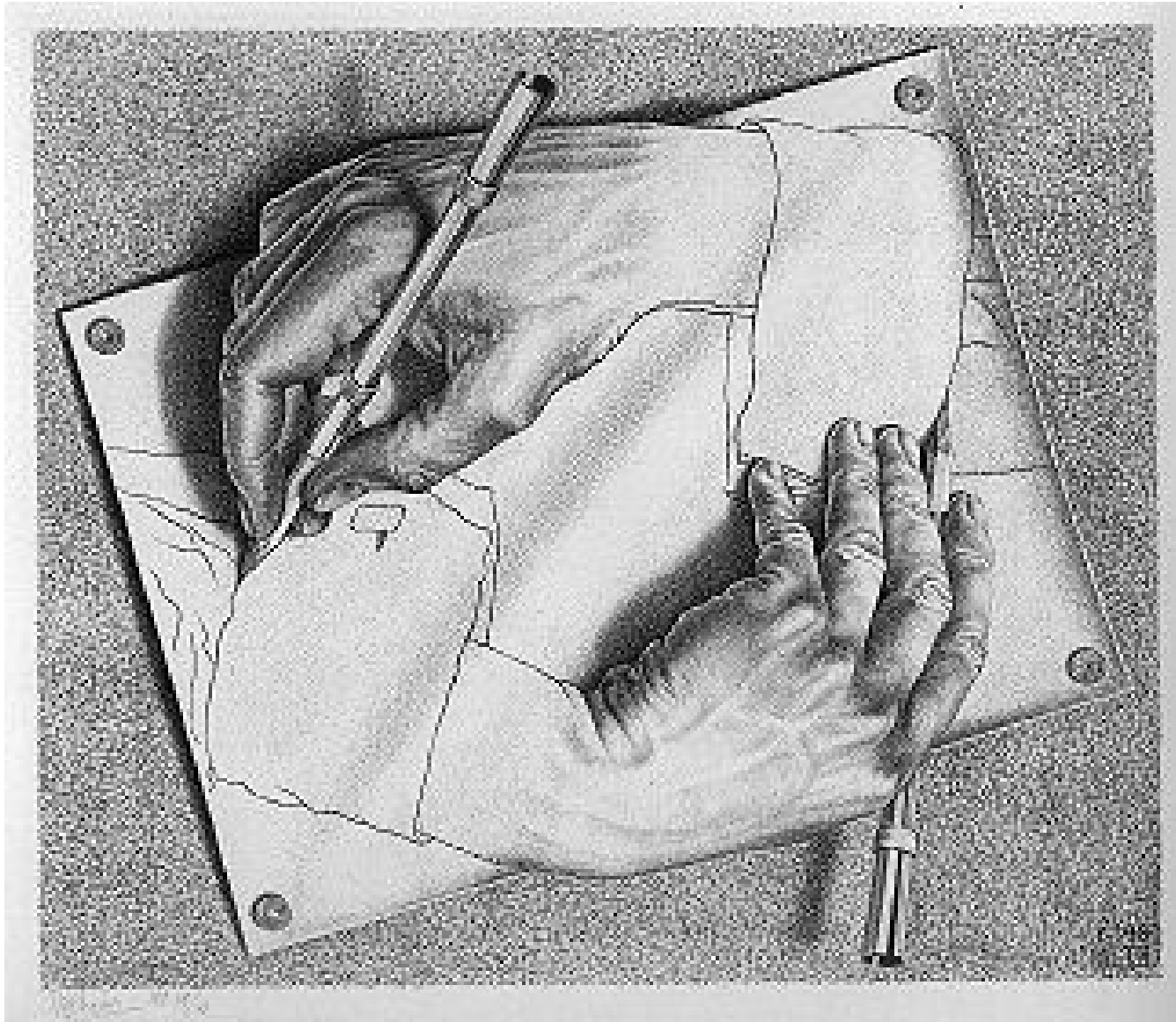
- 1) Markov chain Monte Carlo
- 2) Bootstrap resampling

The above are mainstays. We also use:

- 1) Random permutations
- 2) Random projections
- 3) Sample splitting

Probability/Statistics and Monte Carlo are closely intertwined.

# Statistics and Monte Carlo



M. C. Escher (1948)

MCQMC 2012, February 2012

This talk will show some uses of MC uses in statistics.

# Some statistical notions

$\mathbf{X} \sim F$  random vector  $\mathbf{X}$  has distribution  $F$

$\mathbf{X}_i \stackrel{\text{iid}}{\sim} F$   $\mathbf{X}_i$  are statistically Independent and Identically Distributed (IID) from  $F$

$\mathcal{N}_d(\mu, \Sigma)$  The Gaussian distribution with mean  $\mu \in \mathbb{R}^d$   
and variance covariance matrix  $\Sigma \in \mathbb{R}^{d \times d}$

$\mathbf{X} \sim \mathcal{N}(\mu, \Sigma)$ , means

$$\Pr(\mathbf{X} \in A) = \int_A f(\mathbf{x}) \, d\mathbf{x} \quad \text{where}$$

$$f(\mathbf{x}) = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^\top \Sigma^{-1}(\mathbf{x} - \mu)\right)}{(2\pi)^{d/2} \det(\Sigma)^{1/2}}$$

***p-values***

Observe  $T = t$  and compute  $p = \Pr(T \geq t)$ .

If  $p < 0.01$  then the observed value  $t$  happens 1% or less of the time.

Evidence against the hypothesized distribution of  $T$ .

# Problem one

We have

$$\mathbf{X}_1, \dots, \mathbf{X}_{n_x} \stackrel{\text{iid}}{\sim} F \text{ in } \mathbb{R}^d$$

$$\mathbf{Y}_1, \dots, \mathbf{Y}_{n_y} \stackrel{\text{iid}}{\sim} G \text{ in } \mathbb{R}^d$$

is  $F = G$ ?

We might assume  $F = \mathcal{N}(\mu_1, \Sigma)$  and  $G = \mathcal{N}(\mu_2, \Sigma)$ .

Then we test  $\mu_1 = \mu_2$ . This is an old problem.

Revived interest,  $d \gg n_x + n_y$

DNA microarrays

expression levels of  $d \approx 30,000$  genes

on  $n_x$  healthy and  $n_y$  diseased individuals

$n_x, n_y$  tens or hundreds

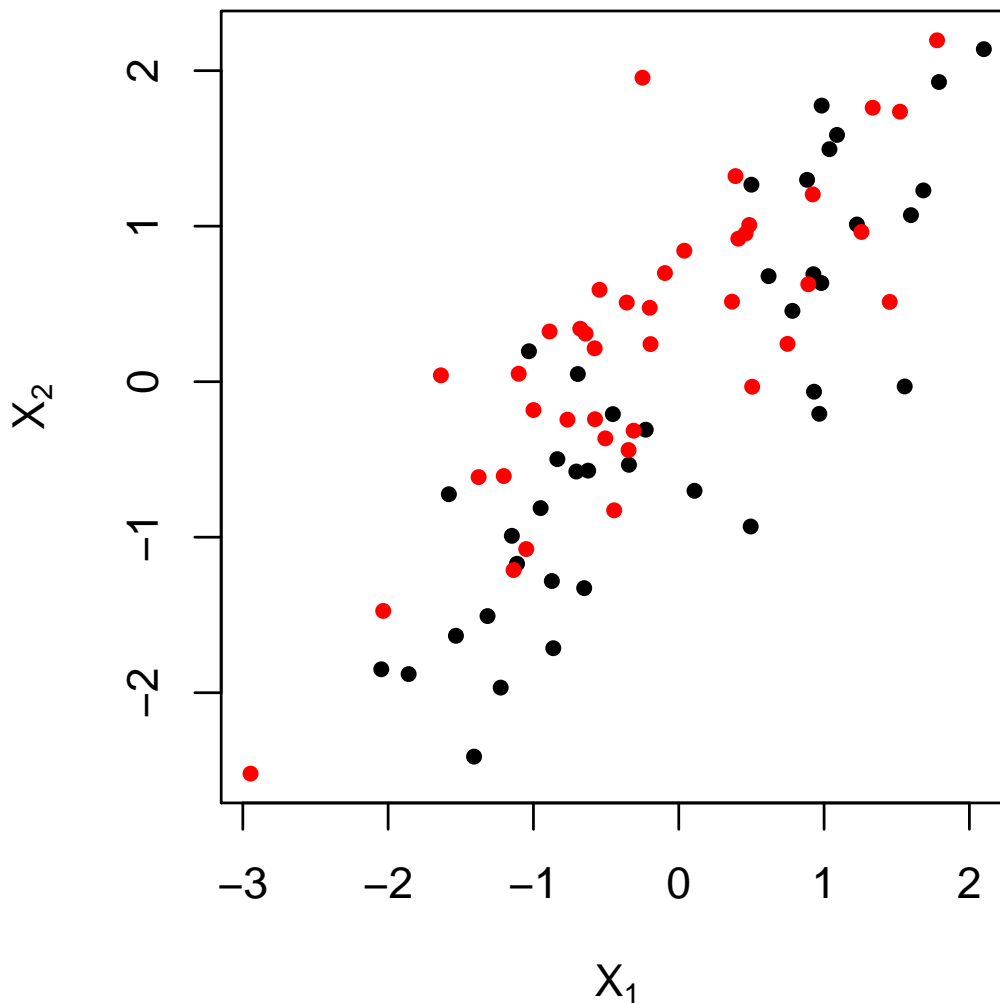
Genome wide association studies

$d \approx 2,000,000$  markers

with  $n_x, n_y$  thousands or more

Also: fMRI, finance

# Illustration



50 red and 50 black points in  $\mathbb{R}^2$

Black points normally distributed

Red points shifted **Northwest** wrt black

$X_1$  not significantly different,  
 $p = 0.47$

$X_2$  not significantly different,  
 $p = 0.09$

$X_1 + X_2$  not significantly different,  
 $p = 0.60$

$X_1 - X_2$  **very** significantly different,  
 $p = 1.7 \times 10^{-4}$

So: how to find the interesting projection?

# Hotelling's $T^2$

Find  $\theta \in \mathbb{R}^d$  with  $\theta^\top \theta = 1$  to maximize the apparent separation between

$$\tilde{X}_i = \theta^\top \mathbf{X}_i \in \mathbb{R} \quad \text{and} \quad \tilde{Y}_i = \theta^\top \mathbf{Y}_i \in \mathbb{R}$$

Answer depends on

$$\begin{aligned} \bar{\mathbf{X}} &= \frac{1}{n_x} \sum_{i=1}^{n_x} \mathbf{X}_i & S_x &= \sum_{i=1}^{n_x} (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^\top \\ \bar{\mathbf{Y}} &= \frac{1}{n_y} \sum_{i=1}^{n_y} \mathbf{Y}_i & S_y &= \sum_{i=1}^{n_y} (\mathbf{Y}_i - \bar{\mathbf{Y}})(\mathbf{Y}_i - \bar{\mathbf{Y}})^\top \end{aligned}$$

Algebraically we get

$$T^2 = \frac{n_x n_y}{n_x + n_y} (\bar{\mathbf{X}} - \bar{\mathbf{Y}})^\top S^{-1} (\bar{\mathbf{X}} - \bar{\mathbf{Y}}) \quad \text{where}$$

$$S = \frac{S_x + S_y}{n_x + n_y - 2} \quad \text{Hotelling (1931)}$$

Get  $T^2 = 18.58$ . Also  $\Pr(T^2 \geq 18.58) = 2.6 \times 10^{-4}$  (p value)



# In high dimensions

When  $d \gg n_x + n_y$  the covariance matrix  $S$  is not invertible. Can't use

$$T^2 = \frac{n_x n_y}{n_x + n_y} (\bar{\mathbf{X}} - \bar{\mathbf{Y}})^\top S^{-1} (\bar{\mathbf{X}} - \bar{\mathbf{Y}})$$

## Geometrically

Some projection  $\theta \in \mathbb{R}^d$  has

$\theta^\top \mathbf{X}_i = \text{constant}$  for  $i = 1, \dots, n_x$  and

$\theta^\top \mathbf{Y}_i = \text{different constant}$ .

IE: we will get perfect separation, even if  $F = G$ .

A classic remedy by [Dempster \(1958\)](#) takes

$$T_{\text{Dempster}}^2 = \frac{n_x n_y}{n_x + n_y} \frac{\|\bar{\mathbf{X}} - \bar{\mathbf{Y}}\|^2}{\text{tr}(S)} = \frac{n_x n_y}{n_x + n_y} \frac{\sum_{j=1}^d (\bar{X}_j - \bar{Y}_j)^2}{\sum_{j=1}^d S_{jj}}$$

but this makes no use of correlations

Recent improvements by: [Bai, Saradanasa, Hall, Fan, Chen, Srivastava](#)

also don't use correlations

# Random projections

Lopes, Jacob, Wainwright (2011)

Choose random  $\Theta \in \mathbb{R}^{d \times k}$  with  $\Theta^\top \Theta = I_k$ .

Put  $\tilde{\mathbf{X}}_i = \Theta^\top \mathbf{X}_i$  and  $\tilde{\mathbf{Y}}_i = \Theta^\top \mathbf{Y}_i$

Then use

$$\tilde{T}_\Theta^2 = \frac{n_x n_y}{n_x + n_y} (\bar{\mathbf{X}} - \bar{\mathbf{Y}})^\top \Theta (\Theta^\top S \Theta)^{-1} \Theta^\top (\bar{\mathbf{X}} - \bar{\mathbf{Y}})$$

Exists if  $k < n_x + n_y - 2$

That is

project data into a random  $k$  dimensional subspace

test means of projected data

this retains some of the correlations

# Uniform random projections

From  $\mathbb{R}^d$  to  $\mathbb{R}$ : normalize a Gaussian vector

$$\theta = \frac{Z}{\|Z\|}, \quad Z \sim \mathcal{N}(0, I_k)$$

To project from  $\mathbb{R}^d$  to  $\mathbb{R}^k$

$$Z = \begin{pmatrix} Z_{11} & Z_{12} & \cdots & Z_{1k} \\ Z_{21} & Z_{22} & \cdots & Z_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ Z_{d1} & Z_{d2} & \cdots & Z_{dk} \end{pmatrix} \in \mathbb{R}^{d \times k} \quad Z_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$$

Gram-Schmidt yields  $Z = QR$

deliver  $\Theta = Q \in \mathbb{R}^{d \times k}$

project  $\widetilde{\mathbf{X}}_i = \Theta^\top \mathbf{X}_i$

Any  $QR$  decomposition with positive  $R_{ii}$  will do.

# Lopes et al. ctd.

They make just **one** random projection of the data

They find that  $k \approx (n_x + n_y - 2)/2$  performs well

## Why just one?

If your one projection is 'unlucky' then you might miss the pattern.

But with just one projection the distribution of  $\tilde{T}^2$  is known.

## Multiple projections

$$\bar{T}^2 = \frac{1}{M} \sum_{i=1}^M \tilde{T}_i^2$$

$$\tilde{T}_i^2 \text{ based on } \Theta_i \in \mathbb{R}^{d \times k}$$

Get some kind of 'average' luck.

But distribution of  $\bar{T}^2$  not known.

# Multiple projections

Work with **S. Emerson**: average over  $M$  independent random  $\Theta_i \in \mathbb{R}^{d \times k}$

$$\bar{T}^2 = \frac{1}{M} \sum_{i=1}^M \tilde{T}_i^2, \quad \text{where}$$

$$\tilde{T}_i^2 = \frac{n_x n_y}{n_x + n_y} (\bar{\mathbf{X}} - \bar{\mathbf{Y}})^\top \Theta_i (\Theta_i^\top S \Theta_i)^{-1} \Theta_i^\top (\bar{\mathbf{X}} - \bar{\mathbf{Y}})$$

## Easily

- 1)  $\mathbb{E}(\bar{T}^2) = \mathbb{E}(\tilde{T}_i^2)$
- 2)  $\text{Var}(\bar{T}^2) < \text{Var}(\tilde{T}_i^2)$ , unless both are infinite! (averaging reduces variance)

## Less easily

- 1) Finite variance requires  $k \leq n_x + n_y - 6$
- 2) Finite mean requires  $k \leq n_x + n_y - 4$

Unfortunately: the distribution of  $\bar{T}^2$  is not known.

# Separation

Simulate 2000 data sets

$$\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \Sigma), \quad \mathbf{Y}_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\delta, \Sigma), \quad \delta \in \mathbb{R}^d$$

Of these:

1000 **null** cases       $\|\delta\| = 0$

1000 **non-null** cases       $\|\delta\| > 0$

Rank 2000  $\bar{T}^2$  scores

See if nulls get smaller  $\bar{T}^2$  values.

## The ROC\* curve

Shown later, shows how well the test separates the two cases

\* Receiver Operating Characteristic (don't ask)

# Simulated case

$$\mathbf{X}_i, \mathbf{Y}_i \in \mathbb{R}^{200} \quad n_x = n_y = 50$$

Pick  $\|\delta\| = 3$  uniform on 200 dimensional sphere

$$\text{Pick } \Sigma = I_d \times 50/\sqrt{d}$$

## Why these

Uniform  $\delta$  means that the group separation is unrelated to the covariance structure.

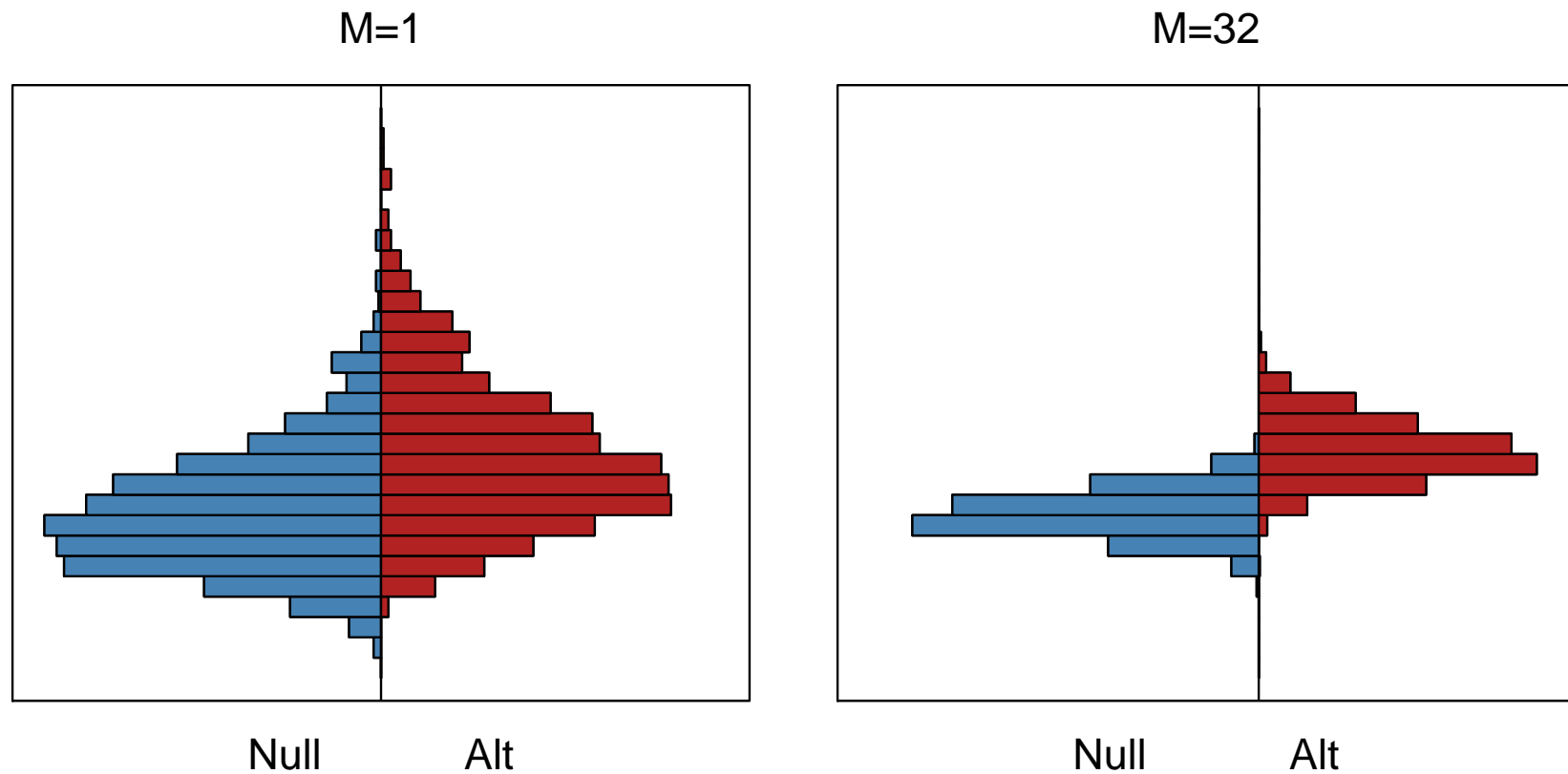
Debatable. We follow [Lopes et al](#) in making this assumption.

WLOG, under uniformity  $\Sigma = \text{diag}(\lambda_1, \dots, \lambda_d)$   $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$

Interesting cases are equal  $\lambda_j$  and rapidly decreasing  $\lambda_j$

# Multiple projections

Simulated T squared



$$n_x = n_y = 50, \quad d = 200, \quad k = 49,$$

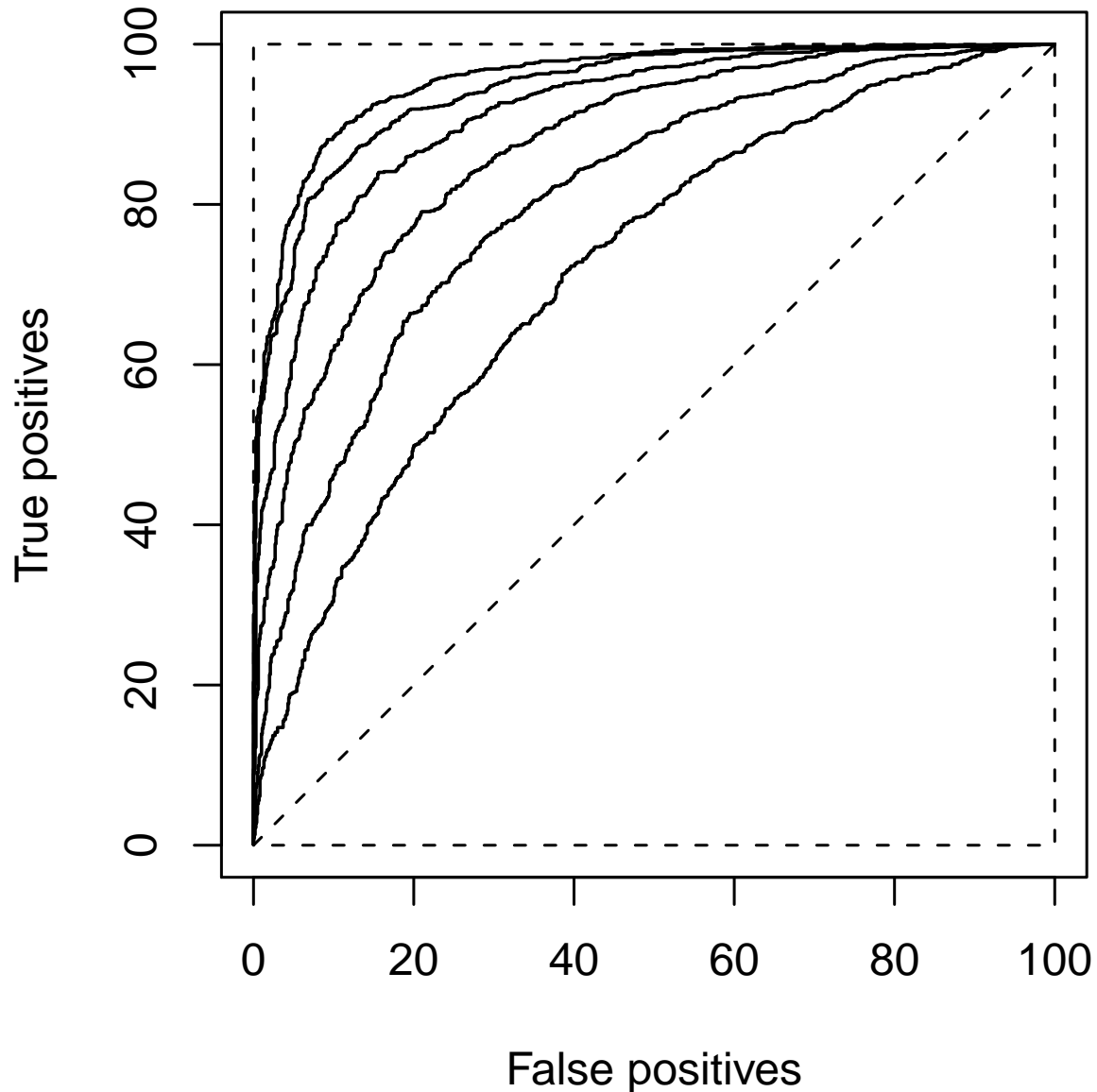
$$\text{Null: } \|\delta\| = 0$$

$$\text{Alt: } \|\delta\| = 3$$



# The ROCs

ROC curves:  $M=1,2,4,8,16,32$



Larger  $M$  has greater area under the curve:

$M$	AUC
1	71.9
2	80.6
4	87.1
8	91.4
16	94.3
32	95.7

## Varying $k$

Lopes et al. prefer  $k \approx (n_x + n_y - 2)/2$

That is not always optimal. But may be a good default.

For the previous scenario: small  $k$  do relatively poorly.

$32 \leq k \leq 56$  all gave  $\text{AUC} \approx 0.95$  with  $M = 32$

## Other scenarios

S. Emerson: advantage of averaging persists in other decay rates for eigenvalues of  $\Sigma$

# Using $\bar{T}^2$

The usual  $p$ -value is  $\Pr(\bar{T}^2 \geq t^2)$  where  $t^2$  is the observed value on our data.

We have no good approximation for this.

Even the moments of  $\bar{T}^2$  involve difficult integrals over  $\Theta \in V_{d,k}$ , the Stiefel manifold, e.g.

$$\begin{aligned} & \int_{\Theta \in V_{d,k}} \Theta(\Theta^\top S \Theta)^{-1} \Theta^\top d\mathbf{U}(\Theta) \\ &= (2\pi)^{-dk/2} \int_{Z \in \mathbb{R}^{d \times k}} Z(Z^\top S Z)^{-1} Z^\top e^{-\text{tr}(Z^\top Z)/2} dZ \end{aligned}$$

Non-negative diagonal  $S \in \mathbb{R}^{d \times d}$  with  $n_x + n_y - 2$  positive entries

$\mathbf{U}(\Theta)$  is the uniform (Haar) measure.

Above is the first moment. Closed forms for first and second moments could lead to useful test statistics.

# Permutation tests

There are  $\binom{n_x+n_y}{n_x}$  ways to allocate  $n_x$  of the pooled observations

$$(\mathbf{X}_1, \dots, \mathbf{X}_{n_x}, \mathbf{Y}_1, \dots, \mathbf{Y}_{n_y})$$

to the first sample (the  $\mathbf{X}$ 's). The re-allocated data:

$$(\mathbf{X}_1^*, \dots, \mathbf{X}_{n_x}^*, \mathbf{Y}_1^*, \dots, \mathbf{Y}_{n_y}^*)$$

have statistic  $\bar{T}^{*2}$ . Then

$$p = \frac{\#\{\bar{T}^{2,*} \mid \bar{T}^{2,*} \geq \bar{T}^2\}}{\binom{n_x+n_y}{n_x}}$$

For  $n_x = n_y = 50$  there are  $\approx 10^{29}$  allocations (permutations).

Use a Monte Carlo sample of random permutations.

Justified in text by [Lehmann & Romano \(2005\)](#), "Testing Statistical Hypotheses"

Combination test might be a better name

# Summary of problem one

If  $d > n_x + n_y - 2$ , then  $T^2$  is not defined

Dempster (almost) used coordinate projections

Lopes et al. use one random projection from  $d$  to  $k$  dimensions

We find benefits from multiple projections

But have to use permutation tests for significance

Monte Carlo enters to compute the test statistic and then to judge its significance

# Problem two

- How can we judge uncertainty in two-way and higher way data?
- We would like to use the bootstrap.
- But McCullagh (2000) proved this is impossible.

## Solution

We will get a Monte Carlo method which mildly overestimates the sampling uncertainty.

## But first

it is necessary to describe the bootstrap, as well as multi-way data.

# Bootstrap sampling

Data are  $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{\text{iid}}{\sim} F$

We compute  $\hat{T} = T(\mathbf{X}_1, \dots, \mathbf{X}_n)$

What is sampling uncertainty in  $\hat{T}$ , e.g.  $\text{Var}(\hat{T}) = \text{Var}(\hat{T} \mid F)$ ?

## Combine two ideas

**Monte Carlo**    Sample from  $F$  to estimate  $\text{Var}(\hat{T} \mid F)$  (but we don't know  $F$ ).

**Plug in**        Use  $\text{Var}(\hat{T} \mid \hat{F})$ , as if  $F = \hat{F}$ , the empirical distribution\*.

From [Efron \(1979\)](#).

∃ extensive variations on the idea.

## \*Empirical distribution

$\hat{F} = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{X}_i}$  puts probability  $1/n$  on each sample observation.

$\delta_{\mathbf{x}}$  is a 'point mass' at  $\mathbf{x}$

To sample  $\mathbf{X} \sim \hat{F}$ , pick one of the original data points

# Bootstrap pseudocode

Given  $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{\text{iid}}{\sim} F$

Resample the data  $B$  times

For  $b = 1, \dots, B$

For  $i = 1, \dots, n$

$$i^* \sim \mathbf{U}\{1, \dots, n\}$$

$$\mathbf{X}_i^* = \mathbf{X}_{i^*}$$

$$T_b^* = T(\mathbf{X}_1^*, \dots, \mathbf{X}_n^*)$$

Compute summary

$$\bar{T}^* = \frac{1}{B} \sum_{b=1}^B T_b^*$$

$$\widehat{\text{Var}}(\hat{T}) = \frac{1}{B-1} \sum_{b=1}^B (T_b^* - \bar{T}^*)^2$$



# First reactions

- 1) How could that ever work?
- 2) How could that ever fail?

The bootstrap works by a continuity argument

Plug-in:

Under conditions,

$\text{Var}(\hat{T} \mid F)$  is continuous in  $F$ , and  $\hat{F} \rightarrow F$ .

So  $\text{Var}(\hat{T} \mid \hat{F}) \rightarrow \text{Var}(\hat{T} \mid F)$

Monte Carlo:

As  $B \rightarrow \infty$ ,  $\widehat{\text{Var}}(\hat{T} \mid \hat{F}) \rightarrow \text{Var}(\hat{T} \mid \hat{F})$

Fails when

Continuity conditions fail or when  $\hat{F}$  does not mimick  $F$  well enough

# Bootstrapping the mean

The base case is for

$$T(\mathbf{X}_1, \dots, \mathbf{X}_n) = \bar{\mathbf{X}} \equiv \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$$

We have non-bootstrap methods for  $\widehat{\text{Var}}(\bar{\mathbf{X}})$  (!)

Correctness for the mean extends to more complicated statistics via Taylor approximations

## Why we like it

No need to assume the Gaussian or any other distributional form

It is explainable to scientific colleagues

# Some variants

The Bayesian bootstrap, Rubin (1981)

$$\hat{F} = \frac{\sum_{i=1}^n W_i \delta_{\mathbf{X}_i}}{\sum_{i=1}^n W_i}$$
$$W_i \stackrel{\text{iid}}{\sim} \text{Exp}(1)$$

That is: place an independent random weight on each observation.

If  $T(\mathbf{X}_1, \dots, \mathbf{X}_n)$  is the sample mean, then under bootstrapping

$$T^* = \frac{\sum_{i=1}^n W_i \mathbf{X}_i}{\sum_{i=1}^n W_i}$$

is a random ratio.

# Weight condition

We need  $\mathbb{E}(W_i) = 1$  and  $\text{Var}(W_i) = 1$

Then the Bayesian bootstrap is comparable to ordinary bootstrap

Can also use  $\text{Poi}(1)$  weights (Poisson distribution)

Used in machine learning [Oza \(2001\)](#)

Half sampling:

$$W_i = \begin{cases} 0 & \text{with probability } 1/2 \\ 2 & \text{with probability } 1/2 \end{cases}$$

# Next

That was the bootstrap

We'll bootstrap 2-way data

# One-way data

Patient	Height	Weight	Age	...	Blood pressure
⋮					
42	1.7	72	32	...	141
43	2.1	97	42	...	109
⋮					

A usual data matrix has  $\mathbf{X}_i$  in row  $i$ .

Data within a row are dependent,

e.g. height and weight of same person are correlated  
as are blood pressure and age

Data in two different rows are independent.

patient 42 and patient 43

# Two way data

patients    ×    medical interns    →    blood pressure measurements  
genes        ×    environments        →    crop yields  
students    ×    exam questions    →    points

## Features

Measurements on the same patient are correlated

Measurements by the same interns are also correlated

(hopefully that effect is smaller)

Neither rows nor columns are IID

# Netflix data

Rating	Viewer 1	Viewer 2	Viewer 3	...	Viewer C
Movie 1	4	4	1	...	4
Movie 2	5	5	NA	...	NA
Movie 3	3	3	NA	...	2
⋮	⋮	⋮	⋮	⋮	⋮
Movie R	NA	5	3	...	4

Ratings on same movie are correlated as are ratings by same person

Danny Deckchair (2003)

Most common rating was 4 stars



# Netflix continued

17,770 movies, 480,189 customers, 100,000,000+ ratings,

## Sample uncertainty

Ratings made on Tuesday came out a little lower than Sunday ratings.

Is it real or a sampling artifact?

## Further problems

- 1) Missing data make the matrix very imbalanced
- 2) Unequal variances, e.g. some customers just give 1s or 5s

# Facebook data

**Alice** (shares a URL) “Hey, check out `http://www.mcqmc2012.unsw.edu.au/`”

**Bob** (comments on it) “Thanks for sharing that, I learned a lot.”

**Data** `url = http://www.mcqmc2012.unsw.edu.au/`

`sharer = Alice`

`commenter = Bob`

`log length X = log(41) ≐ 3.71`

## Data size

18,134,419 comments by 8,078,531 commenters on 2,085,639 URLs

This is 3 way data: `url × sharer × commenter`

## Of interest:

users’ sharing and commenting behaviour,

e.g. who makes longer comments, U.S. or U.K. users?

Probably greater interest for ad clicks, linking and liking activity

# Random effects model

$$X_{ij} = \mu + a_i + b_j + \varepsilon_{ij} \quad i = 1, \dots, R \quad j = 1, \dots, C$$

$$a_i \sim \mathcal{N}(0, \sigma_A^2) \quad \text{e.g. patients}$$

$$b_j \sim \mathcal{N}(0, \sigma_B^2) \quad \text{e.g. interns}$$

$$\varepsilon_{ij} \sim \mathcal{N}(0, \sigma_E^2)$$

## simplest model for two-way data

Used in agriculture

Studied for decades

$\hat{\mu}$  is  $\bar{X}_{\bullet\bullet}$

No bootstrap exists for  $V(\hat{\mu})$

None can exist . . .

. . . McCullagh (2000)

We can't even bootstrap a balanced  $\bar{X}$  !

He rules out resampling, permuting or any “monoid” operations on rows and columns.

# What about classical approaches?

prime reference:

“Variance Components” by  
Searle, Casella, McCulloch (1992)

- Excellent for balanced Gaussian data
- Unbalance  $\implies$  invert large matrices
- Emphasis on homogeneous variances

# McCullagh (2000)

$$\text{For } \hat{\mu} = \bar{X}_{\bullet\bullet} = \frac{1}{R} \frac{1}{C} \sum_{i=1}^R \sum_{j=1}^C X_{ij}$$

**Naive** Resample from  $N = RC$  values

**Product** Resample  $R$  rows and resample  $C$  columns (indep)

$$V_{\text{RE}}(\hat{\mu}) = \frac{\sigma_A^2}{R} + \frac{\sigma_B^2}{C} + \frac{\sigma_E^2}{RC} \quad \text{true var}$$

$$\mathbb{E}_{\text{RE}}(\hat{V}_{\text{Naiv}}(\hat{\mu})) \doteq \left( \sigma_A^2 + \sigma_B^2 + \sigma_E^2 \right) \frac{1}{RC} \quad \text{way too small}$$

$$\mathbb{E}_{\text{RE}}(\hat{V}_{\text{Prod}}(\hat{\mu})) \doteq \frac{\sigma_A^2}{R} + \frac{\sigma_B^2}{C} + \frac{3\sigma_E^2}{RC} \quad \text{not so bad}$$

Naive resampling seriously flawed, product resampling is close

# Notation

Index  $j$  takes values  $i_j = 1, 2, 3, \dots$

Observation  $\mathbf{i}$  has multi-index  $\mathbf{i} = (i_1, \dots, i_r) \in \{1, 2, \dots\}^r$

Random  $X_{\mathbf{i}} \in \mathbb{R}^d$  with  $Z_{\mathbf{i}} = \begin{cases} 1 & X_{\mathbf{i}} \text{ known} \\ 0 & X_{\mathbf{i}} \text{ missing.} \end{cases}$

## Sample size

$$0 < N \equiv \sum_{\mathbf{i} \in \mathbb{N}^r} Z_{\mathbf{i}} < \infty$$

## Sample mean

$$\bar{X} = \frac{\sum_{\mathbf{i}} Z_{\mathbf{i}} X_{\mathbf{i}}}{\sum_{\mathbf{i}} Z_{\mathbf{i}}}$$

We want to estimate  $\text{Var}(\bar{X})$

treating  $Z_{\mathbf{i}}$  as fixed

# $r$ -fold product bootstrap

$$\bar{X}^* = \frac{\sum_i Z_i W_i X_i}{\sum_i Z_i W_i} \quad \text{where}$$

$$W_i = \prod_{j=1}^r W_{j,i_j}$$

$$\mathbb{E}(W_{j,i_j}) = 1$$

$$\text{Var}(W_{j,i_j}) = 1$$

From replications  $\bar{X}^{*1}, \dots, \bar{X}^{*B}$

$$\widehat{\text{Var}}(\bar{X}) = \frac{1}{B-1} \sum_{b=1}^B (X^{*b} - \bar{X}^*)^2.$$

It is operationally much easier to have independent  $W_{j,i_j}$   
(as in the Bayesian bootstrap).

Data for URLs might be scattered over several continents.

Then keeping  $\sum_i W_{j,i_j} = N$  is awkward.

# $r$ -fold random effects

We study  $\widehat{\text{Var}}(\bar{X})$  under this model:

$$X_i = \mu + \sum_{\substack{u \subseteq \{1, \dots, r\} \\ u \neq \emptyset}} \varepsilon_{\mathbf{i}_u, u}$$

If  $u = (j_1, \dots, j_k)$  then  $\mathbf{i}_u = (i_{j_1}, \dots, i_{j_k})$

## Moments

$$\mathbb{E}(\varepsilon_{\mathbf{i}_u, u}) = 0 \quad \text{Var}(\varepsilon_{\mathbf{i}_u, u}) = \sigma_u^2$$

$$\text{Cov}(\varepsilon_{\mathbf{i}_u, u}, \varepsilon_{\mathbf{i}'_{u'}, u'}) = 0 \quad \text{if } u \neq u' \quad \text{or} \quad \mathbf{i}_u \neq \mathbf{i}'_{u'}$$

$\sigma_u^2$  can be allowed to depend on  $\mathbf{i}_u$



# Variance

$$\begin{aligned}
 \text{Var}_{\text{RE}}(\hat{\mu}) &= \frac{1}{N^2} \sum_{u \neq \emptyset} \sum_{u' \neq \emptyset} \sum_{\mathbf{i}} \sum_{\mathbf{i}'} Z_{\mathbf{i}} Z_{\mathbf{i}'} \text{Cov}(\varepsilon_{\mathbf{i},u}, \varepsilon_{\mathbf{i}',u'}) \\
 &= \frac{1}{N^2} \sum_{u \neq \emptyset} \left( \sum_{\mathbf{i}} \sum_{\mathbf{i}'} Z_{\mathbf{i}} Z_{\mathbf{i}'} 1_{\mathbf{i}_u = \mathbf{i}'_u} \right) \sigma_u^2 \\
 &= \frac{1}{N} \sum_{u \neq \emptyset} \nu_u \sigma_u^2,
 \end{aligned}$$

for gain coefficients,

$$\nu_u = \frac{1}{N} \sum_{\mathbf{i}} Z_{\mathbf{i}} N_{\mathbf{i},u}, \quad \text{where}$$

$$N_{\mathbf{i},u} = \sum_{\mathbf{i}'} Z_{\mathbf{i}'} 1_{\mathbf{i}_u = \mathbf{i}'_u}.$$

# Running examples

$$\begin{aligned}
 V_{\text{RE}}(\hat{\mu}) &\equiv \frac{1}{N} \sum_{u \neq \emptyset} \nu_u \sigma_u^2 \\
 &\doteq \frac{1}{N} \left( 56,200 \sigma_{\text{movies}}^2 + 646 \sigma_{\text{viewers}}^2 + \sigma_{\text{interaction}}^2 \right) \quad (\text{for Netflix})
 \end{aligned}$$

## For Facebook

$$\nu_{\text{sh}} \doteq 17.71, \quad \nu_{\text{com}} \doteq 7.71, \quad \nu_{\text{url}} \doteq 26,854.92 \quad !$$

$$\nu_{\text{sh,com}} \doteq 5.92, \quad \nu_{\text{sh,url}} \doteq 12.91, \quad \nu_{\text{com,url}} \doteq 5.19, \quad \text{and}$$

$$\nu_{\text{sh,com,url}} \doteq 4.88.$$

$$\nu_{\text{url}} \geq 26,000$$

# Variances

For gain coefficients  $\nu_u$

$$\text{Var}_{\text{RE}}(\bar{X}) = \frac{1}{N} \sum_{\substack{u \subseteq \{1, \dots, r\} \\ u \neq \emptyset}} \nu_u \sigma_u^2$$

Similarly for  $\gamma_u$  (defined later):

$$\mathbb{E}_{\text{RE}}(\widehat{\text{Var}}_{\text{Prod}}(\bar{X}^*)) \doteq \frac{1}{N} \sum_{\substack{u \subseteq \{1, \dots, r\} \\ u \neq \emptyset}} \gamma_u \sigma_u^2$$

And:

$$\mathbb{E}_{\text{RE}}(\widehat{\text{Var}}_{\text{Naiv}}(\bar{X}^*)) \doteq \frac{1}{N} \sum_{\substack{u \subseteq \{1, \dots, r\} \\ u \neq \emptyset}} \sigma_u^2$$

Typically  $1 \ll \nu_u \ll N$  for  $u \neq \{1, \dots, r\} \implies$  Naive bootstrap unreliable

We want  $\gamma_u = \nu_u$ . We'll get  $\gamma_u \geq \nu_u$ .

# The case $r = 2$

O (2007)

Independent bootstrap of rows and columns

Still get  $\mathbb{E}_{\text{RE}}(\widehat{V}_{\text{Prod}}(\hat{\mu})) \doteq V_{\text{RE}}(\hat{\mu})$ , i.e.

Still get  $\approx 1$   $\times$  the main effect contribution

$\approx 3$   $\times$  the interaction contribution

Sunday vs. Tuesday edge of 0.02 stars is real (about 8 standard errors)

New here

- 1) Arbitrary order  $r \geq 2$
- 2) Independent product weights (vs resampling)

# Product bootstrap

$$\hat{\mu}^* = \frac{\sum_i Z_i W_i X_i}{\sum_i Z_i W_i} \equiv \frac{T^*}{N^*} \quad (\text{ratio estimator})$$

$$V_{\text{Prod}}(\hat{\mu}^*) \approx \tilde{V}_{\text{Prod}}(\hat{\mu}^*) \equiv \frac{1}{N^2} \mathbb{E}_{\text{Prod}}((T^* - \hat{\mu} N^*)^2)$$

## Main result

$$\mathbb{E}_{\text{RE}}(\tilde{V}_{\text{Prod}}(\hat{\mu}^*)) = \frac{1}{N} \sum_{u \neq \emptyset} \gamma_u \sigma_u^2$$

where  $\gamma_u \approx \nu_u$  if  $|u| = 1$ , (i.e. cardinality 1)  
 otherwise small  $\gamma_u/\nu_u > 1$

# Exact formula depends on

# index pairs  $i, i'$  that match in set  $u$ :

$$\sum_{i \in \mathbb{N}^r} \sum_{i' \in \mathbb{N}^r} Z_i Z_{i'} 1_{i_u = i'_u}$$

# index pairs  $i, i'$  that match in **precisely**  $k$  places

# index triples  $i, i', i''$  where  $i$  matches  $i'$  in set  $u$  and matches  $i''$  in precisely  $k$  places

# Duplication indices

$$\text{(level dup)} \quad \epsilon = \frac{\text{greatest item popularity}}{N}$$

$$\text{(variable dup)} \quad \eta = \max_{\emptyset \subsetneq u \subsetneq v} \frac{\nu_v}{\nu_u}$$

## Examples

	$\epsilon$	$\eta$
<b>Netflix</b>	$\frac{232,944}{100,480,507} \doteq 0.00232$	$\frac{1}{646} \doteq 0.00155$
	Miss Congeniality	$\nu_{\text{interaction}} / \nu_{\text{movies}}$
<b>Facebook</b>	$\frac{686,990}{18,134,419} \doteq 0.0379$	$\frac{4.88}{5.19} \doteq 0.94$
	a popular URL	$\nu_{\text{sh,com,url}} / \nu_{\text{com,url}}$

$\eta$  is not small for the Facebook data

bootstrap variances will be somewhat more conservative

# Approximations

**Theorem 1.** *In the homogeneous random effects model, the product weight bootstrap with  $V(W_{j,i_j}) = \tau^2 = 1$ , satisfies*

$$\gamma_u = \nu_u [2^{|u|} - 1 + \Theta_u \epsilon] + \sum_{v \supsetneq u} 2^{|v|} \nu_v,$$

where  $|\Theta_u| \leq 2^{r+1} - 2$ .

*Proof.* [O & Eckles \(2011\)](#), who consider general  $\tau^2$ . □

For small  $\epsilon$  and  $r$  (i.e.  $2^r \epsilon \ll 1$ )

$$\gamma_u \approx (2^{|u|} - 1) \nu_u + \sum_{v \supsetneq u} 2^{|v|} \nu_v$$

If also  $\eta \ll 1$

$$\gamma_u \approx (2^{|u|} - 1) \nu_u$$



# Some specific approximations

For  $r = 2$

$$\begin{aligned}\gamma_{\{j\}} &= \nu_{\{j\}}(1 + \Theta_j \epsilon) + 2 \quad j = 1, 2 \\ \gamma_{\{1,2\}} &= \nu_{\{1,2\}}(3 + \Theta_{\{1,2\}} \epsilon), \quad \text{where} \\ |\Theta_u| &\leq 6.\end{aligned}$$

For  $r = 3$

$$\begin{aligned}\gamma_{\{1\}} &\approx \nu_{\{1\}} + 4\nu_{\{1,2\}} + 4\nu_{\{1,3\}} + 8 \\ \gamma_{\{1,2\}} &\approx 3\nu_{\{1,2\}} + 8 \\ \gamma_{\{1,2,3\}} &\approx 7.\end{aligned}$$

If  $0 < m \leq \min_u \sigma_u^2 \leq \max_u \sigma_u^2 \leq M < \infty$  then

$$\frac{\mathbb{E}_{\text{RE}}(\tilde{V}_{\text{Prod}}(\hat{\mu}^*))}{V_{\text{RE}}(\hat{\mu})} = 1 + O(\eta + \epsilon).$$

# Facebook loquacity

For each commenter, url and sharer, we obtain:

$X = \log(\#char \text{ in comment})$  as well as,  
country  $c \in \{US, UK\}$  of commenter, and  
mode  $m \in \{web, mobile\}$  of commenter.

Now let

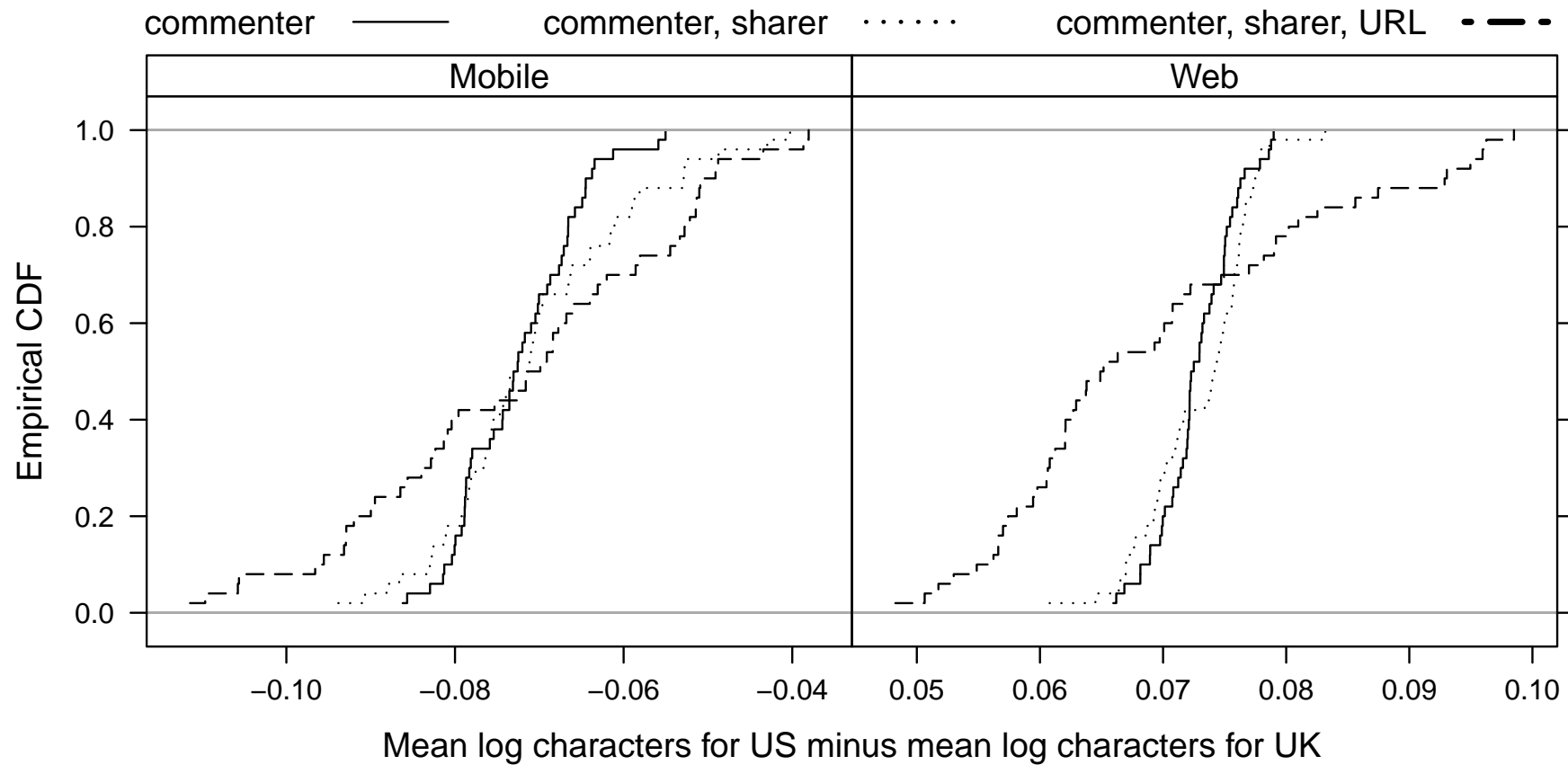
$$\hat{\mu}_{cm} = \frac{\sum_i Z_i X_i 1_{\text{country}=c} 1_{\text{mode}=m}}{\sum_i Z_i 1_{\text{country}=c} 1_{\text{mode}=m}}$$

We see small differences

	US	UK
web	3.62	3.55
mobile	3.50	3.57

but they're larger than sample fluctuations

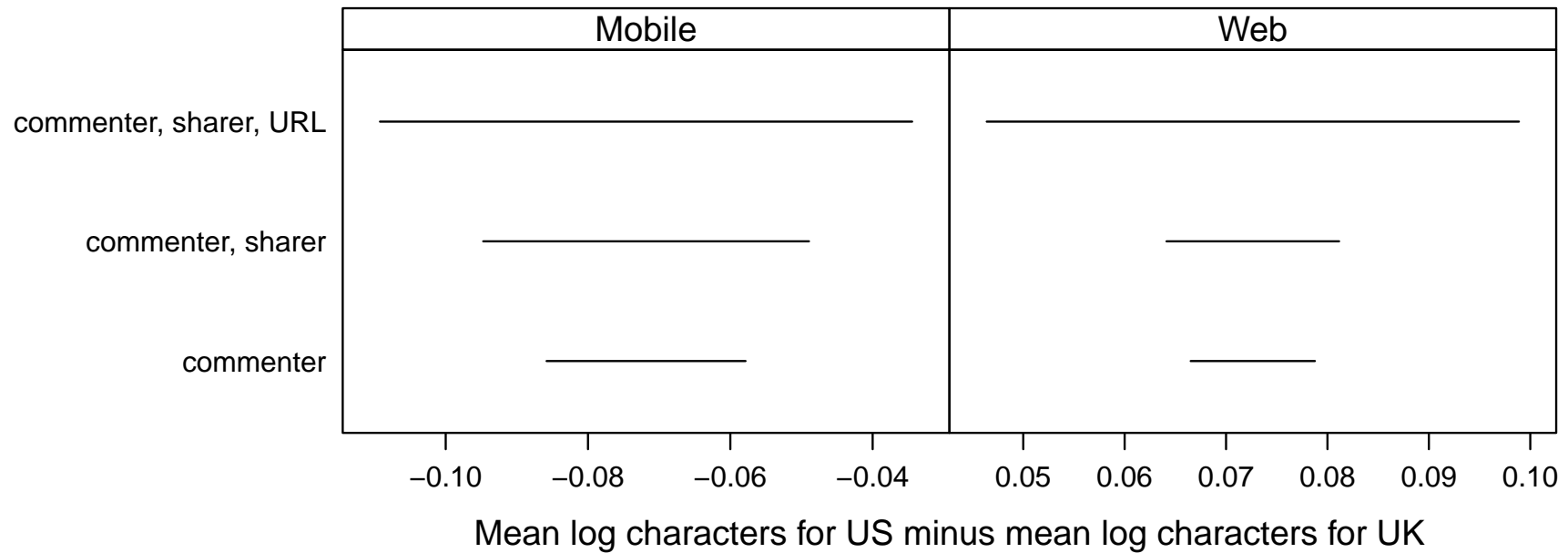
# Loquacity ECDFs



ECDF over 50 bootstraps of  $\hat{\mu}_{USm} - \hat{\mu}_{UKm}$

Reweighting one, two, or three ways

# Loquacity confidence intervals



Central 95% confidence intervals from 50 bootstraps of  $\hat{\mu}_{USm} - \hat{\mu}_{UKm}$

Reweighting one, two, or three ways

# Summary of problem two

Crossed random effects data require special care

No correct bootstrap exists

Product sampling is at least conservative, mildly over-estimating the variance

It works also with unequal variances

what remains to do

better seeding

extensions to more complicated analyses

# Conclusion

Statistics and machine learning are still finding new ways to consume Monte Carlo ideas

In addition to MCMC there are:

- permutations
- rotations
- projections
- subsampling
- reweighting

# Thanks

- Collaborators: [Dean Eckles](#) and [Sarah Emerson](#)
- NSF DMS-0906056
- Data: Netflix and Facebook
- UNSW
- Organizers: [Ian Sloan](#), [Frances Kuo](#), [Josef Dick](#) and [Gareth Peters](#)

# Missingness

The  $\approx 100,000,000$  movie ratings we see are only 1% of  $R \times C$ .

Those we see are probably biased towards higher values. How do we adjust?

**Ans:** we can't, because nobody knows the bias function.

It would not be reasonable to expect a bootstrap method to fix missing data bias. For instance how could one sampling method fix both positive bias and negative bias?

Given an adjustment algorithm (based on assumptions or knowledge from outside the given data) we might be able to bootstrap its predictions.