

# Regularization, optimization, and approximation

The benefits of a convex combination

---

Daniel J. McDonald

Indiana University, Bloomington

[pages.iu.edu/~dajmcdon](http://pages.iu.edu/~dajmcdon)

6 February 2019

## Obligatory “data are big” slide

Modern statistical applications — genomics, neural image analysis, text analysis, weather prediction — have large numbers of covariates  $p$

Also frequently have lots of observations  $n$

Need algorithms which can handle these kinds of data sets — with good statistical properties

### Computational choices impact statistical performance

These choices can take many forms:

- choosing tuning parameters
- different optimization algorithms return different solutions
- how long do we run our MCMC (and which kind do we use)

Statistical theory often neglects these choices:

- LASSO works with oracle tuning parameter
- We have the posterior if our MCMC runs forever
- EM gives us a global solution
- Algorithmic guarantees hold uniformly over ALL datasets

## Today's emphasis

Many statistical methods or optimization algorithms use a singular value decomposition (SVD):

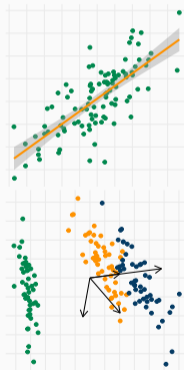
$$X = UDV^T : U^T U = I, V^T V = I, D \text{ diagonal.}$$

- Penalized Least Squares:

$$\min_{\beta} \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \text{Pen}(\beta)$$

- PCA:

$$\max_{V^T V = I_d} \text{tr}(V^T X^T X V)$$



The SVD is computationally expensive.

For a generic  $n \times p$  matrix, the SVD requires  $O(\min\{np^2, n^2p\})$  and storage of the entire matrix in fast memory.

I want to understand the statistical properties of some approximations that speed up computation and save storage.

1. Better algorithmic choices allow us to solve a statistical problem.
2. Algorithmic advances suggest new statistical techniques.
3. Statistical models suggest better algorithms.

1. Introduction
2. Climate science and choosing better algorithms
3. Approximations for many observations and robust statistics
4. Many measurements and better approximate algorithms
5. Ongoing/related/future work

## **Estimating the trend in cloud-top temperature volatility**

---

The scientific consensus is that

1. World-wide climate is changing.
2. This change is mostly driven by human behavior.

Global warming → climate change: the distribution of temperature (and precipitation) is changing

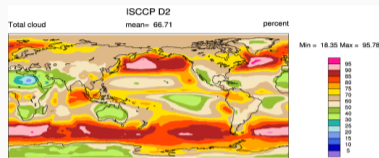
Increasing mean temperature understates the costs:

1. More frequent extremes have severe effects
2. Local discrepancies lead to more storms
3. Temporal dependencies mean persistence



## Drivers of climate variation:

1. Ocean currents
2. Jet stream
3. Annular modes + El Niño/La Niña
4. **Cloudiness**

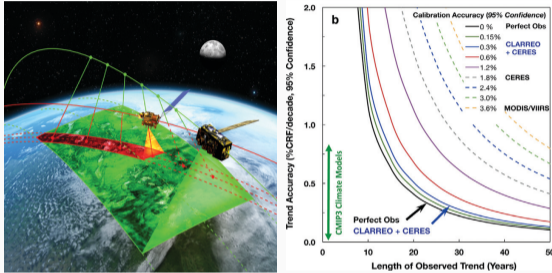


**CLARREO satellite:** monitor cloud top temperature as it relates to climate.

- Has yet to launch, no sooner than 2022
- Defunded in most recent federal budget

Source: NCAR CCSM3 Diagnostic Plots.

# CLARREO vs MetOp/Modis

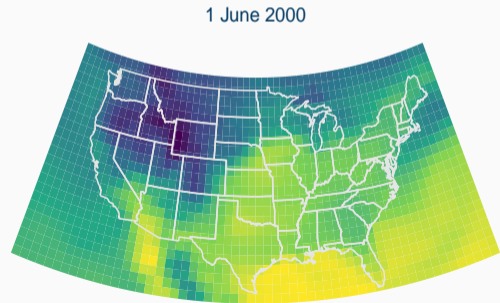
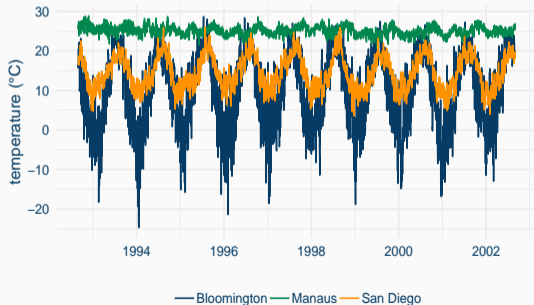


- Weather satellites aren't made for this.
- More information in higher moments than in average?

Source: Wielicki, et al. (2013).

Once collaborators do lots of processing...

- 52,000 time series
- daily records over ~ 40 years
- “trends” are local, nonlinear, not sinusoidal



Let  $y_{ts}$  be the observed temperature at time  $t$  and location  $s$ .

$$y_{ts} \sim N(0, \exp(h_{ts}))$$

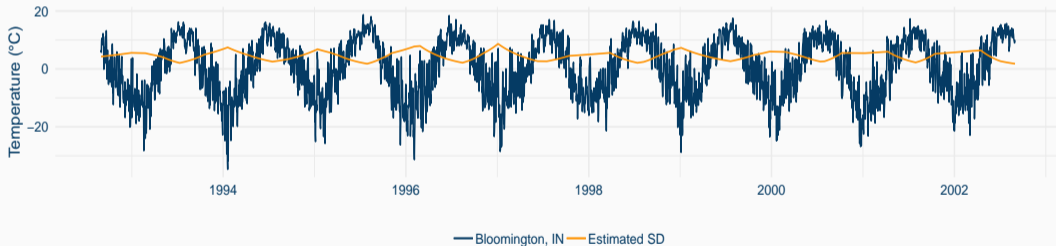
Estimate  $h$ , but it should be “smooth” relative to space and time.

Use a matrix  $D$  + penalty to encode this smoothness.

# Optimization problem

$$\min_h \sum h_{st} + y_{st}^2 e^{-h_{st}} + \lambda \|Dh\|_1$$

Standard optimizer: Primal Dual Interior Point method.



see Tibshirani (2014) or K-K-Boyd-G (2009)

1. Start with a guess  $h^{(1)}$
2. Solve a linear system  $[Au = v]$
3. Calculate a step size
4. Iterate 2 & 3 until convergence

The matrix  $A$  changes each iteration, dense, and roughly  $10^9 \times 10^9$ .

This isn't going to work.

$$\begin{array}{l} \text{Primal} \\ \min_h \quad f(h) + \lambda \|Dh\|_1 \end{array}$$

$$\begin{array}{l} \text{Dual} \\ \min_v \quad f^*(-D^\top v) \\ \text{s.t.} \quad \|v\|_\infty \leq \lambda \end{array}$$

- $f(h) := \sum h_{st} + y_{st}^2 e^{-h_{st}}$
- $f^*(u) := \sum (u_{st} - 1) \log \frac{y_{st}^2}{1 - u_{st}} + u_{st} - 1$

KKT conditions ( $w > 0$ )  $\implies$

$$r_w(v, \mu_1, \mu_2) := \begin{bmatrix} \nabla f^*(-D^\top v) + D(v - \lambda \mathbf{1})^\top \mu_1 - D(v + \lambda \mathbf{1})^\top \mu_2 \\ -\mu_1(v - \lambda \mathbf{1}) + \mu_2(v + \lambda \mathbf{1}) - w^{-1} \mathbf{1} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

- As  $w \rightarrow \infty$ , this converges to the optimum.
- But this is a nonlinear system, can't solve.
- Use Newton steps, which give the  $[Au = v]$  thing
- $A$  is the Jacobian of  $r_w$ .

- New algorithms: Khodadadi and M. (2019)
- Don't have to invert that matrix
- Must repeat for many tuning parameters
- Current work studies statistical properties:
- → Algorithm is “exact”, no statistical loss
- This is the best case





## **Solving big regression problems approximately**

---

## Core techniques

Suppose we have a matrix  $X \in \mathbb{R}^{n \times p}$  and vector  $Y \in \mathbb{R}^n$

Model:

$$Y = X\beta_* + \epsilon$$

**Least squares**  $\hat{\beta} = \operatorname{argmin}_{\beta} \frac{1}{2n} \|X\beta - Y\|_2^2$

**Ridge regression**  $\hat{\beta}_2(\lambda) = \operatorname{argmin}_{\beta} \frac{1}{2n} \|X\beta - Y\|_2^2 + \frac{\lambda}{2} \|\beta\|_2^2$

**LASSO**  $\hat{\beta}_1(\lambda) = \operatorname{argmin}_{\beta} \frac{1}{2n} \|X\beta - Y\|_2^2 + \lambda \|\beta\|_1$

**Dantzig selector**

**Group LASSO**

...

## Typical results for approximation

The general philosophy: Find an approximation  $\tilde{\beta}$  that is as close as possible to the solution of the original problem **UNIFORMLY**.

OLS:

$$\|X\tilde{\beta} - Y\|_2^2 \leq (1 + \epsilon)\|X\hat{\beta} - Y\|_2^2$$

Ridge:

$$\|X\tilde{\beta} - Y\|_2^2 + \frac{\lambda}{2}\|\tilde{\beta}\|_2^2 \leq (1 + \epsilon)\left(\|X\hat{\beta}_2(\lambda) - Y\|_2^2 + \frac{\lambda}{2}\|\hat{\beta}_2(\lambda)\|_2^2\right)$$

etc.

For an approximation  $\tilde{\theta}$  of  $\hat{\theta}$ ,

$$\begin{aligned}\mathbb{E}\|\theta - \tilde{\theta}\|_2^2 &= \mathbb{E}\|\tilde{\theta} - \hat{\theta} + \hat{\theta} - \theta\|_2^2 \\ &= \mathbb{E}\text{Approx error}^2 + \text{MSE}(\hat{\theta}) + 2\mathbb{E}\left[(\text{Approx error})(\hat{\theta} - \theta)\right] \\ &\leq \sup(\text{Approx error}^2) + \text{MSE}(\hat{\theta}) \\ &\quad + 2\sup(\text{Approx error}) \times \text{RMSE}(\hat{\theta})\end{aligned}$$

For an approximation  $\tilde{\theta}$  of  $\hat{\theta}$ ,

$$\begin{aligned}\mathbb{E}\|\theta - \tilde{\theta}\|_2^2 &= \mathbb{E}\|\tilde{\theta} - \hat{\theta} + \hat{\theta} - \theta\|_2^2 \\ &= \mathbb{E}\text{Approx error}^2 + \text{MSE}(\hat{\theta}) + 2\mathbb{E}\left[(\text{Approx error})(\hat{\theta} - \theta)\right] \\ &\leq \sup(\text{Approx error}^2) + \text{MSE}(\hat{\theta}) \\ &\quad + 2\sup(\text{Approx error}) \times \text{RMSE}(\hat{\theta}) \\ &\leq C\left(\sup(\text{Approx error}^2) + \text{MSE}(\hat{\theta})\right)\end{aligned}$$

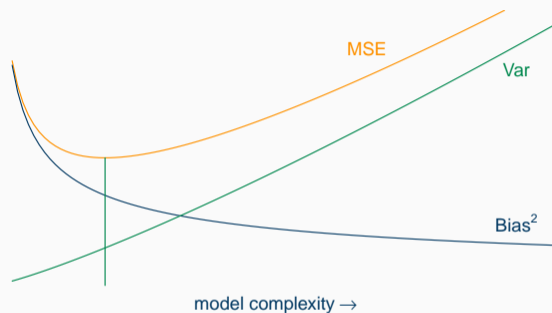
- Previous work examines Approx error, assumes  $\text{MSE}(\hat{\theta})$  small

For an approximation  $\tilde{\theta}$  of  $\hat{\theta}$ ,

$$\mathbb{E}\|\theta - \tilde{\theta}\|_2^2 = \text{Var}(\tilde{\theta}) + \text{Bias}(\tilde{\theta})^2$$

$$\begin{aligned}\mathbb{E}\|\theta - \tilde{\theta}\|_2^2 &= \mathbb{E}\|\tilde{\theta} - \hat{\theta} + \hat{\theta} - \theta\|_2^2 \\ &= \mathbb{E}\text{Approx error}^2 + \text{MSE}(\hat{\theta}) + 2\mathbb{E}\left[(\text{Approx error})(\hat{\theta} - \theta)\right] \\ &\leq \sup(\text{Approx error}^2) + \text{MSE}(\hat{\theta}) \\ &\quad + 2\sup(\text{Approx error}) \times \text{RMSE}(\hat{\theta}) \\ &\leq C\left(\sup(\text{Approx error}^2) + \text{MSE}(\hat{\theta})\right)\end{aligned}$$

- Previous work examines Approx error, assumes  $\text{MSE}(\hat{\theta})$  small
- We examine the MSE of the procedure



- We examine  $\mathbb{E}\|\theta - \tilde{\theta}\|_2^2$  where the expectation is over everything.
- Only other similar is [Ma, Mahoney, and Yu \(JMLR, 2015\)](#).

If  $X$  fits into RAM, there exist excellent algorithms in LAPACK that are

- Double precision
- Very stable
- $O(np^2)$  when  $n \gg p$ .
- $O(n^2p)$  when  $n \ll p$ .
- require extensive random access to matrix

There is a lot of interest in finding and analyzing techniques that extend these approaches to large( $r$ ) problems



## Out-of-core techniques for least squares

Many techniques focus on randomized compression

This is sometimes known as **sketching** or **preconditioning**

- Rokhlin, Tygert, (2008) "A fast randomized algorithm for overdetermined linear least-squares regression."
- Drineas, Mahoney, et al., (2011) "Faster least squares approximation."
- Woodruff, (2014) "Sketching as a Tool for Numerical Linear Algebra."
- Wang, Lee, Mahdavi, Kolar, Srebro, (2017) "Sketching meets random projection in the dual."
- Ma, Mahoney, and Yu, (2015), "A statistical perspective on algorithmic leveraging."
- Pilanci and Wainwright, (2015-2016). Multiple papers.
- Others.

## Basic Idea ( $n \gg p$ ):

- Choose some matrix  $Q \in \mathbb{R}^{q \times n}$ .
- Under many conditions, sufficient to choose  $q = \Omega(p)$ .
- Use  $QX$  (and)  $QY$  instead in the optimization.
- $O(np^2) \rightarrow O(p^3)$ .

## Basic Idea ( $n \gg p$ ):

- Choose some matrix  $Q \in \mathbb{R}^{q \times n}$ .
- Under many conditions, sufficient to choose  $q = \Omega(p)$ .
- Use  $QX$  (and)  $QY$  instead in the optimization.
- $O(np^2) \rightarrow O(p^3)$ .

## Basic Idea ( $n \ll p$ ):

- Choose some matrix  $Q \in \mathbb{R}^{p \times q}$ .
- Under many conditions, sufficient to choose  $q = \Omega(n)$ .
- Use  $XQ$  instead in the optimization.
- $O(n^2p) \rightarrow O(n^3)$ .

### Basic Idea ( $n \gg p$ ):

- Choose some matrix  $Q \in \mathbb{R}^{q \times n}$ .
- Under many conditions, sufficient to choose  $q = \Omega(p)$ .
- Use  $QX$  (and)  $QY$  instead in the optimization.
- $O(np^2) \rightarrow O(p^3)$ .

### Basic Idea ( $n \ll p$ ):

- Choose some matrix  $Q \in \mathbb{R}^{p \times q}$ .
- Under many conditions, sufficient to choose  $q = \Omega(n)$ .
- Use  $XQ$  instead in the optimization.
- $O(n^2p) \rightarrow O(n^3)$ .

Finding  $QX$  for arbitrary  $Q$  and  $X$  takes  $O(qnp)$  computations.

## Large $n$ smaller $p$ :

1. Introduce 2 new versions of compression.
2. General theoretical techniques.
3. Show how to choose tuning parameters (without extra computation).
4. Compression a bit worse under “nice” model. But **better** with outliers.

## Large $p$ smaller $n$ :

1. Procedure for low-rank (sparse) regression.
2. Statistical performance is near-optimal under the right model.
3. Algorithm is approximate, but doesn't lose anything.

**Lots of observations**

---

## Family of 4 (current SOTA)

### 1. Full compression:

$$\begin{aligned}\tilde{\beta}_{FC} &= \underset{\beta}{\operatorname{argmin}} \|Q(X\beta - Y)\|_2^2 \\ &= \underset{\beta}{\operatorname{argmin}} \|QX\beta\|_2^2 - 2Y^T Q^T QX\beta \\ &= (X^T Q^T QX)^{-1} X^T Q^T QY\end{aligned}$$

### 2. Partial compression:

$$\begin{aligned}\tilde{\beta}_{PC} &= \underset{\beta}{\operatorname{argmin}} \|QX\beta\|_2^2 - 2Y^T X\beta \\ &= (X^T Q^T QX)^{-1} X^T Y\end{aligned}$$

Also called "Hessian Sketching".

## Family of 4 (our versions)

Write:

$$B = [ \tilde{\beta}_{FC} \tilde{\beta}_{PC} ]$$

$$W = XB$$

3. Linear combination compression:

$$\hat{\alpha}_{lin} = \underset{\alpha}{\operatorname{argmin}} \|W\alpha - Y\|_2^2$$

$$\tilde{\beta}_{lin} = B\hat{\alpha}_{lin}$$

4. Convex combination compression:

$$\hat{\alpha}_{con} = \underset{\substack{0 \leq \alpha \\ \sum \alpha = 1}}{\operatorname{argmin}} \|W\alpha - Y\|_2^2$$

$$\tilde{\beta}_{con} = B\hat{\alpha}_{con}$$



## Why these?

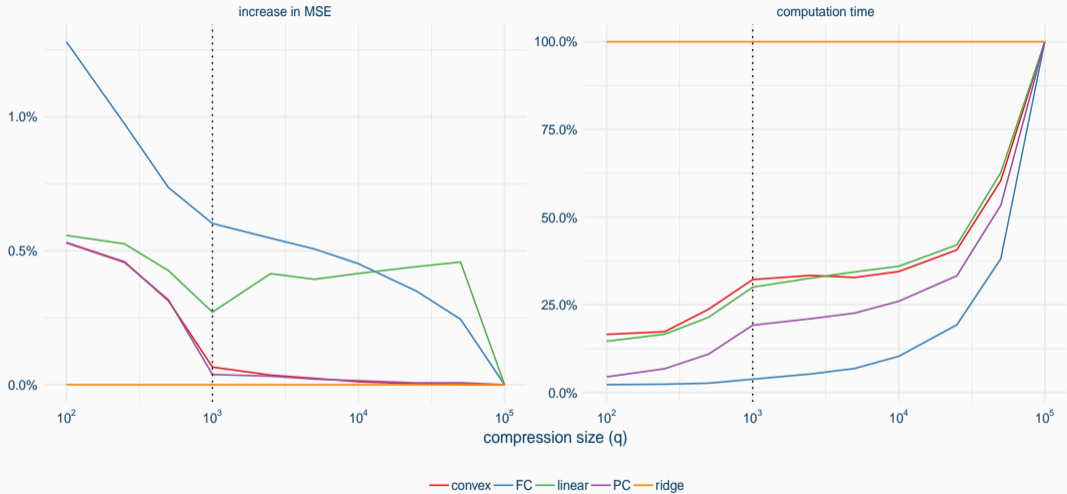
Turns out that FC is unbiased  $\implies$  worse than OLS (has high variance)

PC is biased and empirics demonstrate low variance

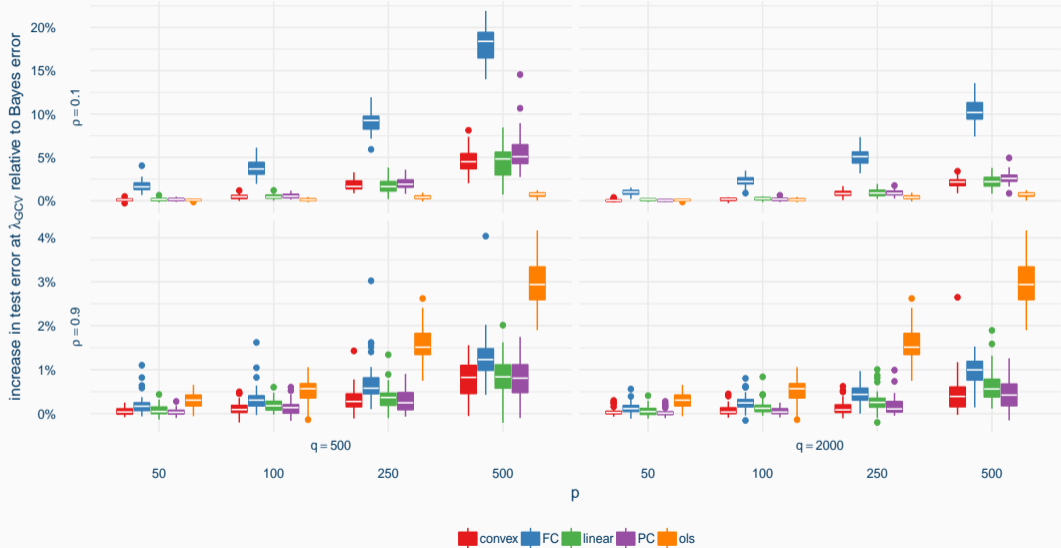
Combination should (and does) give better statistical properties

We do everything with an  $\ell_2$  penalty

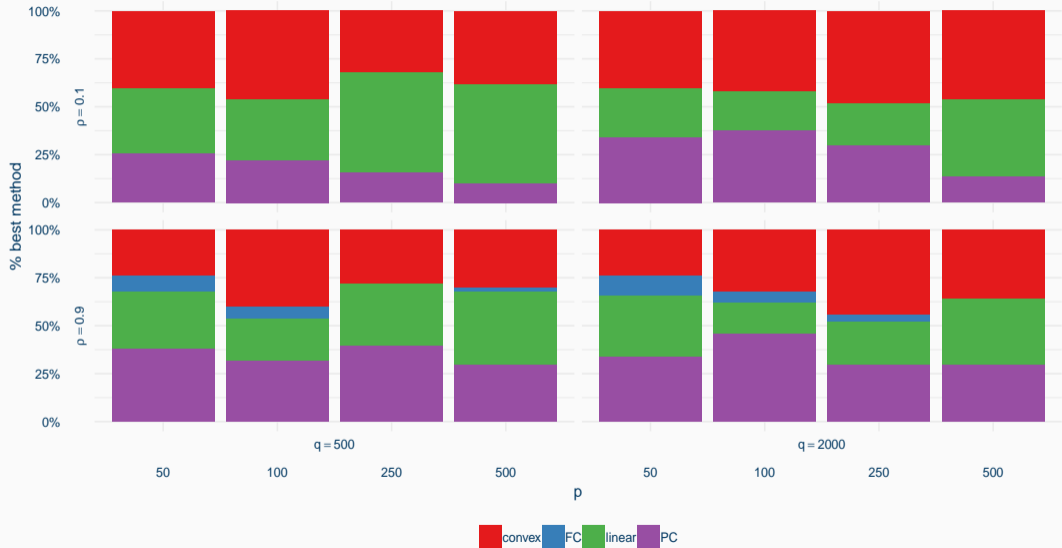
# Timing vs. accuracy



# Relative prediction error



# Which one wins?



- We use GCV with the degrees of freedom:

$$\text{GCV}(\lambda) = \frac{\frac{1}{n} \|X\tilde{\beta}(\lambda) - Y\|_2^2}{(1 - \text{df}/n)^2}$$

- df is easy for full or partial compression (though ignored in literature)
- For the other cases, an ad hoc approximation works, but has no justification.
- We derive an estimate via Stein's method.
- Easy to calculate both for a range of  $\lambda$  without extra computations.

- The **degrees of freedom** for a generic predictor  $f$  is

$$\text{df}(f) := \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(Y_i, f_i(Y)).$$

- For normal linear model (OLS), we have

$$\text{df} = \frac{1}{\sigma^2} \mathbb{E} \left[ \left( \widehat{Y} - \mathbb{E}[\widehat{Y}] \right)^\top (Y - \mu_Y) \right] = \frac{1}{\sigma^2} \mathbb{E} [Y^\top H Y] = p.$$

- In general, Stein's Lemma gives us the following: if  $Y_i - \mathbb{E}[Y_i|X_i] \sim N(0, \sigma^2)$

$$\text{Cov}(Y_i, f_i(Y)) = \mathbb{E} [(Y_i - \mathbb{E}[Y_i|X_i]) f_i(Y)] = \mathbb{E} [\nabla f_i(Y)].$$

Because linear and convex combination compression have closed form expressions:

e.g. **Linear combination:**

$$f(Y) = \begin{bmatrix} XB \end{bmatrix} \left[ (B^T X^T XB)^{-1} B^T X^T Y \right]$$

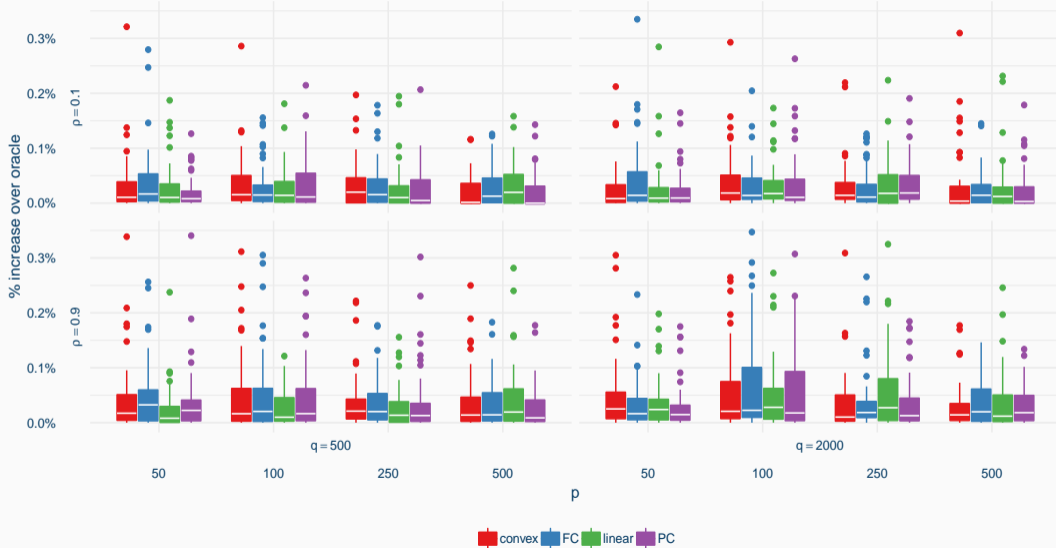
where

$$B = \begin{bmatrix} (X^T Q^T QX + \lambda I)^{-1} X^T \end{bmatrix} \begin{bmatrix} Q^T QY \\ Y \end{bmatrix}.$$

Apply every calculus rule you can find to yield a nasty expression which won't fit on this slide.

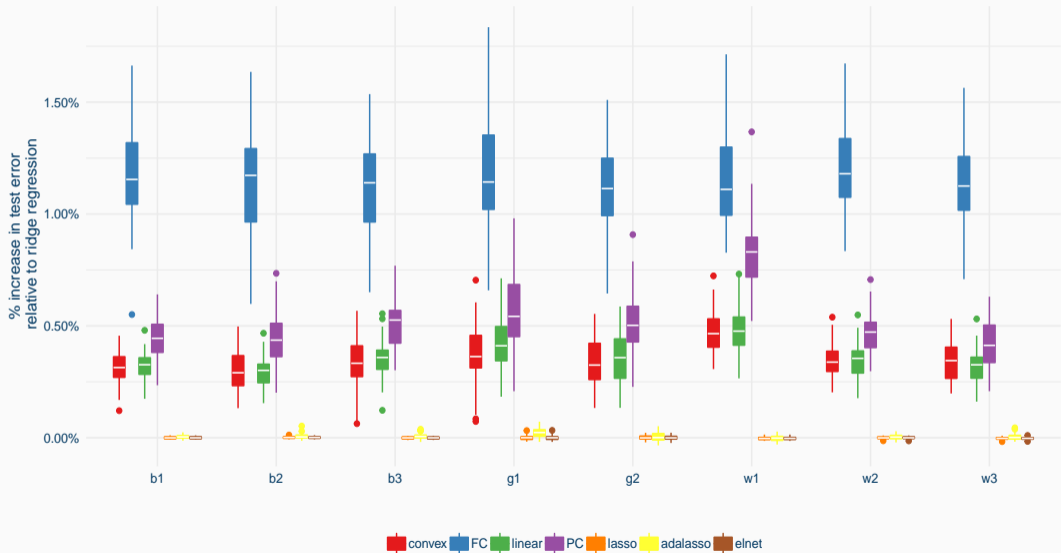
This is an unbiased estimate of df.

# Risk estimation accuracy





# Predicting read counts



### Theorem

$$\text{bias}^2 \left( \widehat{\beta}_2(\lambda) \mid X \right) = \lambda^2 \beta_*^\top V (D^2 + \lambda I_p)^{-2} V^\top \beta_*.$$
$$\text{tr} \left( \text{Var} \left[ \widehat{\beta}_2(\lambda) \mid X \right] \right) = \sigma^2 \sum_{i=1}^p \frac{d_i^2}{(d_i^2 + \lambda)^2}.$$

## Theorem

$$\begin{aligned}\text{bias}^2 [\tilde{\beta}_{FC} | X] &= \lambda^2 \beta_*^\top V (D^2 + \lambda I_p)^{-2} V^\top \beta_* + o_p(1) \\ \text{tr} (\text{Var} [\tilde{\beta}_{FC} | X]) &= \sigma^2 \sum_{i=1}^p \frac{d_i^2}{(d_i^2 + \lambda)^2} + o_p(1) \\ &\quad + \frac{(s-3)_+}{q} \text{tr} (\text{diag}(\text{vec}(I_n)) M^\top M \otimes (I-H) M \beta_* \beta_*^\top M^\top (I-H)) \\ &\quad + \frac{\beta_*^\top M^\top (I-H)^2 M \beta_*}{q} \text{tr}(M M^\top) \\ &\quad + \frac{1}{q} \text{tr} ((I-H) M \beta_* \beta_*^\top M^\top (I-H) M^\top M) .\end{aligned}$$

Note:  $M = (X^\top X + \lambda I_p)^{-1} X^\top$  and  $H = XM$

### Corollary

If  $\frac{1}{n}X^T X = I_p$ ,

$$\begin{aligned}\text{MSE}(\widehat{\beta}_2) &= b^2 \left( \frac{\theta}{1+\theta} \right)^2 + \frac{p\sigma^2}{n(1+\theta)^2} \\ \text{MSE}(\tilde{\beta}_{FC}) &= b^2 \left( \frac{\theta}{1+\theta} \right)^2 + \frac{p\sigma^2}{n(1+\theta)^2} + \frac{b^2 p \theta^2 (s-2)_+}{q(1+\theta)^4} + \frac{p^2 \theta^2 b^2}{q(1+\theta)^4} \\ \text{MSE}(\tilde{\beta}_{PC}) &= b^2 \left( \frac{\theta}{1+\theta} \right)^2 + \frac{p\sigma^2}{n(1+\theta)^2} + \frac{p(s-2)_+ b^2}{q(1+\theta)^2} + \frac{p b^2}{q(1+\theta)^4}\end{aligned}$$

where  $b^2 := \|\beta_*\|_2^2$ , and  $\theta := \lambda/n$

## What's the trick?

- All the estimators depend (at least) on

$$(X^T Q^T Q X + \lambda I_p)^{-1}$$

- We derived properties of  $Q^T Q$

$$\mathbb{E} \left[ \frac{s}{q} Q^T Q \right] = I_n$$
$$\text{Var} \left[ \text{vec} \left( \frac{s}{q} Q^T Q \right) \right] = \frac{(s-3)_+}{q} \text{diag}(\text{vec}(I_n)) + \frac{1}{q} I_{n^2} + \frac{1}{q} K_{nn}$$

- So the technique is to do a Taylor expansion around  $\frac{s}{q} Q^T Q = I_n$ .

## The takeaway message

Compression works in a fraction of the time.

Combining standard (FC and PC) is better.

Their genesis is an examination of the statistical performance.

Are there models under which compressed estimators are better?

## The takeaway message

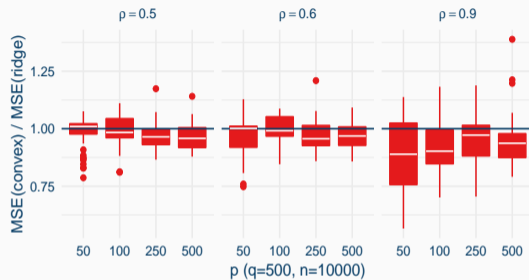
Compression works in a fraction of the time.

Combining standard (FC and PC) is better.

Their genesis is an examination of the statistical performance.

Are there models under which compressed estimators are better?

**Answer: YES! If there are outliers that are also high leverage.**



$$\widehat{\beta}_{n+1} - \widehat{\beta}_n = \left( \frac{y_{n+1} - x_{n+1}^\top \widehat{\beta}_n}{1 + x_{n+1}^\top (X^\top X)^{-1} x_{n+1}} \right) (X^\top X)^{-1} x_{n+1}$$



## Lots of measurements

---

“Random projection”

OLS:

$$\hat{\gamma} = \operatorname{argmin}_{\gamma} \frac{1}{2n} \|XQ\gamma - Y\|_2^2$$
$$\hat{\beta}_{RP} \leftarrow Q\hat{\gamma}$$

Before, sketched the covariance  $X^T X \rightarrow X^T Q^T Q X$ .

RP sketches the Gram matrix  $XX^T \rightarrow XQ^T QX^T$ .

So existing analyses are equivalent: as long as the sketching grabs most of the action.

## Grabbing the action

Before, we needed to get the high-leverage rows.

Now, we need “the high-leverage columns”.

Means most of the predictive information comes from a few variables.

Essentially, a low-rank predictor, driven by a few variables.

This is just sparse PCA regression.

## Grabbing the action

Before, we needed to get the high-leverage rows.

Now, we need “the high-leverage columns”.

Means most of the predictive information comes from a few variables.

Essentially, a low-rank predictor, driven by a few variables.

This is just sparse PCA regression.

We do that instead.

Two stages:

1.  $\widehat{V} = \operatorname{argmax}_{V^T V = I_d} \operatorname{tr}(V^T X^T X V)$
2.  $\widehat{\gamma}_{pcr} = \operatorname{argmin}_{\gamma \in \mathbb{R}^d} \frac{1}{2n} \|Y - X \widehat{V} \gamma\|_2^2$   
 $\widehat{\beta}_{pcr} \leftarrow \widehat{V} \widehat{\gamma}_{pcr}$

Here,  $\widehat{V}$  is analogous to  $Q$  but it uses the structure of  $X$ .

But 1. is hard.

**Algorithmically:** Need the SVD of a  $n \times p$  matrix  $\rightarrow O(np^2)$ .

**Statistically:**  $p \gg n \implies \widehat{V}$  is inconsistent for the population analogue.

**"Supervised" PCA** First screen away most of the variables using  $Y$ .

Solves both problems if  $\mathbb{E} [X_j Y] = 0 \Rightarrow \beta_{*,j} = 0$ .

See work of Bair+Paul+Hastie (2004, 2006, 2008) or Tay, Friedman, Tibshirani (2018)

**"Sparse" PCA** Solve a constrained version of 1.

Good statistical properties for 1. but ignores  $Y$

See d'Aspremont et al. (2007) Johnstone and Lu (2009) Zhang and Ghaoui (2011), among others

Two stages:

1.  $\widehat{V} = \underset{\substack{0 \leq W^T \leq I \\ \text{tr}(V)=d}}{\text{argmax}} \text{tr}(V^T X^T X V) - \lambda \|V V^T\|_{1,1}$
2.  $\widehat{\gamma}_{FR} = \underset{\gamma \in \mathbb{R}^d}{\text{argmin}} \frac{1}{2n} \|Y - X \widehat{V} \gamma\|_2^2$   
 $\widehat{\beta}_{FR} \leftarrow \widehat{V} \widehat{\gamma}_{FR}$

$\|V V^T\|_{1,1}$  forces rows of  $V$  to be 0

$V_j = 0 \Rightarrow V_j \gamma = 0 \Rightarrow X_j$  doesn't predict  $Y$ .

Vu et al. (2013) call 1. “Fantope projection”, nearly minimax optimal

Bad algorithmic properties:  $O(\# \text{ iterations} \times p^3)$

## Alternating direction method of multipliers

Restate your optimization

Original	Equivalent
$\min_x f(x) + g(x)$	$\min_{x,z} f(x) + g(z)$
	s.t. $x - z = 0$

Then, iterate the following with  $\rho > 0$

$$x \leftarrow \operatorname{argmin}_x f(x) + \frac{\rho}{2} \|x - z + u\|_2^2$$

$$z \leftarrow \operatorname{argmin}_z g(z) + \frac{\rho}{2} \|x - z + u\|_2^2$$

$$u \leftarrow u + x - z$$



## Why would you do this?

- It decouples  $f$  and  $g$
- If  $f$  and  $g$  have the right structure  $\implies$  parallelizable
- There are often many ways to decouple a problem
- The individual minimizations don't have to be solved in closed form

ADMM for Fantope projection:

$$A \leftarrow \Pi_{\mathcal{F}^d} \left( V - U + \frac{1}{n\rho} X^T X \right)$$

$$V \leftarrow \mathcal{S}_{\lambda/\rho}(A + U)$$

$$U \leftarrow U + A - V$$

$$[\mathcal{S}_a(b)]_k = \text{sgn}(b_k)(|b_k| - a)_+$$

Given an eigen decomposition of  $Z = \sum_i \gamma_i z_i z_i^T$ .

$$\Pi_{\mathcal{F}^d}(Z) = \sum_i \gamma_i^+(\theta) z_i z_i^T$$

$$\gamma_i^+(\theta) = \min(\max(\gamma_i - \theta, 0), 1),$$

$$\theta \text{ s.t. } \sum_i \gamma_i^+(\theta) = d$$

- The  $\gamma$ - $\theta$  stuff solves a monotone, piecewise linear equation.
- But we have to do the decomposition at every iteration.

## Conditions for convergence of ADMM

- When the updates are exact, all you need for convergence is
  1.  $f, g$  are convex, extended real valued.
  2.  $f(x) + g(z) + u^T(x - z)$  has a saddle point.
- The convergence rate is not well understood (seems linear).
- It turns out, you can solve the minimizations approximately.

$$\sum_{k=1}^{\infty} \|\Pi(y^k) - \tilde{\Pi}(y^k)\|_2 < \infty$$

- We do that using earlier work (Homrighausen and McDonald, 2016)

# Results in simulations



Assume many conditions,  $s := |\beta_*|$ ,  $\text{supp}(v) := \{j : v_j \neq 0\}$ ,

### Theorem

$$\|\widehat{\beta}_{FR} - \beta_*\|_2 = o_P\left(\sigma\sqrt{\frac{(s^2 + d)\log p}{n}}\right),$$

and

$$\left|\text{supp}(\widehat{\beta}_{FR}) \Delta \text{supp}(\beta_*)\right| = o_P\left(\sigma\sqrt{\frac{s^2 \log p}{n}}\right).$$

## The big picture

---

When you have large datasets and complicated estimators

1. Algorithms can avoid hard computations, enable inference.
2. Understanding how algorithms work can motivate statistical advances.
3. Understanding statistics may motivate new algorithms.

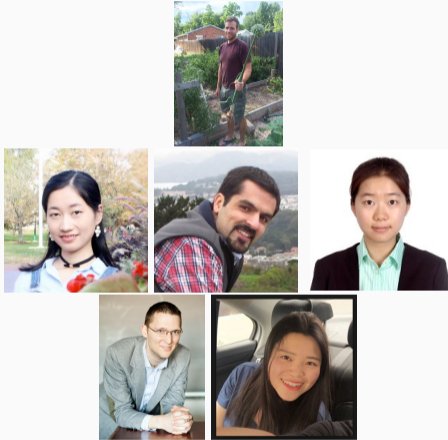
- Approximation for Least Squares ( $n \gg p$ )
  - Homrighausen and McDonald. (2019+). Under review.
- Approximation for dimension reduction ( $p \gg n$ )
  - Homrighausen and McDonald. (2016). JCGS.
  - Ding and McDonald. (2017). Bioinformatics.
  - Ding and McDonald. (2019). Under review.
- Algorithms for large data
  - ADMM for large constrained kernel PCA.
  - McDonald and Khodadadi. (2019). AAAI.
  - Trend filtering for Spatio-temporal exponential families.



- Approximation for Least Squares ( $n \gg p$ )
  - Homrighausen and McDonald. (2019+). Under review.
- Approximation for dimension reduction ( $p \gg n$ )
  - Homrighausen and McDonald. (2016). JCGS.
  - Ding and McDonald. (2017). Bioinformatics.
  - Ding and McDonald. (2019). Under review.
- Algorithms for large data
  - ADMM for large constrained kernel PCA.
  - McDonald and Khodadadi. (2019). AAAI.
  - Trend filtering for Spatio-temporal exponential families.

- GCV and SURE for compressed regression
- SURE for exponential families.
- CV and  $\ell_1$  regularization
  - Homrighausen and McDonald. (2013). ICML.
  - Homrighausen and McDonald. (2014). Machine Learning.
  - Homrighausen and McDonald. (2017). Stat. Sinica.
- Dependence and high dimensions
  - Homrighausen and McDonald. (2018). JSCS.
  - McDonald and Shalizi. (2018+). Under review.
  - McDonald, Shalizi and Schervish. (2017). JMLR.

# Collaborators and funding



Institute for  
**New Economic  
Thinking**

# Appendix

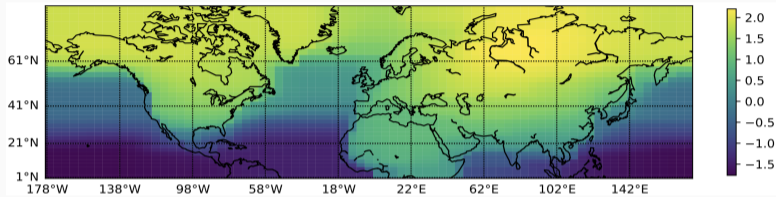
---

## Algorithm 1 Linearized ADMM

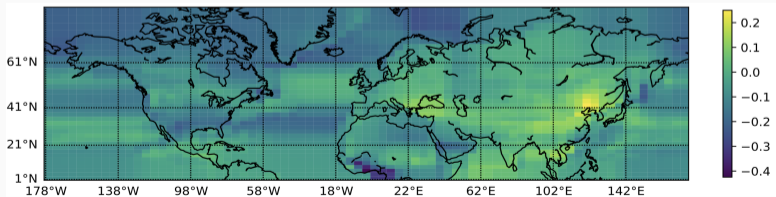
---

- 1: **Input** data  $y$ , penalty matrix  $D$ ,  $\epsilon, \rho, \lambda_t, \lambda_s > 0$ .
  - 2: **Set**  $h \leftarrow 0, z \leftarrow 0, u \leftarrow 0$ . ▷ Initialization
  - 3: **repeat**
  - 4:    $h_k \leftarrow \mathcal{W}\left(\frac{y_k^2}{\mu} \exp\left(\frac{1-\mu u_k}{\mu}\right)\right) + \frac{1-\mu u_k}{\mu}$  for all  $k = 1, \dots, TS$ . ▷ Primal update
  - 5:    $z \leftarrow S_{\rho\lambda}(u)$ . ▷ Elementwise soft thresholding
  - 6:    $u \leftarrow u - z$ . ▷ Dual update
  - 7: **until**  $\max\{\|Dh - z\|, \|z^{m+1} - z^m\|\} < \epsilon$
  - 8: **Return**  $z$ .
-

## Average variance



## Change in average variance from 1961–2011



**Gaussian** Well behaved distribution and eas(ier) theory. Dense matrix

### Fast Johnson-Lindenstrauss Methods

**Randomized Hadamard** (or Fourier) transformation. Allows for  $O(np \log(p))$  computations.

**Subsampling**  $Q = \pi\tau$  for  $\pi$  a permutation of  $I$  and  $\tau = [I_q \ 0]$ .  $QX$  means “read  $q$  (random) rows”

### Sparse Bernoulli

$$Q_{ij} \stackrel{i.i.d.}{\sim} \begin{cases} 1 & \text{with probability } 1/(2s) \\ 0 & \text{with probability } 1 - 1/s \\ -1 & \text{with probability } 1/(2s) \end{cases}$$

This means  $QX$  takes  $O\left(\frac{qnp}{s}\right)$  “computations” on average.

## Why combine?

- $Y = X\hat{\beta} + \hat{e}$  with  $P(\hat{e} \in \text{col}(X)) = 0$ , and  $\mathbb{E}[\hat{e}] = 0$ .

$$\begin{aligned}\hat{\beta}_{FC}(0) &= (X^T Q^T Q X)^{\dagger} X^T Q^T Q Y = (X^T Q^T Q X)^{\dagger} X^T Q^T Q (X\hat{\beta} + \hat{e}) \\ &= \hat{\beta} + (X^T Q^T Q X)^{\dagger} X^T Q^T Q \hat{e}\end{aligned}$$

$$\Rightarrow \mathbb{E}[\hat{\beta}_{FC}(0)] = \mathbb{E}[\hat{\beta}] = \beta_*$$

$$\begin{aligned}\hat{\beta}_{PC}(0) &= (X^T Q^T Q X)^{\dagger} X^T Y = (X^T Q^T Q X)^{\dagger} X^T (X\hat{\beta} + \hat{e}) \\ &= (X^T Q^T Q X)^{\dagger} X^T X \hat{\beta}\end{aligned}$$

$$\Rightarrow \mathbb{E}[\hat{\beta}_{PC}(0)] = (X^T Q^T Q X)^{\dagger} X^T X \beta_*$$



## Simulation setup ( $n \gg p$ )

- Draw  $X_i \sim \text{MVN}(\mathbf{0}, (1 - \rho)I_p + \rho\mathbf{1}\mathbf{1}^\top)$
- Draw  $\beta \sim \text{N}(\mathbf{0}, \tau^2 I_p)$
- Draw  $Y_i = X_i^\top \beta_* + \epsilon_i$  with  $\epsilon_i \sim \text{N}(\mathbf{0}, \sigma^2)$ .
- For this model, the optimal estimator (in MSE) is

$$\widehat{\beta}_B = (X^\top X + \lambda_* I_p)^{-1} X^\top Y$$

with  $\lambda_* = \frac{\sigma^2}{n\tau^2}$

Define the linear combination prediction as  $\widehat{Y}_{LC} = XB(\lambda)\widehat{\alpha}_{LC} = XB(\lambda)(B^T(\lambda)X^T XB(\lambda))^\dagger B^T(\lambda)X^T Y$  and let  $Z := XB(\lambda)$  so that  $Z(Z^T Z)^\dagger Z^T Y =: P_Z Y$ .

**Theorem:**

$$\begin{aligned} \text{div}_{lin}(\lambda) = & \widehat{\alpha}_{FC} df_{FC} + \widehat{\alpha}_{PC} df_{PC} + \text{tr}(P_Z) - \widehat{\alpha}_{PC} \text{tr}(P_Z H_{PC}) - \widehat{\alpha}_{FC} \text{tr}(P_Z H_{FC}) \\ & - \text{tr} \left( Z^\dagger [H_{PC} \widehat{Y}_{LC} H_{FC} \widehat{Y}_{LC}] \right). \end{aligned}$$

## Theorem:

$$\begin{aligned} \operatorname{div}_{\text{con}}(\lambda) &= \widehat{\alpha}_{PC} df_{PC} + (1 - \widehat{\alpha}_{PC}) df_{FC} \\ &+ \frac{(Y - \widehat{Y}_{FC})^\top H_{PC} (\widehat{Y}_{PC} - \widehat{Y}_{FC})}{\left(\widehat{Y}_{PC} - \widehat{Y}_{FC}\right)^\top \left(\widehat{Y}_{PC} - \widehat{Y}_{FC}\right)} + \frac{\widehat{Y}_{PC}^\top (I_n - H_{FC}) (\widehat{Y}_{PC} - \widehat{Y}_{FC})}{\left(\widehat{Y}_{PC} - \widehat{Y}_{FC}\right)^\top \left(\widehat{Y}_{PC} - \widehat{Y}_{FC}\right)} \\ &- \frac{2 \left(\widehat{Y}_{PC} - \widehat{Y}_{FC}\right)^\top (Y - \widehat{Y}_{FC})}{\left(\left(\widehat{Y}_{PC} - \widehat{Y}_{FC}\right)^\top \left(\widehat{Y}_{PC} - \widehat{Y}_{FC}\right)\right)^2} \left(\widehat{Y}_{PC} - \widehat{Y}_{FC}\right)^\top (H_{PC} - H_{FC}) \left(\widehat{Y}_{PC} - \widehat{Y}_{FC}\right). \end{aligned}$$

## Flop count for different algorithms

Approach	$O(\cdot)$
Linear system	$np^2 + p^3$
Low-rank linear system	$npr + r^3$
Gradient descent	$n^{3/2}p^2 \log \epsilon^{-1}$
Acc. gradient descent	$n^{5/4}p^{3/2} \log \epsilon^{-1}$
Coordinate descent	$n^{3/2}p \log \epsilon^{-1}$
SVRG, SDCA, SAG	$(np + n^{1/2}p^2) \log \epsilon^{-1}$
Catalyst, APPA	$(np + n^{3/4}p^{3/2}) \log \epsilon^{-1}$
DSPDC	$npr + (nr + n^{3/4}p^{3/2}r) \log \epsilon^{-1}$
Iterative Hessian sketch	$np \log p + n^{1/4}p^{3/2}r \log^2 \epsilon^{-1}$
Dual random projection	$np \log n + (nr^2 + r^3) \log \epsilon^{-1}$

## Genes and the SPC model

sparsity of $\Sigma_{XX}^{-1}$	1.0000	0.9999	0.9998	0.9995	0.9991	0.9984	0.9975	0.9963	0.9946	0.9922
% non-zero $\beta_*$ 's	0.0162	0.0216	0.0287	0.0418	0.0618	0.0843	0.1193	0.1803	0.2645	0.3699
False Negative Rate	0.0000	0.2500	0.4340	0.6117	0.7374	0.8077	0.8641	0.9100	0.9387	0.9562

