# COMPRESSED AND PENALIZED LINEAR REGRESSION

### Daniel J. McDonald
### Indiana University, Bloomington
mypage.iu.edu/ dajmcdon

2 June 2017

# OBLIGATORY "DATA IS 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

# LESSON OF THE TALK

Many statistical methods use (perhaps implicitly) a singular value decomposition (SVD) to solve an optimization problem.

The SVD is computationally expensive.

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

Spoiler: sometimes approximations actually improve the statistical properties

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

## LEAST SQUARES:

$$\widehat{\beta} = \underset{\beta}{\mathrm{argmin}} \, ||\mathbb{X}\beta - Y||_2^2$$

## $\ell_2$ REGULARIZATION:

$$\widehat{\beta}(\lambda) = \underset{\beta}{\mathrm{argmin}} \, ||\mathbb{X}\beta - Y||_2^2 + \lambda \, ||\beta||_2^2$$

# CORE TECHNIQUES

If $\mathbb{X}$ fits into RAM, there exist excellent algorithms in LAPACK that are

- Double precision
- Very stable
- cubic complexity, $O(np^2 + p^3)$, with small constants
- 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

If $\mathbb{X}$ is too large to manipulate in RAM, there are other approaches

- (Stochastic) gradient descent
- Conjugate gradient
- Rank-one QR updates
- Krylov subspace methods

# OUT-OF-CORE TECHNIQUES

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, (2016) "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:

- Choose some matrix $Q \in \mathbb{R}^{q \times n}$ .
- Use $Q\mathbb{X}$ (and) $QY$ instead in the optimization.

Finding $Q\mathbb{X}$ for arbitrary $Q$ and $\mathbb{X}$ takes $O(qnp)$ computations.

This can be expensive.

To get this approach to work, we need some structure on $Q$

# THE $Q$ MATRIX

- 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.
- $Q = \pi\tau$ for $\pi$ a permutation of $I$ and $\tau = [I_q\ 0]$:
  $Q\mathbb{X}$ means "read $q$ (random) rows"
- Sparse Bernoulli:

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

This means $Q\mathbb{X}$ takes $O\left(\frac{qnp}{s}\right)$ "computations" on average.

The general philosophy: Find an approximation algorithm that is as close as possible to the solution of the original problem.

For OLS, typical results would be to produce an $\tilde{\beta}$ such that

$$\left|\left|\mathbb{X}\tilde{\beta} - Y\right|\right|_2^2 \le (1 + \epsilon) \left|\left|\mathbb{X}\beta_* - Y\right|\right|_2^2,$$
$$\left|\left|\mathbb{X}(\tilde{\beta} - \beta_*)\right|\right|_2^2 \le \epsilon \left|\left|\mathbb{X}\beta_*\right|\right|_2^2,$$
$$\left|\left|\tilde{\beta} - \beta_*\right|\right|_2^2 \le \epsilon \left|\left|\beta_*\right|\right|_2^2,$$

Here, $\tilde{\beta}$ should be 'easier' to compute than

$$\beta_* = \underset{\beta}{\operatorname{argmin}} \left|\left|\mathbb{X}\beta - Y\right|\right|_2^2$$

# COLLABORATOR & GRANT SUPPORT

Collaborator:
Darren Homrighausen
Department of Statistics
Southern Methodist University

Grant Support:
NSF, Institute for New Economic
Thinking

An alternate perspective

# RISK

Form a loss function $\ell : \Theta \times \Theta \to \mathbb{R}^+$

The quality of an estimator is given by its risk
$$R(\widehat{\theta}) = \mathbb{E}\left[\ell(\widehat{\theta}, \theta)\right]$$

We could use $\ell_2$ estimation risk:
$$R(\widehat{\theta}) = \mathbb{E}\left|\left|\theta - \widehat{\theta}\right|\right|_2^2$$

or excess $\ell_2$ prediction risk

$$R(\widehat{\theta}) = \mathbb{E}\left|\left|Y - \mathbb{X}\widehat{\theta}\right|\right|_2^2 = \mathbb{E}\left|\left|Y - \mathbb{X}\theta + \mathbb{X}\theta - \mathbb{X}\widehat{\theta}\right|\right|_2^2$$
$$\propto \text{constant} + \mathbb{E}\left|\left|\mathbb{X}(\theta - \widehat{\theta})\right|\right|_2^2$$

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

$$\mathbb{E}\left\|\theta - \tilde{\theta}\right\|_2^2 = \mathbb{E}\left\|\tilde{\theta} - \widehat{\theta} + \widehat{\theta} - \mathbb{E}\widehat{\theta} + \mathbb{E}\widehat{\theta} - \theta\right\|_2^2$$
$$\leq \mathbb{E}\text{Approx. error}^2 + \text{Variance} + \text{Bias}^2$$

- Previous analyses focus only on the approximation error (with expectation over the algorithm, not the data).

# BIAS-VARIANCE TRADEOFF



model complexity $\rightarrow$

- Typical result compares to the zero-bias estimator which is assumed to have small variance.
- We examine $\mathbb{E} \left\lVert \theta - \tilde{\theta} \right\rVert_2^2$ where the expectation is over everything random.
- Closest similar analysis is Ma, Mahoney, and Yu (JMLR, 2015).

# COMPRESSED REGRESSION

Let $Q \in \mathbb{R}^{q \times n}$

Call the fully compressed least squares estimator

$$\widehat{\beta}_{FC} = \underset{\beta}{\operatorname{argmin}} \, ||Q(\mathbb{X}\beta - Y)||_2^2$$

$\rightarrow$ A common way to solve least squares problems that are:

- Very large or
- Poorly conditioned

The numerical/theoretical properties generally depend on $Q$, $q$, $\mathbb{X}$

1. Full compression:

$$\hat{\beta}_{FC} = \underset{\beta}{\mathrm{argmin}} \, ||Q(\mathbb{X}\beta - Y)||_2^2$$
$$= \underset{\beta}{\mathrm{argmin}} \, ||QY||_2^2 + ||Q\mathbb{X}\beta||_2^2 - 2Y^\top Q^\top Q\mathbb{X}\beta$$
$$= (\mathbb{X}^\top Q^\top Q\mathbb{X})^{-1}\mathbb{X}^\top Q^\top QY$$

2. Partial compression:[1]

$$\hat{\beta}_{PC} = \underset{\beta}{\mathrm{argmin}} \, ||Y||_2^2 + ||Q\mathbb{X}\beta||_2^2 - 2Y^\top \mathbb{X}\beta$$
$$= (\mathbb{X}^\top Q^\top Q\mathbb{X})^{-1}\mathbb{X}^\top Y$$

[1] Also called "Hessian Sketching".

## WE ALSO COMBINE THESE

Write:

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

$$W = \mathbb{X}B$$

3. Linear combination compression:

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

$$\widehat{\beta}_{lin} = B\widehat{\alpha}_{lin}$$

4. Convex combination compression:

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

$$\widehat{\beta}_{con} = B\widehat{\alpha}_{con}$$

- These are simple to calculate given FC and PC.

# WHY THESE?

- Turns out that FC is (approximately) unbiased, and therefore worse than OLS (has high variance)
- On the other hand, PC is biased and empirics demonstrate low variance
- Combination should give better statistical properties
- We do everything with an $\ell_2$ penalty

# Evidence from simulations

# SIMULATION SETUP

- Draw $\mathbb{X}_i \sim \text{MVN}(0, (1 - \rho)I_p + \rho \mathbf{1}\mathbf{1}^\top)$
  - We use $\rho = \{0.2, \ 0.8\}$.
- Draw $\beta \sim \text{N}(0, \tau^2 I_p)$
- Draw $Y_i = \mathbb{X}_i^\top \beta + \epsilon_i$ with $\epsilon_i \sim \text{N}(0, \sigma^2)$.

# BAYES ESTIMATOR

- For this model, the optimal estimator (in MSE) is

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

- In particular, with $\lambda_* = \frac{\sigma^2}{n\tau^2}$
- This is the posterior mode of the Bayes estimator under conjugate normal prior
- It is also the ridge regression estimator for a particular $\lambda$

With $p < n$

1. If $\lambda_*$ is too big, we will tend to shrink all coefficients to 0.
   - This problem is too hard.
2. If $\lambda_*$ is too small, OLS will be very close to the optimal estimator.
   - This problem is too easy.
3. Need $\tau^2$, $\sigma^2$ "just right".

- Take $\tau^2 = \pi/2$. This implies $\mathbb{E}[|\beta_i|] = 1$ (convenient)
- Take $n = 5000$. Big but computable.
- Take $\sigma = 50 \Rightarrow \log(\lambda_*) \approx -1.14$ (reasonable)
- Take $p \in \{50, \ 100, \ 250, \ 500\}$

$\rho = 0.8$, p = 50, q/n = 0.2

log₁₀(estimation risk) vs log(λ)

convex   linear   FC   PC   ols   bayes

22

ρ = 0.2, p = 500, q/n = 0.2

convex  linear  FC  PC  ols  bayes

In that case, ridge was optimal.

We did it again with $\beta_i \equiv 1$.

$\rho = 0.8$, p = 50, q/n = 0.2

convex    linear    FC    PC    ols    ridge

$\rho = 0.2, p = 500, q/n = 0.2$

x-axis: $\log(\lambda)$
y-axis: $\log_{10}(\text{estimation risk})$

legend: convex, linear, FC, PC, ols, ridge

28

# SELECTING TUNING PARAMETERS

- We use GCV with the degrees of freedom:

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

- $df$ is easy for full or partial compression.
- For the other cases, we do a dumb approximation: weighted average using $\widehat{\alpha}$.
- Easy to calculate for a range of $\lambda$ without extra computations.
- Also calculate the divergence to derive Stein's Unbiased Risk Estimate.
- After tedious algebra, it had odd behavior in practice (can be huge, or negative!?!)

$$n = 5000, \quad p = 50, \quad \rho = 0.8, \quad \beta \sim \mathbf{N}(0, \pi/2)$$

# Theoretical results (sketch)

Theorem

$$\text{bias}^2\left(\widehat{\beta}_{ridge}(\lambda) \mid \mathbb{X}\right) = \lambda^2 \beta_*^\top V (D^2 + \lambda I_p)^{-2} V^\top \beta_*.$$

$$\text{tr}\left(\mathbb{V}[\widehat{\beta}_{ridge}(\lambda) \mid \mathbb{X}]\right) = \sigma^2 \sum_{i=1}^{p} \frac{d_i^2}{(d_i^2 + \lambda)^2}.$$

# WHAT'S THE TRICK?

- Similar results are hard for compressed regression.
- All the estimators depend (at least) on

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

- We derived properties of $Q^\top Q$

$$\mathbb{E}\left[\frac{s}{q} Q^\top Q\right] = I_n$$

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

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

Theorem

$$\text{bias}^2[\widehat{\beta}_{FC} \mid \mathbb{X}] = \lambda^2 \beta_*^\top V (D^2 + \lambda I_p)^{-2} V^\top \beta_* + o_p(1)$$

$$\text{tr}(\mathbb{V}[\widehat{\beta}_{FC} \mid \mathbb{X}]) = \sigma^2 \sum_{i=1}^{p} \frac{d_i^2}{(d_i^2 + \lambda)^2} + o_p(1)$$

$$+ \frac{(s-3)_+}{q} \text{tr} \left( \text{diag}(\text{vec}\,(I_n)) M^\top M \otimes (I - H) M \beta_* \beta_*^\top M^\top (I - H) \right)$$

$$+ \frac{\beta_*^\top M^\top (I - H)^2 M \beta_*}{q} \text{tr}(M M^\top)$$

$$+ \frac{1}{q} \text{tr} \left( (I - H) M \beta_* \beta_*^\top M^\top (I - H) M^\top M \right).$$

Note: $M = (\mathbb{X}^\top \mathbb{X} + \lambda I_p)^{-1} \mathbb{X}^\top$ and $H = \mathbb{X} M$ (hat matrix for ridge regression)

# Applications

# RNA-SEQ

- short-read RNA sequence data
- using a (Poisson) linear model to predict read counts based on the surrounding nucleotides
- 8 different tissues
- data is publicly available, preprocessing already done

- Test error averaged over 10 replications.
- On each replication, tuning parameters were chosen with GCV.
- At the best tuning parameter, the test error was computed and then these were averaged.
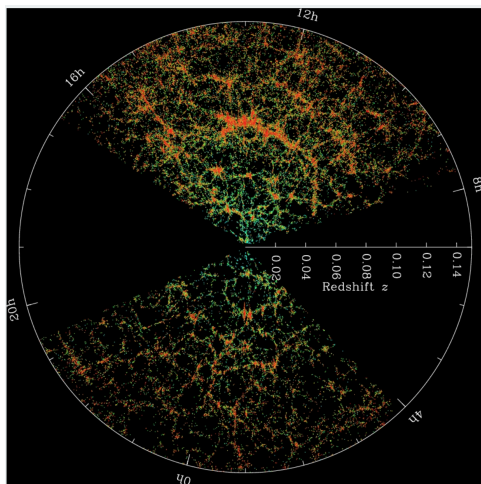- For each dataset, we randomly chose 75% of the data as training on each replication.

# TAKE-AWAY MESSAGE

- $q = 10000$ results in data reductions between 74% and 93%
- $q = 20000$ gives reductions between 48% and 84%
- ridge and ordinary least squares give equivalent test set performance (differing by less than .001%)
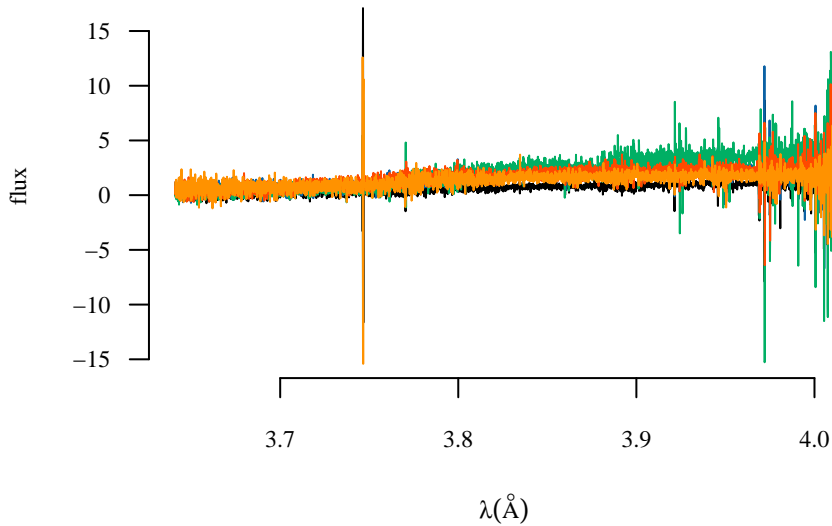- This is an "easy" problem.
- Full compression is the worst.

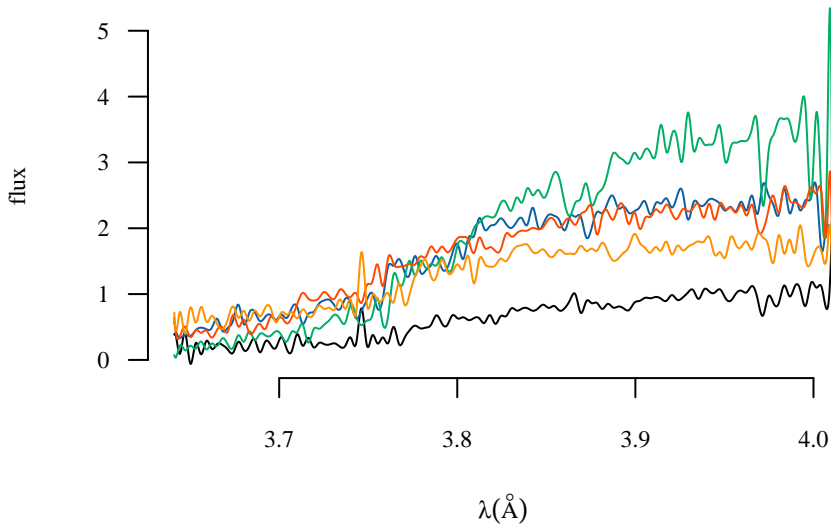This image and the next are from the Sloan Digital Sky Survey

- Repeatedly "photograph" regions of the sky
- Collect information on piles of objects
- Around 200 million galaxy spectra at the moment
- Each spectra contains around 5000 wavelengths
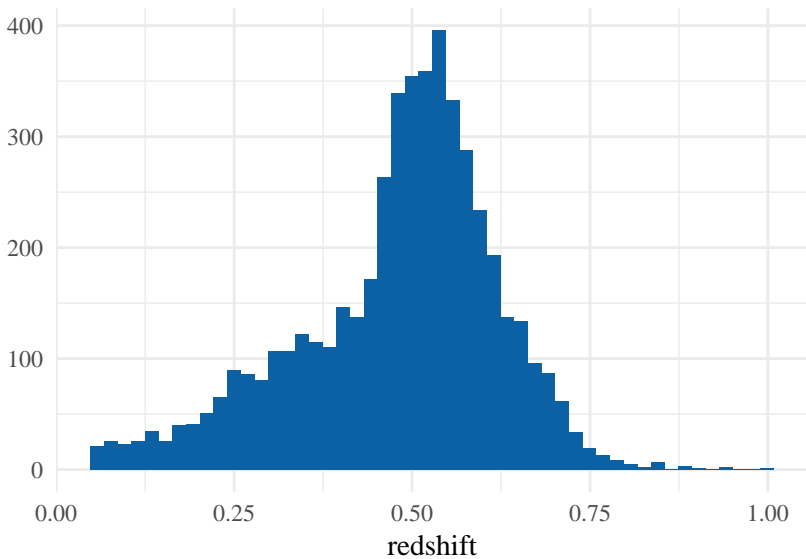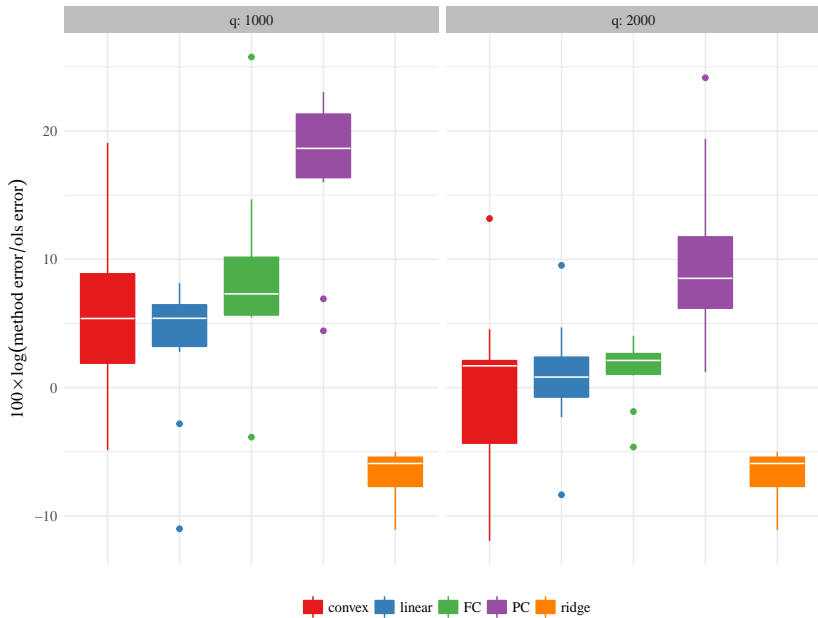- Want to predict the redshift from the spectra

A FEW GALAXIES

# SAME GALAXIES AFTER SMOOTHING

5000 REDSHIFTS

# Conclusions

# CONCLUSIONS

- Lots of people are looking at approximations to standard statistical methods (OLS, PCA, etc.)
- They characterize the approximation error
- We have been looking at whether we can actually benefit from the approximation
- Today looked at compressed regression

# ONGOING AND FUTURE WORK

- We want to generalize our results here to other penalties
- Also generalized linear models
- We're also looking at how compression interacts with PCA
- Combine compression and random projection
- Connections to sufficient dimension reduction?
- Iterative versions?