

Backfitting and crossed random effects

Art B. Owen

and

Trevor Hastie

and

Swarnadip Ghosh

<https://arxiv.org/abs/2007.10612>

Summary

Crossed random effects are common:

Plant varieties \times environments

Customers \times products

Hospitals \times dialysis centers

He, Kalbfleisch, Li, Li (2013) on readmissions

Simple models cost $O(N^{3/2})$ (or worse)

Thesis and papers of Katelyn Gao

Our contribution

We can get $O(N)$ cost via backfitting

Buja, Hastie, Tibshirani (1989)

Linear models today

GLMMs “real soon now”

Example: clothing approval, $N = 5,000,000$ ratings from Stitch Fix

major thanks to Brad Klingenberg

Categorical variables

- treated vs untreated patients
- 3 kinds of iris flower
- 50 US states

Some factors have many more levels

Factors

- Product SKU
e.g., dust filter for a Hoover Max Extract Pressure Pro model 60
- Query string:
“covid”, “Biden”, . . . , “heteroscedasticity”
- Customer ID, URL, IP address
these ‘churn’

Some differences

- 1) can have millions of levels
- 2) power law frequency
- 3) continually changing set of levels

We **might** want to treat factors as random effects.

Crossed random effects

Customers × products

Students × questions

Genotype × environments

Diners × restaurants

Review × book × mailing address × credit card #

I.E., two or more of these variables

More generally: many \leftrightarrow many relationships

Simple model

$$Y_{ij} = \mathbf{x}_{ij}^T \beta + a_i + b_j + \varepsilon_{ij}$$

$$a_i \sim (0, \sigma_A^2) \quad b_j \sim (0, \sigma_B^2) \quad \varepsilon_{ij} \sim (0, \sigma_E^2) \quad \text{indep.}$$

Modest because

Only one layer

No latent factors to discover communities / genres

Already challenging

- Gaussian likelihood costs $O(N^{3/2})$ to evaluate once
Gao & O (2019)
- Gibbs takes $O(N^{1/2})$ iterations at cost $O(N)$ each
Gao & O (2017)

Generalization

- Statistics runs on replication
- Easy from IID data
- Ok for hierarchical data
 - dependent with clusters
 - independent between

Crossed effects

Hold out rows

⇒ dependence comes in via columns

Borders on an $n = 1$ setting

Subjective difficulty ladder

IID \preceq hierarchical \preceq time series \preceq crossed effects \preceq networks

Balanced and unbalanced

Simplest if all ij pairs observed.

Our use cases are unbalanced.

Notation

'Rows' $i = 1, \dots, R$

'Columns' $j = 1, \dots, C$

$Z_{ij} = 1 \iff (\mathbf{x}_{ij}, Y_{ij})$ observed (0 else)

More indices

$$1 \leq r \leq R \quad 1 \leq s \leq C$$

Sample sizes

$$N_{i\bullet} = \sum_{j=1}^C Z_{ij} \quad \text{'size' of row } i$$

$$N_{\bullet j} = \sum_{i=1}^R Z_{ij} \quad \text{'size' of col } j$$

$$N = \sum_{i=1}^R N_{i\bullet} = \sum_{j=1}^C N_{\bullet j} \quad \text{sample size}$$

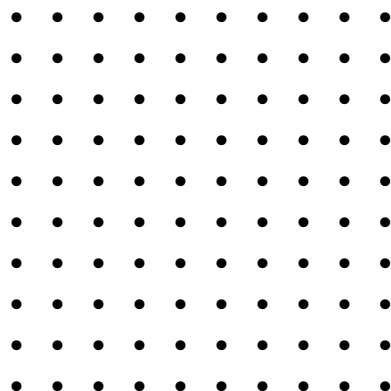
Sparsity

$$1 \ll R, C \ll N \ll R \times C$$

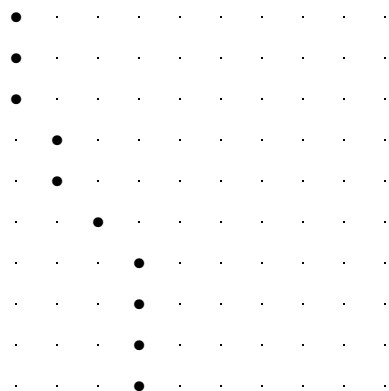
Observation patterns

Solid for $Z_{ij} = 1$ dot/invisible for $Z_{ij} = 0$

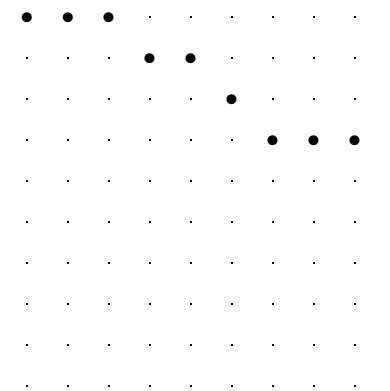
Balanced



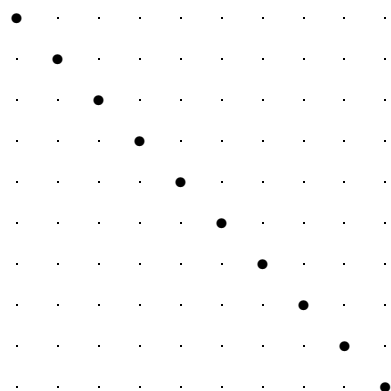
Row nested in col



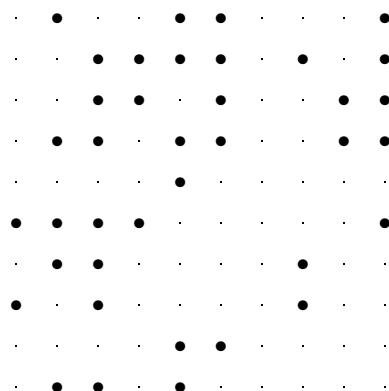
Col nested in row



IID



Arbitrary



Informative missingness

Movie / TV ratings biased high

Restaurant ratings biased towards extremes

Handling informative missingness

- Requires information from outside the data
- Pass for now
- Linear mixed models are hard enough already
- Maybe propensity methods will help

OLS and GLS

Ordinarily least squares and generalized least squares

$$\hat{\beta}_{\text{OLS}} = (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \mathcal{Y}$$

$$\hat{\beta}_{\text{GLS}} = (\mathcal{X}^T \mathcal{V}^{-1} \mathcal{X})^{-1} \mathcal{X}^T \mathcal{V}^{-1} \mathcal{Y}$$

In compatible order

$$\mathcal{X} \in \mathbb{R}^{N \times p}$$

rows are \mathbf{x}_{ij}

$$\mathcal{Y} \in \mathbb{R}^N$$

elements are Y_{ij}

$$\mathcal{V} \in \mathbb{R}^{N \times N}$$

Cov(\mathcal{Y})

Nota Bene

$\mathbf{x}_{ij}, \beta \in \mathbb{R}^p$ p not large and not growing with N

Two problems with OLS

OLS is **inefficient**:

$$\text{Var}(\hat{\beta}_{\text{OLS}}) \succcurlyeq \text{Var}(\hat{\beta}_{\text{GLS}})$$

low power

OLS is **naive**:

$$\widehat{\text{Var}}_{\text{OLS}}(\hat{\beta}_{\text{OLS}}) \preccurlyeq \text{Var}_{\text{GLS}}(\hat{\beta}_{\text{GLS}}) \quad (\text{in expectation})$$

false discoveries

Toy example of $\mathcal{V} = \text{Cov}(\mathcal{Y})$

$R = 3$ rows and $C = 4$ columns with $N = 8$ observations

$$\begin{array}{c} Z \\ 1 \\ 2 \\ 3 \end{array} \begin{array}{cccc} 1 & 2 & 3 & 4 \\ \left[\begin{array}{cccc} 1 & 1 & \cdot & \cdot \\ \cdot & 1 & 1 & 1 \\ 1 & 1 & \cdot & 1 \end{array} \right] \end{array}$$

Labels $\ell = 1, \dots, N$

$$\begin{array}{c} Z \\ 1 \\ 2 \\ 3 \end{array} \begin{array}{cccc} 1 & 2 & 3 & 4 \\ \left[\begin{array}{cccc} 1 & 2 & \cdot & \cdot \\ \cdot & 3 & 4 & 5 \\ 6 & 7 & \cdot & 8 \end{array} \right] \end{array}$$

E.g. observation $\ell = 7$ is in row $i = 3$ and column $j = 2$.

Correlation structure

Correlations come from common rows or columns

Row correlations of a_i , $R = 3$

$$\begin{array}{c}
 \begin{array}{cccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\
 1 & \left[\begin{array}{cccccccc}
 1 & 1 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
 1 & 1 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
 \cdot & \cdot & 1 & 1 & 1 & \cdot & \cdot & \cdot \\
 \cdot & \cdot & 1 & 1 & 1 & \cdot & \cdot & \cdot \\
 \cdot & \cdot & 1 & 1 & 1 & \cdot & \cdot & \cdot \\
 \cdot & \cdot & \cdot & \cdot & \cdot & 1 & 1 & 1 \\
 \cdot & \cdot & \cdot & \cdot & \cdot & 1 & 1 & 1 \\
 \cdot & \cdot & \cdot & \cdot & \cdot & 1 & 1 & 1
 \end{array} \right]
 \end{array}
 \end{array}$$

Column correlations of b_j , $C = 4$

$$\begin{array}{c}
 \begin{array}{cccccccc}
 & 1 & 6 & 2 & 3 & 7 & 4 & 5 & 8 \\
 1 & \left[\begin{array}{cccccccc}
 1 & 1 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
 1 & 1 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
 \cdot & \cdot & 1 & 1 & 1 & \cdot & \cdot & \cdot \\
 \cdot & \cdot & 1 & 1 & 1 & \cdot & \cdot & \cdot \\
 \cdot & \cdot & 1 & 1 & 1 & \cdot & \cdot & \cdot \\
 \cdot & \cdot & \cdot & \cdot & \cdot & 1 & \cdot & \cdot \\
 \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & 1 & 1 \\
 \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & 1 & 1
 \end{array} \right]
 \end{array}
 \end{array}$$

Cov(\mathcal{Y}) in row order

$$\mathcal{V} = \sigma_A^2 \begin{pmatrix} 1 & 1 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & 1 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & 1 & 1 & 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & 1 & 1 & 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & 1 & 1 & 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & 1 & 1 & 1 \\ \cdot & \cdot & \cdot & \cdot & \cdot & 1 & 1 & 1 \\ \cdot & \cdot & \cdot & \cdot & \cdot & 1 & 1 & 1 \end{pmatrix} + \sigma_B^2 \begin{pmatrix} 1 & \cdot & \cdot & \cdot & \cdot & 1 & \cdot & \cdot \\ \cdot & 1 & 1 & \cdot & \cdot & \cdot & 1 & \cdot \\ \cdot & 1 & 1 & \cdot & \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & 1 & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 & \cdot & \cdot & 1 \\ 1 & \cdot & \cdot & \cdot & \cdot & 1 & \cdot & \cdot \\ \cdot & 1 & 1 & \cdot & \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 & \cdot & \cdot & 1 \end{pmatrix} + \sigma_E^2 I$$

We need $\mathcal{V}^{-1}\mathcal{X}$.

Sherman-Morrison-Woodbury does not do it.

Sad trombone sound.

Linear mixed models

Naive algebra costs $O(N^3)$

Actual algebra costs $O((R + C)^3)$

$$RC \geq N \implies \max\{R, C\} \geq \sqrt{N} \implies (R + C)^3 > N^{3/2}$$

Upshot

Crossed: superlinear cost.

Nested: blocks \implies linear cost.

LMM computation

The best is [Doug Bates'](#) most recent Julia code.

It costs $O(N^{3/2})$ to evaluate the likelihood **once**.

What about MCMC?

- Plain Gibbs takes $O(\sqrt{N})$ iterations at $O(N)$ cost Gao & O (2017)
- quite unlike successes in the nested case, e.g.
Yu & Meng (2011) interweaving, Gelman et al. STAN
- Problems with: block Gibbs, reparameterization, Langevin, MALA, Indep sampler, RWW, RWM subsampling, pCN
- Bates et al. (2015) took Bayes out of `lme4` ... comput'n deemed unreliable

Literature check

Nested: Lots of MCMC papers, theory and applied, hierarchical models.

Crossed: Very few MCMC papers.

State of the art

1) Cameron, Gelbach & Miller (2011)

OLS, non-naive via Huber-White

2) Gao & O method of moments, inefficient but not naive.

Efficient if $\sigma_A^2 \approx 0$ or $\sigma_B^2 \approx 0$

3) Papaspiliopoulos, Roberts, & Zanella (2020). hereafter PRZ

Collapsed Gibbs sampler. [Right way to do Bayes.]

Similar idea [Johndrow](#) (personal communication).

PRZ (2020)

✓ Analytically integrate out intercept.

Huge improvement.

✗ Assume all $N_{i\bullet} = N/R$ and all $N_{\bullet j} = N/C$.

Then mixing time is $O(1 \times \text{unknown})$

Unknown mixing time of Gibbs walk on Z_{ij} :

on later slide

Backfitting

Robinson (1991) proves GLS \equiv

$$\min_{\beta, \mathbf{a}, \mathbf{b}} \|\mathcal{Y} - \mathcal{X}\beta - \mathcal{Z}_A \mathbf{a} - \mathcal{Z}_B \mathbf{b}\|^2 + \lambda_A \|\mathbf{a}\|^2 + \lambda_B \|\mathbf{b}\|^2$$

for observation matrices

$$\mathcal{Z}_A \in \{0, 1\}^{N \times R} \quad \mathcal{Z}_B \in \{0, 1\}^{N \times C}$$

and 'ridge' penalties

$$\lambda_A = \frac{\sigma_E^2}{\sigma_A^2} \quad \lambda_B = \frac{\sigma_E^2}{\sigma_B^2}$$

Basic backfit

Update \mathbf{a} then \mathbf{b} then \mathbf{a} etc.

= Block coordinate descent = Gauss-Seidel

For one random effect

$$\mathcal{Y} = \mathcal{X}\beta + \mathcal{Z}_A\mathbf{a} + \mathbf{e}$$

So solve

$$\min_{\beta, \mathbf{a}} \|\mathcal{Y} - \mathcal{X}\beta - \mathcal{Z}_A\mathbf{a}\|^2 + \lambda_A \|\mathbf{a}\|$$

Normal equations

$$0 = \mathcal{X}^\top (\mathcal{Y} - \mathcal{X}\hat{\beta} - \mathcal{Z}_A\hat{\mathbf{a}})$$

$$0 = \mathcal{Z}_A^\top (\mathcal{Y} - \mathcal{X}\hat{\beta} - \mathcal{Z}_A\hat{\mathbf{a}}) - \lambda_A \hat{\mathbf{a}}$$

We get

$$\mathcal{Z}_A\hat{\mathbf{a}} = \left[\mathcal{Z}_A (\mathcal{Z}_A^\top \mathcal{Z}_A + \lambda_A I_R)^{-1} \mathcal{Z}_A^\top \right] (\mathcal{Y} - \mathcal{X}\hat{\beta})$$

$$\equiv \mathcal{S}_A (\mathcal{Y} - \mathcal{X}\hat{\beta}) \quad \text{“smoother matrix” } \mathcal{S}_A$$

$$\hat{\beta} = (\mathcal{X}^\top (I_N - \mathcal{S}_A) \mathcal{X})^{-1} \mathcal{X}^\top (I_N - \mathcal{S}_A) \mathcal{Y}$$

NB: \mathcal{S}_A does shrunken row averages

Two effects

For generic response $\mathcal{R} \in \mathbb{R}^N$, we alternate

$$\mathcal{Z}_A \hat{\mathbf{a}} \leftarrow \mathcal{S}_A(\mathcal{R} - \mathcal{Z}_B \hat{\mathbf{b}})$$

$$\mathcal{Z}_B \hat{\mathbf{b}} \leftarrow \mathcal{S}_B(\mathcal{R} - \mathcal{Z}_A \hat{\mathbf{a}})$$

It converges

Buja, Hastie, Tibshirani (1990).

Limit equals two-factor smoother matrix that we can apply to \mathcal{Y} and columns of \mathcal{X} .

Details in

Ghosh, Hastie & Owen arXiv:2007.10612

Cost

$O(N)$ per iteration.

\implies We must bound the # iterations.

Centering

Usually \mathcal{X} has an intercept

Then the solution has

$$\sum_{i=1}^R \hat{a}_i = 0 = \sum_{j=1}^C \hat{b}_j. \quad (*)$$

Why make the iterations discover $(*)$?

let's just bake it in

Analagous to collapsed sampling

Simple centering

After each iteration:

$$a_i \leftarrow a_i - \bar{a}, \quad \text{for } \bar{a} = \frac{1}{R} \sum_{i=1}^R a_i$$
$$b_j \leftarrow b_j - \bar{b}$$

Principled centering

$$\min_{\mathbf{a}} \|\mathcal{R} - \mathcal{Z}_A \mathbf{a}\|^2 + \lambda_A \|\mathbf{a}\|^2 \quad \text{subject to } \sum_{i=1}^R a_i = 0$$

Theorem 3.2

$$\mathcal{R}_{i\bullet} = \sum_{j=1}^C Z_{ij} \mathcal{R}_{ij} \quad \text{row sums}$$

$$w_i = \frac{(N_{i\bullet} + \lambda_A)^{-1}}{\sum_r (N_{r\bullet} + \lambda_A)^{-1}} \quad \text{normalized weights}$$

Then

$$\hat{a}_i \leftarrow \frac{\mathcal{R}_{i\bullet} - \sum_r w_r \mathcal{R}_{r\bullet}}{N_{i\bullet} + \lambda_A}$$

Stitch Fix

Stylists select clothing and send 5 items to clients

Clients buy some and return others

Variables

- Customer i
- Garment j
- Features \mathbf{x}_{ij} price, size, materials, brand, ZIP code \dots

Model

$$Y_{ij} \sim \mathbf{x}_{ij}^T \beta + a_i + b_j + \varepsilon_{ij}$$

Response

Y_{ij} : size ok 0/1 or liked (10 point scale)

Enormous thanks to [Brad Klingenberg](#) for data.

Stitch Fix data

$N = 5,000,000$ ratings by $R = 762,752$ clients on $C = 6,318$ items.

This a subset of their customer / inventory base.

we got no identifying features

Data is ≥ 5 years old.

Ratings Y_{ij} are on a 10 point scale.

Predictors

$\text{Match}_{ij} \in [0, 1]$, a prediction from some baseline model

(not representative of all their algos).

Whether item is 'Edgy' or 'Boho'.

Same for client.

Material type: leather, fur, acrylic, \dots , wool.

$p = 30$, including intercept.

Analysis

- We only look at one regression model
- Explore OLS naivete
- Explore OLS inefficiency

Variance components

Gao & O (2019)

Moments \rightarrow consistent $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$ and $\hat{\sigma}_E^2$ in $O(N)$ work

We get λ_A and λ_B

Thesis Gao (2017)

$\hat{\sigma}_A^2$ etc asymptotically normal

Model for ratings

For each observed client-item pair (i, j) :

$$\begin{aligned}
 Y_{ij} = & \beta_0 + \beta_1 \text{Match}_{ij} + \beta_2 \mathbb{I}\{\text{client edgy}\}_i + \beta_3 \mathbb{I}\{\text{item edgy}\}_j \\
 & + \beta_4 \mathbb{I}\{\text{client edgy}\}_i \times \mathbb{I}\{\text{item edgy}\}_j + \beta_5 \mathbb{I}\{\text{client boho}\}_i \\
 & + \beta_6 \mathbb{I}\{\text{item boho}\}_j + \beta_7 \mathbb{I}\{\text{client boho}\}_i \times \mathbb{I}\{\text{item boho}\}_j \\
 & + \beta_8 \text{Material}_{ij} + a_i + b_j + e_{ij}
 \end{aligned}$$

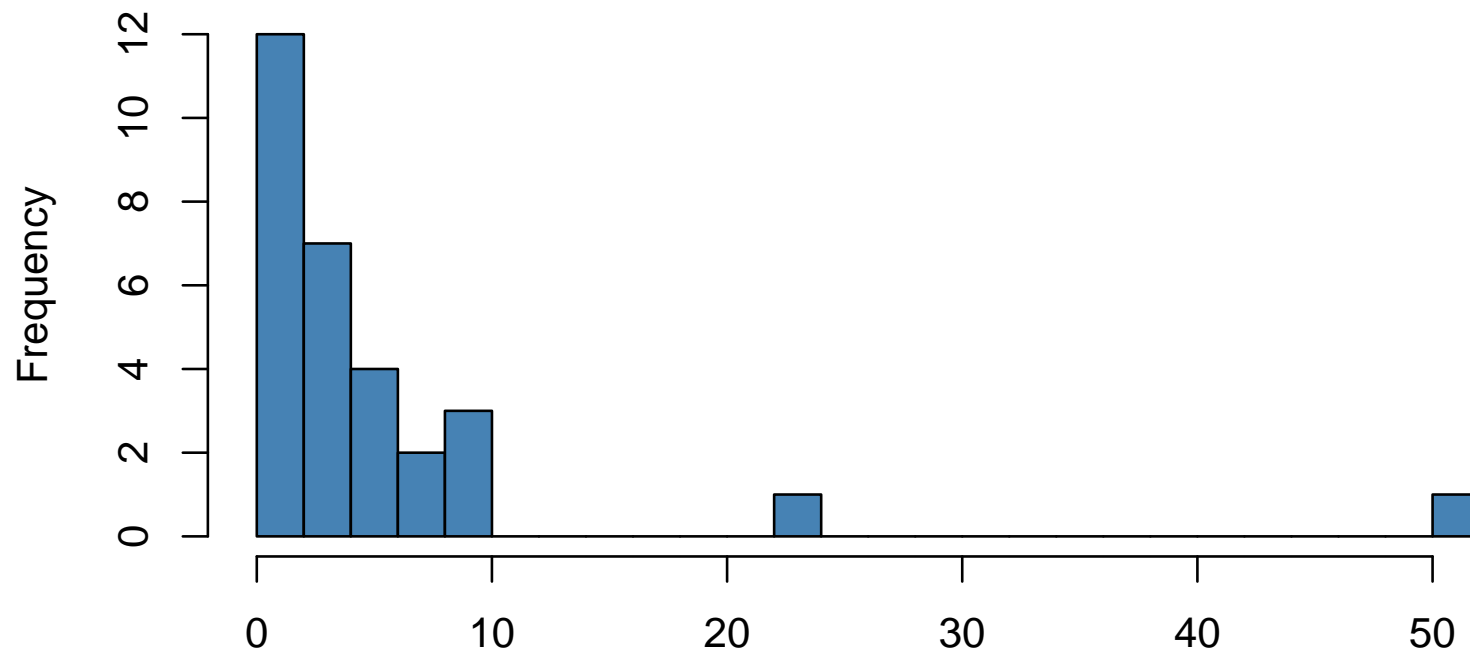
Notes

- Categorical $\text{Material}_{ij} \implies$ Indicator variables (baseline = Polyester)
- $p = 30$
- Gao & O found edgy items to edgy clients worked
(but even boho clients tended not to like boho items)

Inefficiency of OLS

$$1 < \frac{\widehat{\text{Var}}_{\text{GLS}}(\hat{\beta}_{\text{OLS},j})}{\widehat{\text{Var}}_{\text{GLS}}(\hat{\beta}_{\text{GLS},j})} < 51 \quad j = 0, \dots, 29$$

Inefficiency of OLS by coefficient



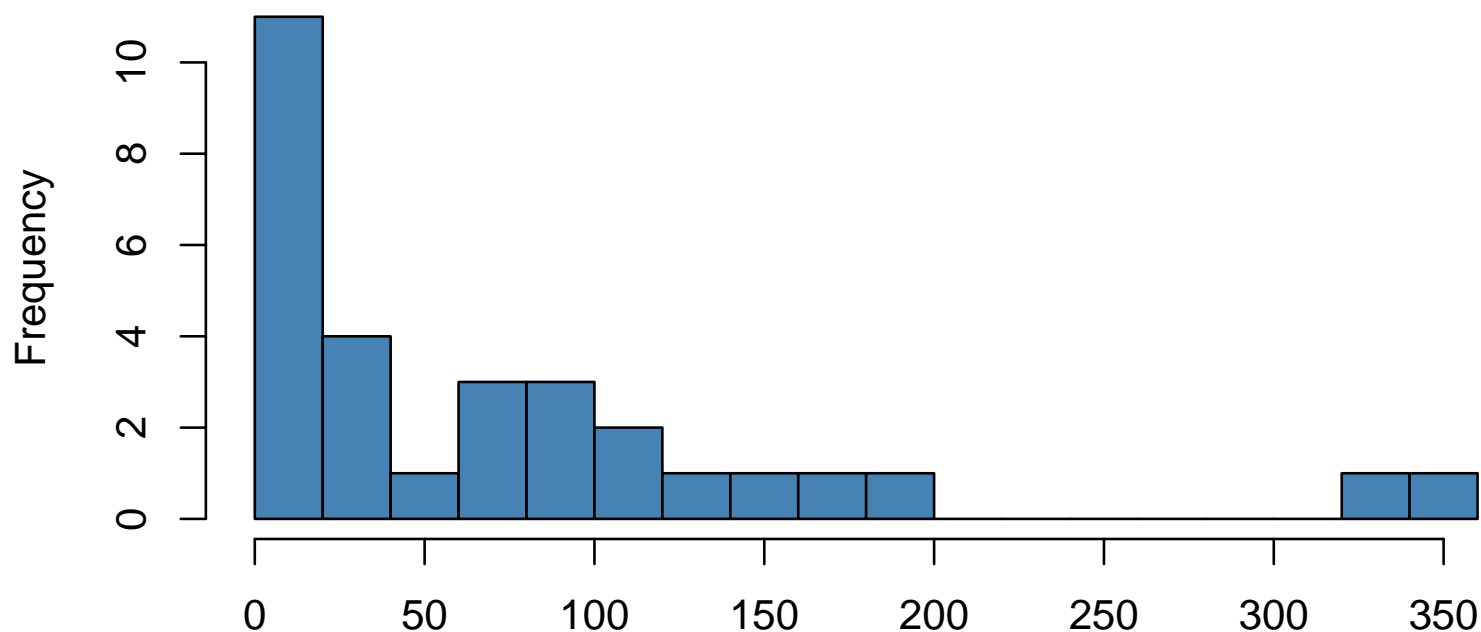
E.g., ratio = 5 is like ignoring 80% of information

Naivete of OLS

$$1.75 < \frac{\widehat{\text{Var}}_{\text{GLS}}(\hat{\beta}_{\text{OLS},j})}{\widehat{\text{Var}}_{\text{OLS}}(\hat{\beta}_{\text{OLS},j})} < 350 \quad j = 1, \dots, 30$$

Worst is material = Modal, next is Tencel.

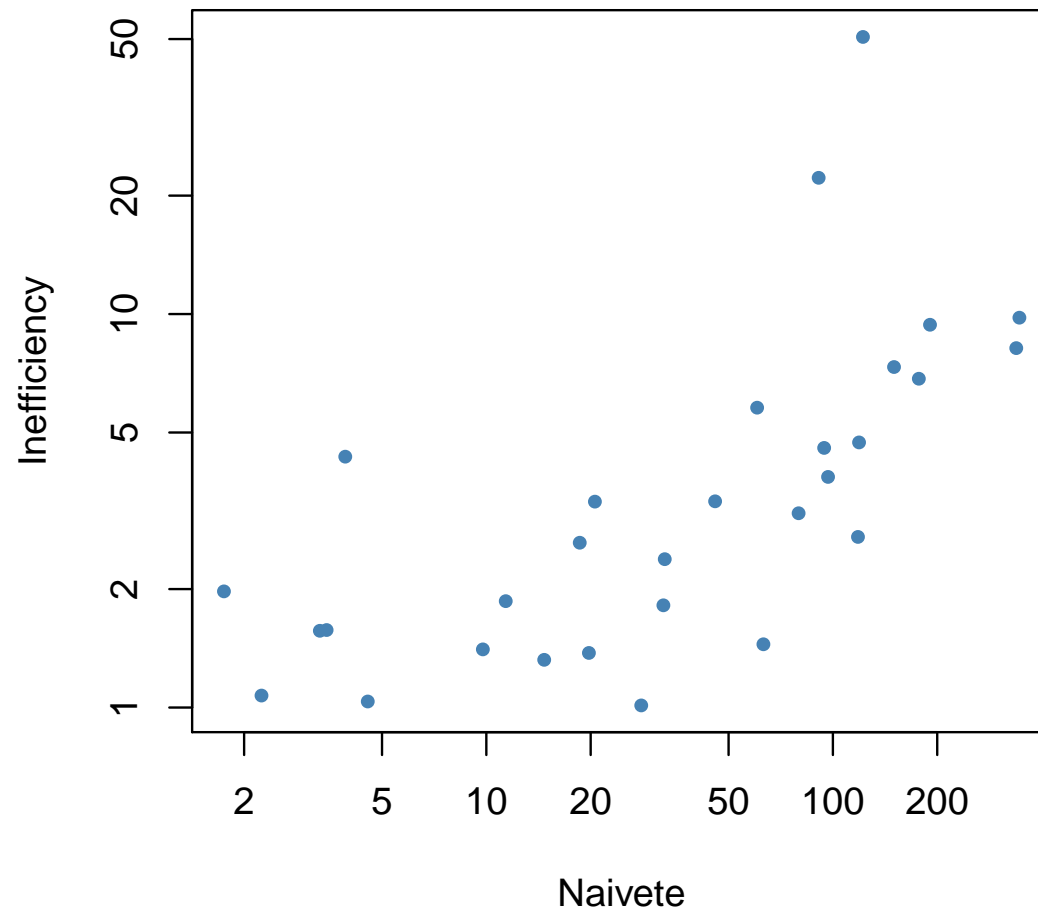
Naivete of OLS by coefficient



Naive by 100 \implies CIs 10x too narrow \implies Want 95% get 15.5%

Inefficiency vs naivete

Flaws of OLS by coefficient



Convergence

Backfitting update to \mathbf{b} :

$$\mathbf{b} \leftarrow M\mathbf{b} + \eta$$

If it converges, so does \mathbf{a} . Solution:

$$\mathbf{b} = \eta + \sum_{k=1}^{\infty} M^k \eta$$

For finite iterations

$$\|M\|_p \equiv \sup_{\mathbf{b} \neq 0} \frac{\|M\mathbf{b}\|_p}{\|\mathbf{b}\|_p} \leq 1 - \delta \quad \text{some } 1 \leq p \leq \infty \text{ and } \delta > 0$$

Keeps spectral radius below $1 - \delta$.

Matrix M

$$a_i \leftarrow \frac{\sum_s Z_{is}(Y_{is} - b_s)}{N_{i\bullet} + \lambda_A} \quad b_j \leftarrow \frac{\sum_r Z_{rj}(Y_{rj} - a_r)}{N_{\bullet j} + \lambda_B}$$

Leads to

For $1 \leq j, s \leq C$

$$b_j \leftarrow \eta_j + \sum_{s=1}^C M_{js} b_s$$

$$M_{js} = \frac{1}{N_{\bullet j} + \lambda_B} \sum_{r=1}^R \frac{Z_{rs}(Z_{rj} - N_{\bullet j}/R)}{N_{r\bullet} + \lambda_A}$$

Centering

Above is $M^{(0)}$.

Simple centering $M^{(1)}$.

Principled centering $M^{(2)}$.

Main results

Conditions on Z_{ij} for which

$$\mathbb{P}(\|M^{(2)}\|_1 < 1 - \delta) \rightarrow 1$$

as sampling increases. Similar for $M^{(1)}$.

For that we need

Model with $R, C \rightarrow \infty$

$N \ll RC$

Z_{ij} to control $\|M\|_1$

Why $\|\cdot\|_1$?

More tractable than $\|\cdot\|_2$ or spectral norm

The model

Problem size is $S \rightarrow \infty$

$$R = S^\rho, \quad C = S^\kappa$$

$$N \ll RC = S^{\rho+\kappa}$$

Sampling pattern

$$Z_{ij} \stackrel{\text{ind}}{\sim} \text{Bern}(p_{ij}) \quad 1 \leq i \leq R \quad 1 \leq j \leq C$$

$$\frac{S}{RC} \leq p_{ij} \leq \Upsilon \frac{S}{RC} \quad 1 \leq \Upsilon < \infty$$

The good

Unequal random $N_{i\bullet}$ and $N_{\bullet j}$

Unequal $\mathbb{E}(N_{i\bullet})$ and $\mathbb{E}(N_{\bullet j})$

The disappointing

Still very close to equal

Does not include tiny $N_{i\bullet}$ and $N_{\bullet j}$

$$S \leq \mathbb{E}(N) \leq \Upsilon S$$

Main result

Theorem 4.3. If

$$0 < \rho, \kappa < 1$$

$$1 \ll R, C \ll N$$

$$\rho + \kappa > 1$$

$$N \ll RC$$

$$2\rho + \kappa < 2$$

controls $N_{i\bullet}$

$$3\rho + 4\kappa < 4$$

gets column overlaps & controls $N_{\bullet j}$

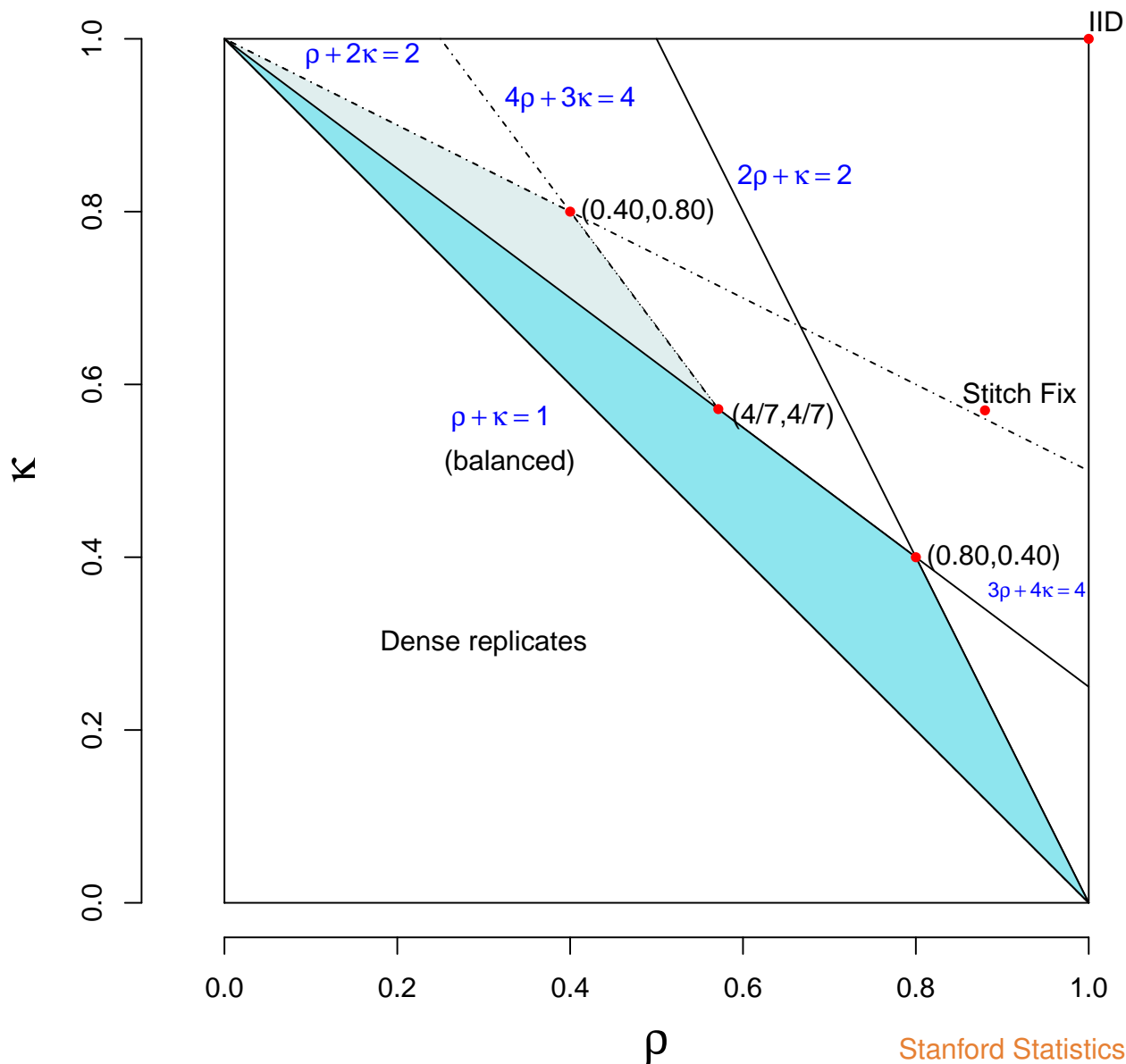
Then for any $\epsilon > 0$

$$\mathbb{P}(\|M^{(1 \text{ or } 2)}\|_1 \leq \Upsilon^2 - \Upsilon^{-2} + \epsilon) \rightarrow 1$$

Best Υ

$$\Upsilon < \sqrt{\frac{1 + \sqrt{5}}{2}} \doteq 1.27$$

Asymptotic domain



The proof

For $M \in \mathbb{R}^{C \times C}$

$$\|M\|_1 = \max_{1 \leq s \leq C} \sum_{j=1}^C |M_{js}|$$

For $j \neq s$

$$M_{js} = \frac{1}{N_{\bullet j} + \lambda_B} \sum_{r=1}^R \frac{Z_{rj}}{N_{r\bullet} + \lambda_A} (Z_{rs} - \bar{Z}_{\bullet s}) \quad \text{where}$$

$$\bar{Z}_{\bullet s} = \sum_{i=1}^R \frac{Z_{is}}{N_{i\bullet} + \lambda_A} / \sum_{i=1}^R \frac{1}{N_{i\bullet} + \lambda_A}$$

Upshot

It's a bit of a mess

Sketch

Hoeffding inequalities keep

$$(1 - \epsilon)S^{1-\rho} \leq N_{i\bullet} \leq (\Upsilon + \epsilon)S^{1-\rho}$$

$$(1 - \epsilon)S^{1-\kappa} \leq N_{\bullet j} \leq (\Upsilon + \epsilon)S^{1-\kappa}$$

and co-observations for $j \neq s$

$$(1 - \epsilon)S^{2-\rho-2\kappa} \leq \sum_{i=1}^R Z_{ij}Z_{is} \leq (\Upsilon^2 + \epsilon)S^{2-\rho-2\kappa}$$

and then

much interval arithmetic to propagate bounds to $\|M\|_1$

$$[a, A] + [b, B] \subseteq [a + b, A + B]$$

$$\text{For } a, b > 0 \quad [a, A]/[b, B] \subseteq [a/B, A/b]$$

$$|[a, A]| \subseteq [0, \max(|a|, |A|)]$$

Actual Stitch Fix norms

$$Z \in \{0, 1\}^{762,752 \times 6318} \quad M \in \mathbb{R}^{6318 \times 6318}$$

$\|M\|_1 < 1 - \delta$ sufficient but not necessary

$$\begin{pmatrix} \|M^{(0)}\|_1 & \|M^{(0)}\|_2 & |\lambda_{\max}(M^{(0)})| \\ \|M^{(1)}\|_1 & \|M^{(1)}\|_2 & |\lambda_{\max}(M^{(1)})| \\ \|M^{(2)}\|_1 & \|M^{(2)}\|_2 & |\lambda_{\max}(M^{(2)})| \end{pmatrix} = \begin{pmatrix} 31.9525 & 1.4051 & 0.6403 \\ 11.2191 & 0.4512 & 0.3338 \\ 8.9178 & 0.4541 & 0.3341 \end{pmatrix}$$

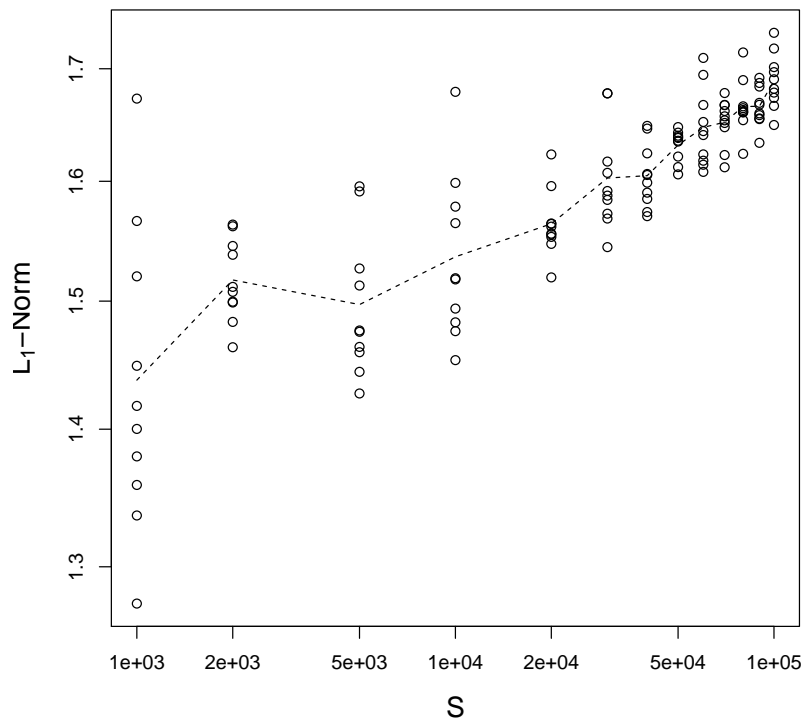
Iterations

6 iterations with threshold 10^{-8}

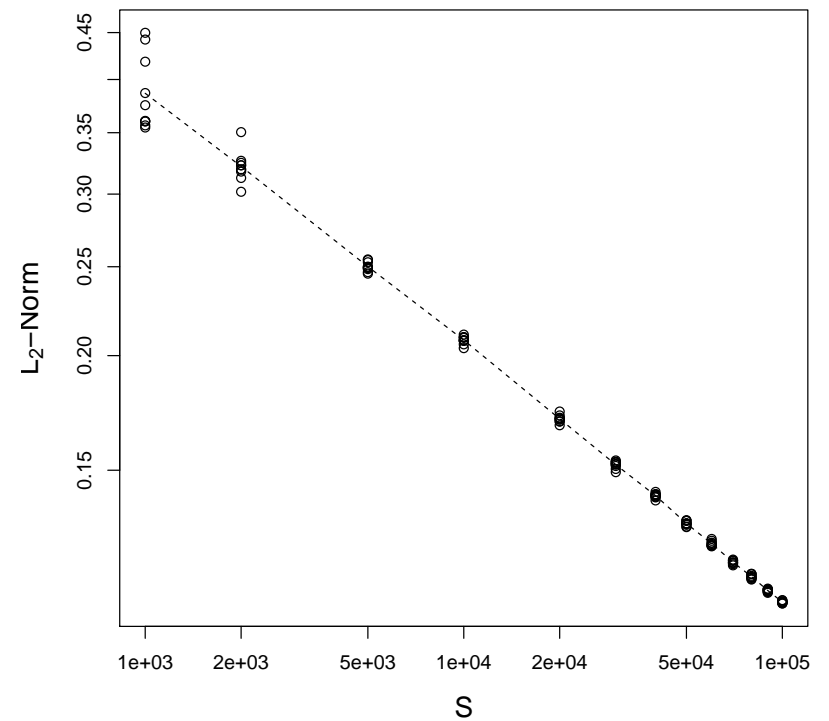
Simulated norms

$\rho = \kappa = 0.70$ outside our domain

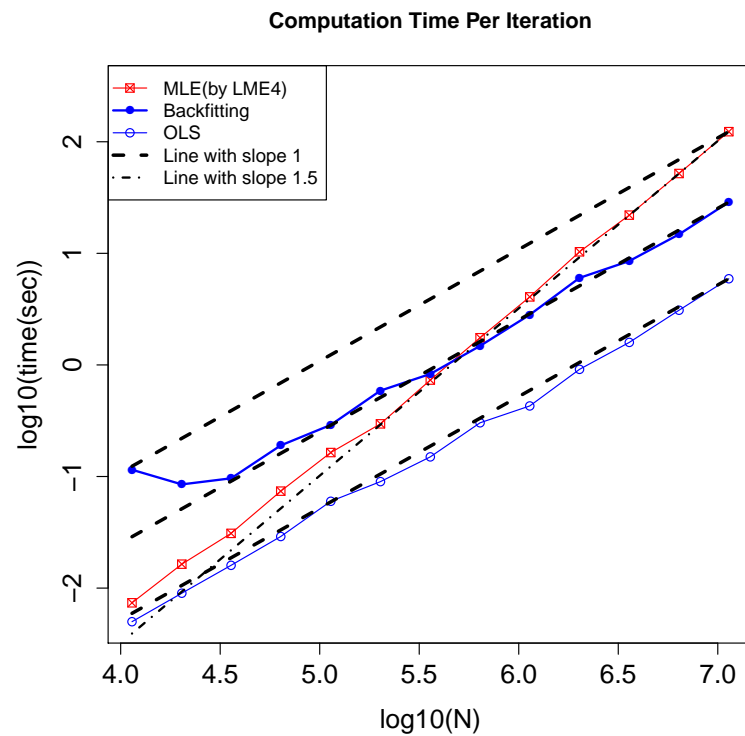
L_1 -Norm vs S for $\rho = 0.70, \kappa = 0.70$



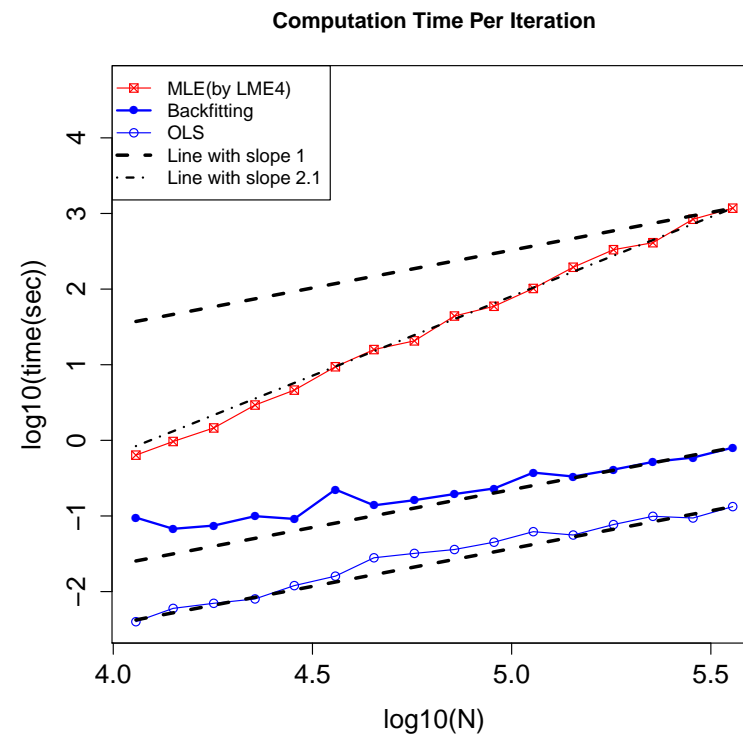
L_2 -Norm vs S for $\rho = 0.70, \kappa = 0.70$



Timings



(a) $(\rho, \kappa) = (0.52, 0.52)$



(b) $(\rho, \kappa) = (0.70, 0.70)$

$$R^3, C^3 = O(N^{2.1})$$

PRZ results

Mixing time

$$\rho_{\text{PRZ}} = \frac{N\sigma_A^2}{N\sigma_A^2 + R\sigma_E^2} \times \frac{N\sigma_B^2}{N\sigma_B^2 + C\sigma_E^2} \times \rho_{\text{AUX}}$$

$\rightarrow \rho_{\text{AUX}}$

Auxilliary ρ

Mixing time for Gibbs sampler on Z :

choose i given j with probability $Z_{ij}/N_{\bullet j}$

choose j given i with probability $Z_{ij}/N_{i\bullet}$

Further steps

- logistic regression (soon)
- higher way tables
- SVD-like interactions
- heteroscedastic effects

Thanks

- Co-authors Swarnadip Ghosh and Trevor Hastie
- Brad Klingenberg (Stitch Fix) data and discussions
- NSF IIS-1837931
- Giacomo Zanella, discussions
- Paromita Dubey, host
- Cindy Kirby, logistics

Backup slides

- 1) is OLS a reasonable comparison?
- 2) why Bayes and frequentist rates are often the same

OLS, really?

It seems like a straw man.

More likely to use a_i and b_j as fixed effects.

Cost

$$O(N(R + C + p)^2) = O(N^2) \quad \text{or worse}$$

So OLS on p variables is feasible.

OLS treatment of effects as fixed is not.

Hunch

people use learning algos that don't distinguish fixed vs random

Bayes & frequentist iterations

Sampling $\mathcal{N}(0, \Sigma)$:

e.g. draw one x_j at a time

convergence rate ρ_Σ

Minimizing $\mathbf{x}^\top Q \mathbf{x}$:

e.g. minimize over one x_j at a time

convergence rate ρ_Q

Very generally $\Sigma = Q \implies \rho_S = \rho_Q$

Colin Fox also Amit & Grenander

Conjectures

- this duality is why $N^{3/2}$ keeps popping up
- progress on Bayes \implies progress on minimization
- and vice versa

Ratio

Inefficiency

$$\frac{\widehat{\text{Var}}_{\text{GLS}}(\hat{\beta}_{\text{OLS},j})}{\widehat{\text{Var}}_{\text{GLS}}(\hat{\beta}_{\text{GLS},j})}$$

Naivete

$$\frac{\widehat{\text{Var}}_{\text{GLS}}(\hat{\beta}_{\text{OLS},j})}{\widehat{\text{Var}}_{\text{OLS}}(\hat{\beta}_{\text{OLS},j})}$$

Inefficiency / naivete

$$\frac{\widehat{\text{Var}}_{\text{OLS}}(\hat{\beta}_{\text{OLS},j})}{\widehat{\text{Var}}_{\text{GLS}}(\hat{\beta}_{\text{GLS},j})}$$

is the ratio of squared confidence interval widths