

Topic 2: Scaling up

Eric B. Laber

Department of Statistical Science, Duke University

Statistics 561



On big data

Across all of its platforms, Google collects 25 petabytes of data every single day. That amount of data was unfathomable just a few years ago. To put this into perspective, imagine that each bit were a grain of rice. In just 8 hours, we would have enough rice to cover...you know what, if you're the kind of person that has better intuition for astronomically large piles of rice than digital data, what's the point of telling you anything? Do you have better intuition about how that rice should be stored? Manipulated? Used to predict new kinds of rice? Get the hell out of my office.
—Cassie Kozyrkov, Chief Decision Scientist at Google



On computing

*One thing a computer can do that most humans can't is
be sealed up in a cardboard box and sit in a warehouse.
—Avie Tevanian, CTO Apple Inc.*



Warm-up (5 minutes)

- ▶ Explain to your group
 - ▶ What is stochastic gradient descent?
 - ▶ What is streaming data? What about a streaming algorithm?
 - ▶ What is an online prediction problem?
- ▶ True or false
 - ▶ The Johnsonville-Liverwurst theorem is a key result in dimension reduction
 - ▶ Googling "SIS" when searching for Sure Independent Screening yields a bunch of terrible YouTube videos that should never be viewed by anyone
 - ▶ Valentine's day evolved from the ancient Roman festival of Lupercalia in which lovers showed their affection through gifts and songs. The new holiday was created because Lupercalia was often confused with the festival of Lupinlavatio in which peasants, chosen by lottery, were forced to bath wolves.



Roadmap

- ▶ It's time to relax: running sums
- ▶ Random projections
- ▶ Stochastic approximation



Roadmap

- ▶ **It's time to relax: running sums**
- ▶ Random projections
- ▶ Stochastic approximation



Review: fitting a linear model

- ▶ Ordinary least squares estimator as

$$\hat{\beta}_n = \arg \min_{\beta \in \mathbb{R}^p} \mathbb{P}_n(Y - \mathbf{X}^\top \beta)^2 = \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{X}_i^\top \beta)^2$$

- ▶ Alternatively, view $\hat{\beta}_n$ as solution to

$$\mathbb{P}_n \mathbf{X} (Y - \mathbf{X}^\top \beta) = 0$$

- ▶ If $\mathbb{P}_n \mathbf{X} \mathbf{X}^\top$ is invertible $\hat{\beta}_n = (\mathbb{P}_n \mathbf{X} \mathbf{X}^\top)^{-1} \mathbb{P}_n \mathbf{X} Y$

Scaling up

- ▶ Collecting, storing, accessing, manipulating data is increasingly easy
 - ▶ More data, faster
 - ▶ Large p , lots of measurements per obs.
 - ▶ Large n , many obs.
- ▶ Need to balance statistical and computational efficiency

Large p (Part I)

- ▶ Methods so far designed for setting where p is smaller than n
- ▶ If $p \gg n$ we may want to screen out seemingly unimportant variables first, then apply the methods we've seen so far
 - ▶ In first step we want to make sure we include important ones, don't worry too much about accidentally including unimportant ones



Sure independence screening (SIS)

- ▶ Rank variables based on their marginal association with Y and choose subset that are ranked highest
 - ▶ Assume \mathbb{X} , Y are centered and scaled

$$W_j = \mathbb{P}_n X_j Y$$

select set of variables $\hat{\mathcal{J}}_M = \{j : |W_j| \text{ is among the } M \text{ largest}\}$

- ▶ R code example: `screen.R`

Sure independence screening (SIS) cont'd

- ▶ How do choose number of variables M ?
 - ▶ Domain knowledge/expert judgment
 - ▶ Could treat M as tuning parameter opt using CV etc.
 - ▶ Error rate control
 - ▶ Many methods exist to avoid accidentally including noise variables (this is not exactly what we want)
 - ▶ Typically based on p -values (FDR control, etc.)

Big n (and possibly big p too)

- ▶ In some settings it may not be feasible to fit the model to all the data in one batch
 - ▶ n is very large, e.g., several GBs or PBs
 - ▶ Data are streaming in over time
- ▶ Trade-off statistical and computational/memory efficiency
- ▶ Key ideas
 - ▶ Divide, distribute, and conquer
 - ▶ Streaming, online updates

Divide, distribute, and conquer

- ▶ Simplest strategy: draw a sample of size $m \ll n$
 - ▶ Fit model, perform inference/validation etc.
 - ▶ Pros: representative sample, computationally feasible
 - ▶ Cons: fail to identify weak signals, lower quality model, may underestimate predictive accuracy
- ▶ Second simplest strategy: partition data into manageable chunks
 - ▶ Fit model on each partition (can do this in parallel)
 - ▶ Aggregate estimators
 - ▶ With parametric model, average parameter estimates
 - ▶ More generally, with a predictive model you can average predicted values from each model



Divide, distribute, and conquer: example 1

- ▶ Sample mean $\mathbb{P}_n \mathbf{X}$
 - ▶ Create partition C_1, \dots, C_K of $\{1, \dots, n\}$, compute

$$S_j = \sum_{i \in C_j} \mathbf{X}_i \quad (\text{serially or in parallel})$$

$$\text{then } \mathbb{P}_n \mathbf{X} = n^{-1} \sum_{j=1}^K S_j$$

- ▶ Note that $\{\mathbf{X}_i\}_{i \in C_j, j=1, \dots, K}$ do not need to be loaded into memory at the same time

Divide, distribute, and conquer: example 2

- ▶ Least squares $\hat{\beta}_n = (\mathbb{P}_n \mathbf{X} \mathbf{X}^\top)^{-1} \mathbb{P}_n \mathbf{X} Y$
 - ▶ Create partition C_1, \dots, C_K of $\{1, \dots, n\}$, compute

$$\Sigma_j = \sum_{i \in C_j} \mathbf{x}_i \mathbf{x}_i^\top$$

$$\Gamma_j = \sum_{i \in C_j} \mathbf{x}_i Y_i$$

$$\text{then } \hat{\beta}_n = \left\{ \sum_{j=1}^K \Sigma_j \right\}^{-1} \sum_{j=1}^K \Gamma_j$$

Divide, distribute, and conquer: example 3

- ▶ Lasso estimator $\hat{\beta}_n^\tau = \arg \min_{\beta} \mathbb{P}_n(Y - \mathbf{X}^\top \beta)^2 + \tau \|\beta\|_1$
 - ▶ Fact: $\hat{\beta}_n^\tau$ is completely determined by $\mathbb{P}_n \mathbf{X} \mathbf{X}^\top$ and $\mathbb{P}_n \mathbf{X} Y$
 - ▶ We can use distd least squares algorithm for lasso

PF:

$$\begin{aligned} \mathbb{P}_n(Y - \mathbf{X}^\top \beta)^2 + \tau \|\beta\|_1 &= \mathbb{P}_n(Y - \mathbf{X}^\top \hat{\beta}_n)^2 + \mathbb{P}_n(\mathbf{X}^\top \hat{\beta}_n - \mathbf{X}^\top \beta)^2 + \tau \|\beta\|_1 \\ &= \text{Const} + (\hat{\beta}_n - \beta)^\top \mathbb{P}_n \mathbf{X} \mathbf{X}^\top (\hat{\beta}_n - \beta) + \tau \|\beta\|_1 \end{aligned}$$



Divide, distribute, and conquer: example 3 cont'd

- ▶ Write $\hat{\beta}_n^\tau = \arg \min_{\beta} (\hat{\beta}_n - \beta)^\top \mathbb{P}_n \mathbf{X} \mathbf{X}^\top (\hat{\beta}_n - \beta) + \tau \|\beta\|_1$
- ▶ Create partition C_1, \dots, C_K of $\{1, \dots, n\}$, compute $\Sigma_j = \sum_{i \in C_j} \mathbf{X}_i \mathbf{X}_i^\top$, $\Gamma_j = \sum_{i \in C_j} \mathbf{X}_i Y_i$, and $\hat{\beta}_n$ using distd least squares algorithm
- ▶ Compute

$$\hat{\beta}_n^\tau = \arg \min_{\beta} n^{-1} (\hat{\beta}_n - \beta)^\top \left(\sum_{j=1}^K \Sigma_j \right) (\hat{\beta}_n - \beta) + \tau \|\beta\|_1$$

- ▶ Can be solved as quadratic program or using subgrad descent

Divide, distribute, and conquer: code

- ▶ R code example: distributed.R
- ▶ In class exercise: with your neighbor design and implement distributed ridge regression
- ▶ Group discussion: how to implement distributed screening?

Warm-up quiz (5 Minutes)

- ▶ Explain to the your stats group
 - ▶ What is right censored data? Give examples.
 - ▶ What is a the Kaplan-Meyer estimator?
 - ▶ What is the Cox proportional hazards (Cox-PH) model?
- ▶ True or false
 - ▶ Stochastic gradient descent was created to fit neural networks
 - ▶ Survival data is not of much interest outside of biomedical applications
 - ▶ The best way to learn about string theory is to watch 'Spatula Madness'



Survival analysis background

- ▶ Censoring: observe partial information about an event
 - ▶ Survival time in study with limited follow-up \Rightarrow patients that are still alive at end of study are said to be right-censored¹
 - ▶ Failure time of components in a system with sporadic inspection \Rightarrow components that have failed at first inspection are said to be left-censored
 - ▶ Victims of an avalanche \Rightarrow when you dig them up they're alive or dead so the death time is said to be interval-censored
- ▶ Why are we talking about this?
 - ▶ Common in application
 - ▶ Illustrate more complex accounting in streaming estimators

¹Other canonical examples include time to loan default, time to machine breakdown, time to reaching goal for robot exploring, ...

Sketch of survival: right-censoring

Divide, distribute, and conquer: KM

- ▶ Suppose we observed right censored data $\{T_i, \delta_i\}_{i=1}^n$
 - ▶ T is (discrete) censoring or failure time
 - ▶ δ denotes censoring indicator ($\delta = 1$ if failure)
 - ▶ Failure has density $f(t)$ and survival fn $S(t)$
- ▶ Kaplan-Meier estimator

$$\hat{S}(t) = \prod_{\ell: t_\ell \leq t} \left(1 - \frac{d_\ell}{n_\ell}\right),$$

where $t_1 < t_2 < \dots < t_k$ are distinct failure times²

- ▶ $d_\ell = \sum_{i=1}^n 1_{T_i=t_\ell} \delta_i$ is the number of failures at t_ℓ
- ▶ $n_\ell = \sum_{i=1}^n 1_{T_i \geq t_\ell}$ is the number at risk at time t_ℓ

²Not censoring times

Distributed KM

- ▶ Distributed computation of KM
 - ▶ Create partition C_1, \dots, C_K of $\{1, \dots, n\}$, compute

$$n_{\ell,j} = \sum_{i \in C_j} 1_{T_i \geq t_\ell} \quad \text{at risk at time } t_\ell$$

$$d_{\ell,j} = \sum_{i \in C_j} 1_{T_i = t_\ell} \delta_i \quad \text{fail at time } t_\ell,$$

$$\text{then } n_\ell = \sum_{j=1}^K n_{\ell,j} \text{ and } d_\ell = \sum_{j=1}^K d_{\ell,j}$$

- ▶ What challenges do you foresee with this implementation?

Distributed KM cont'd

- ▶ Don't know the unique failure times at onset
 - ▶ Use dictionary with keys corresponding to unique failure times
 - ▶ Dynamically add keys as new failure times are encountered
- ▶ Go to `km_streaming.py`

Exact distributed computing

- ▶ So far, we've developed distributed algorithms that give us *exactly* the same solution as if we had used all the data at once
 - ▶ Linear estimators (or fns of linear estimators)
 - ▶ Counts
- ▶ Unfortunately, not all estimators lend themselves to nice analytic decompositions that are easily distributed
- ▶ However, many statistical estimators minimize (maximize) an objective that sums over data
 - ▶ Exploit this structure to construct general distd algorithms
 - ▶ Our primary tool will be approximate gradient descent



Quick review (sub)grad. descent

- ▶ Goal: minimize the function $f(\theta)$ over $\theta \in \mathbb{R}^P$
 - ▶ $f(\theta)$ increases most rapidly in the direction of $\nabla f(\theta)$
 - ▶ Gradient descent update: $\theta^{(k+1)} = \theta^{(k)} - \alpha^{(k)} \nabla f(\theta^{(k)})$
- ▶ Choosing the step-size $\alpha^{(k)}$ is non-trivial
 - ▶ Line-search $\alpha^{(k)} = \arg \min_{\alpha} f \left\{ \theta^{(k)} - \alpha \nabla f(\theta^{(k)}) \right\}$
 - ▶ Barzilai-Borwein

$$\alpha^{(k)} = \frac{(\theta^{(k)} - \theta^{(k-1)})^\top \{ \nabla f(\theta^{(k)}) - \nabla f(\theta^{(k-1)}) \}}{\| \nabla f(\theta^{(k)}) - \nabla f(\theta^{(k-1)}) \|^2}$$

- ▶ If $f(\theta)$ convex but not differentiable, replace gradient with a subgradient (this is what we did with distd lasso)



M-estimators

- ▶ Many statistical estimators have the form

$$\hat{\theta}_n = \arg \min_{\theta} g \{ \mathbb{P}_n f(\mathbf{Z}_i; \theta) \}$$

- ▶ Canonical examples

- ▶ $g(u) = u$, $\mathbf{Z} = (\mathbf{X}, Y)$ and $f(\mathbf{Z}; \theta) = (Y - \mathbf{X}^\top \theta)^2$
- ▶ $g(u) = u$, $f(\mathbf{Z}; \theta)$ is negative log-density (max LH)
- ▶ $g(u) = ||u||^2$, and $\mathbb{P}_n f(\mathbf{Z}; \theta) = 0$ is est. eqn.

Distributed approximate gradient descent

- ▶ Gradient using all the data

$$\frac{\partial}{\partial \theta^\top} g \{ \mathbb{P}_n f(\mathbf{Z}_i; \theta) \} = g' \{ \mathbb{P}_n f(\mathbf{Z}_i; \theta) \} \mathbb{P}_n \nabla f(\mathbf{Z}_i; \theta)$$

- ▶ Distributed approximate gradient descent

1. Create partition C_1, \dots, C_K of $\{1, \dots, n\}$
2. Set $\theta^{(1)}$ to starting value, set $k = 1$
3. Randomly select partition j from $\{1, \dots, K\}$, compute

$$\delta^{(k)} = g' \left\{ n_j^{-1} \sum_{i \in C_j} f(\mathbf{Z}_i; \theta^{(k)}) \right\} n_j^{-1} \sum_{i \in C_j} \nabla f(\mathbf{Z}_i; \theta^{(k)})$$

- ▶ Update $\theta^{(k+1)} = \theta^{(k)} - \alpha^{(k)} \delta^{(k)}$

Distributed Cox-PH

- ▶ Full data are $\{(\mathbf{X}_i, T_i, \delta_i)\}_{i=1}^n$ where \mathbf{X} are covariates, T is obs. time, and δ is censoring indicator
 - ▶ Recall that the Cox-PH model postulates

$$\lambda(t|\mathbf{x}) = \lambda_0(t) \exp \{\mathbf{x}^\top \theta\}$$

- ▶ Define $Y_\ell(t) = 1_{T_\ell \geq t}$, Cox-PH estimator

$$\hat{\theta}_n = \arg \max_{\theta} \sum_{i=1}^n \delta_i \left[\mathbf{X}_i^\top \theta - \log \left\{ \sum_{\ell=1}^n Y_\ell(T_i) \exp(\mathbf{X}_\ell^\top \theta) \right\} \right],$$

Recall that the above is the log partial likelihood

Distributed Cox-PH cont'd

- Evaluate partial likelihood on subset C_j using

$$\sum_{i \in C_j} \delta_i \left[\mathbf{X}_i^T \theta - \log \left\{ \sum_{\ell \in C_j} Y_\ell(T_i) \exp(\mathbf{X}_\ell^T \theta) \right\} \right],$$

- Gradient of negative partial log-likelihood on subset C_j is

$$U_j(\theta) = - \sum_{i \in C_j} \delta_i \left\{ \mathbf{X}_i - \bar{\mathbf{X}}_{C_j}(T_i, \theta) \right\},$$

where

$$\bar{\mathbf{X}}_{C_j}(t, \theta) = \sum_{\ell \in C_j} \mathbf{X}_\ell Y_\ell(t) \exp(\mathbf{X}_\ell^T \theta) / \sum_{\ell \in C_j} Y_\ell(t) \exp(\mathbf{X}_\ell^T \theta)$$

Distributed Cox-PH cont'd

- ▶ Distributed approximate gradient descent Cox-PH
 1. Create partition C_1, \dots, C_K of $\{1, \dots, n\}$
 2. Set $\theta^{(1)}$ to starting value, set $k = 1$
 3. Randomly select partition j from $\{1, \dots, K\}$
 - ▶ Update $\theta^{(k+1)} = \theta^{(k)} - \alpha^{(k)} U_j^{(k)}(\theta^{(k)})$
- ▶ R code break: `coxph.R`

Distributed computing discussion

- ▶ Reviewed basic distd optimization
 - ▶ Exact distd programming
 - ▶ Approximate gradient descent
- ▶ These methods can be applied to data streaming in real-time (the partition elements C_1, \dots, C_k can be arriving as data accumulates)
- ▶ Alternative strategy is to fit separate models and aggregate/average estimators

Break: algorithm quiz

- ▶ Recall that the median satisfies $\hat{m}_n = \arg \min_m \mathbb{P}_n |Y - m|$
 - ▶ With your stats group, design an algorithm that estimates the median using only one data point at a time
 - ▶ At home: implement in R or Python, how does your algorithm compare with batch computation of the median?

The lecture your parents never gave you...

- ▶ Where does big data come from?
 - ▶ High-volume: website exhaust, high-resolution images/video, sensors running in near continuous time, etc.
 - ▶ High-dim: lots of measurements³ and **feature construction**

³What a deep statement!

Basis expansions

- ▶ Linear model viewed as a first-order approximation

- ▶ I.e., $f(\mathbf{x}) = f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^\top (\mathbf{x} - \mathbf{x}_0) + \text{HOT}$

- ▶ Natural approach add flexibility is to include HOT

- ▶ Monomial regression

- ▶ No interactions: $f(\mathbf{x}) \approx \beta_0 + \sum_{j=1}^p \sum_{k=1}^r \beta_{j,k} x_i^k$

- ▶ Pairwise interactions:

$$f(\mathbf{x}) \approx \beta_0 + \sum_{j=1}^p \sum_{k=1}^r \beta_{j,k} x_i^k + \sum_{j=1}^p \sum_{\ell=1}^p \sum_{k=1}^{r-1} \beta_{j,\ell,k} x_j^k x_\ell^{r-k}$$

- ▶ Intuitive but unstable as degree of polynomial increases

On instability

When the wheat started undulating in the wind, the farmer did too, in a way that he had been warned about before. When the police dragged him off, in handcuffs, he was still undulating. —John Stuart Mill

Basis expansions cont'd

- ▶ Consider approximations of the form

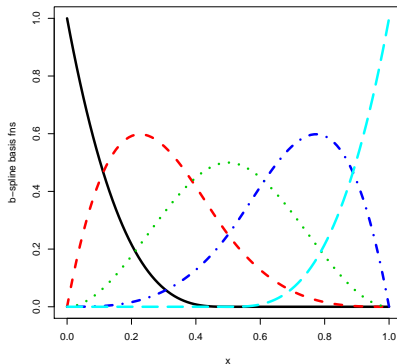
$$f(\mathbf{X}) \approx \beta_0 + \sum_{j=1}^q \beta_j b_j(\mathbf{X}),$$

where $b_j : \mathbb{R}^P \rightarrow \mathbb{R}$ are basis functions

- ▶ E.g., $b_j(\mathbf{X}) = X_{k(j)}^{d(j)}$ in polynomial w/o interactions
- ▶ Typically choose basis $\{b_j(\cdot)\}_{j \geq 1}$ to be dense in some space of interest, e.g., the space of continuous functions $[0, 1]$ etc.

B-splines

- ▶ Polynomial segments between knots
 $0 = a_0 < a_1 < \dots < a_K = 1$
- ▶ Constrained to be smooth at knots
- ▶ Specification of a B-spline basis requires order of polynomial segments and knots



B-splines cont'd

- ▶ Univariate R code example: `nonlinearReg.R`

Other basis systems

- ▶ Fourier basis $\{(\sin(k\omega x), \cos(k\omega x))\}_{k \geq 1}$
 - ▶ Periodic may not be good for growth curves etc.
- ▶ Wavelet basis (many variants)
 - ▶ Widely used for “denoising” in signal processing
 - ▶ A number of deep theoretical results exist for nonparametric regression with wavelet bases

Multiple predictors

- ▶ Given basis system $\{b_j(\cdot)\}_{j \geq 1}$ we create 'features' $\{b_j(X_k)\}_{j \geq 1}$, $k = 1, \dots, p$
 - ▶ We can construct a model with interactions among basis functions etc.
 - ▶ Number of terms in the model can grow quickly, e.g., with J basis fns and pairwise interactions $O(p^2 J^2)$ terms
- ▶ Use all of the model building tools we already talked about to deal with large p

Roadmap

- ▶ It's time to relax: running sums
- ▶ **Random projections**
- ▶ Stochastic approximation



On engaging the scientific community

*Buckle up buttercups. This sh*t is going to blow your little minds. —Vladimir Vapnik, addressing the Royal Society, London 1991*

*It's easy to look like you're standing on the shoulders of giants when you're among tiny people.
—Speech of stats faculty member at reception celebrating their election to NAS⁴*

⁴This one is true; it was not a Duke faculty member :)



Background: dual form of linear regression

- ▶ Overview of this section
 - ▶ Show that ridge can be formulated in terms of inner products
 - ▶ Quick aside: yet another characterization of ridge regression
 - ▶ Show inner products under random projections are close to inner products in original space \Rightarrow regression on projections close to original

Dual ridge regression

- ▶ Observe $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ drawn i.i.d. from P
- ▶ Let $\hat{\boldsymbol{\beta}}_n^\lambda$ denote ridge estimator with penalty λ , i.e.,

$$\hat{\boldsymbol{\beta}}_n^\lambda = \arg \min_{\boldsymbol{\beta}} \|\mathbb{Y} - \mathbb{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|^2$$

then for a new observation $\mathbf{X} = \mathbf{x}$ we have

$$\mathbf{x}^\top \hat{\boldsymbol{\beta}}_n^\lambda = \frac{1}{2\lambda} \sum_{i=1}^n \hat{\alpha}_{n,i}^\lambda \mathbf{X}_i^\top \mathbf{x} = \mathbb{Y}^\top (K + \lambda I_n)^{-1} \phi(\mathbf{x}),$$

where $K \in \mathbb{R}^{n \times n}$ is matrix of inner products, i.e., $K_{i,j} = \mathbf{X}_i^\top \mathbf{X}_j$, and $\phi(\mathbf{x}) \in \mathbb{R}^n$ is given by $\phi_i(\mathbf{x}) = \mathbf{X}_i^\top \mathbf{x}$

- ▶ This is known as the dual version of ridge regression

Strategy for deriving dual ridge regression⁵

- ▶ Express as (convenient) constrained optimization problem
- ▶ Introduce Lagrange multipliers to make unconstrained
- ▶ Minimize over primal variables
- ▶ Maximize over dual variables

⁵Proof taken almost directly from Saunders et al., 1998, ICML



Blank page for notes: Lagrange form



Blank page for notes: minimize over primal



Blank page for notes: maximize over dual



Blank page for notes: spillover



Discussion dual form

- ▶ Ridge (and thus OLS) predictions depend only on inner products which is intuitive but not obvious in primal form
 - ▶ Basis for kernel regression methods (you saw in lab)
 - ▶ It p -large working with inner-products may reduce computational burden, however, if n large may need sparse approx to K
- ▶ Keep dual form in mind when we cover random projections

But first! Yet another view of ridge!

- Suppose we want to predict the outcome Y at a new $\mathbf{X} = \mathbf{x}$ and we want our prediction to be of the form

$$\sum_{i=1}^n \hat{\omega}_{n,i}(\mathbf{x}) Y_i,$$

where $\omega_i(\mathbf{x})$ captures the similarity between \mathbf{X}_i and \mathbf{x}

- Idea! Let $\boldsymbol{\omega} \in \mathbb{R}^n$ solve

$$\hat{\boldsymbol{\omega}}_n(\mathbf{x}) = \arg \min_{\boldsymbol{\omega} \in \mathbb{R}^n} \|\mathbf{x} - \mathbb{X}^T \boldsymbol{\omega}\|^2 + \lambda \|\boldsymbol{\omega}\|^2,$$

i.e., we're representing the new input as a linear combination of training inputs \Rightarrow prediction $\mathbb{Y}^T (K + \lambda I_n)^{-1} \phi(\mathbf{x}) = \text{Ridge!}$



Blank page for notes: deriving ridge estimator



Random projections

- ▶ Let $\Omega \in \mathbb{R}^{k \times p}$ with $k \ll p$
 - ▶ If $(\Omega \mathbf{X}_i)^\top (\Omega \mathbf{X}_j) \approx \mathbf{X}_i^\top \mathbf{X}_j$ for all $i, j \Rightarrow$ fitted regression on reduced data $\{(\Omega \mathbf{X}_i, Y_i)\}_{i=1}^n$ similar to fitted on original data
 - ▶ Can store and use reduced data $\{(\Omega \mathbf{X}_i, Y_i)\}_{i=1}^n \Rightarrow$ faster computation + lower storage
- ▶ How to construct Ω ?
- ▶ To `simple_random_projection.R`

Normies are all the same

Theorem

Let $\mathbf{x} \in \mathbb{R}^p$ and assume $\Omega \in \mathbb{R}^{k \times p}$ is populated with i.i.d. standard normal random variables. Then,

$$P \left\{ (1 - \epsilon) \|\mathbf{x}\|^2 \leq \left\| \frac{1}{\sqrt{k}} \Omega \mathbf{x} \right\|^2 \leq (1 + \epsilon) \|\mathbf{x}\|^2 \right\} \geq 1 - 2 \exp \left\{ -\frac{\epsilon^2 (1 - \epsilon) k}{4} \right\}.$$



Blank page for notes: proving normies are all the same



Blank page for notes: spillover



Warm up: Fact about χ^2 random variables

Theorem

Let χ_k^2 denote a χ^2 random variable with k degrees of freedom. Then for any $\epsilon \in (0, 1)$ it follows that

$$P \{ \chi_k^2 \geq (1 + \epsilon)k \} \leq \exp \{ -\epsilon^2(1 - \epsilon)k/4 \}$$

and

$$P \{ \chi_k^2 \leq (1 - \epsilon)k \} \leq \exp \{ -\epsilon^2(1 - \epsilon)k/4 \}.$$

Proof is based on Markov's inequality. We won't go through derivation here.

Blank page for notes



Blank page for notes



Johnson-Lindenstrauss Lemma

Theorem

Let $\epsilon \in (0, 1/2)$ be given. Let $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ and $k = 20 \log(n)/\epsilon^2$. Then there exists a Lipschitz mapping⁶ $f : \mathbb{R}^p \rightarrow \mathbb{R}^k$ such that for all i, j

$$(1 - \epsilon) \|\mathbf{x}_i - \mathbf{x}_j\|^2 \leq \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|^2 \leq (1 + \epsilon) \|\mathbf{x}_i - \mathbf{x}_j\|^2.$$

⁶Recall that f is Lipschitz if there exists a constant c such that $\|f(\mathbf{x}) - f(\mathbf{x}')\| \leq c \|\mathbf{x} - \mathbf{x}'\|$ for all \mathbf{x}, \mathbf{x}' .

Blank page for notes: Johsonville-Liverwurst theorem



Blank page for notes: Johsonville-Liverwurst theorem



Amuse-bouche: 3 minutes

- ▶ Explain to your stats group:
 - ▶ What is a sparse matrix?
 - ▶ How can we store a sparse matrix efficiently?

Projections: no longer just for normies

- Recall: Rademacher random variable is uniform on $\{-1, 1\}$

Theorem

Assume $\Omega \in \{-1, 1\}^{p \times p}$ be populated with Rademacher random variables. Then

$$P \left\{ \left\| \frac{1}{\sqrt{k}} \Omega \right\|^2 \geq (1 + \epsilon) \|\mathbf{x}\|^2 \right\} \leq \exp \left\{ -\frac{\epsilon^2(1 - \epsilon)k}{4} \right\}$$

and

$$P \left\{ \left\| \frac{1}{\sqrt{k}} \Omega \right\|^2 \leq (1 - \epsilon) \|\mathbf{x}\|^2 \right\} \leq \exp \left\{ -\frac{\epsilon^2(1 - \epsilon)k}{4} \right\},$$

thus the preceding theorems will go through with Rademacher variables in place of standard normals.

Sparse random projections

- ▶ Proj inputs $\Omega \mathbf{X} \in \mathbb{R}^k$ save space and computation as $k \ll p$
- ▶ Can obtain further space and computation savings if the projected inputs are sparse, i.e., if they have many zero entries
- ▶ Consider a matrix Ω comprising i.i.d. random variables

$$\Omega_{i,j} = \begin{cases} \sqrt{\kappa} & \text{with probability } 1/(2\kappa) \\ -\sqrt{\kappa} & \text{with probability } 1/(2\kappa) \\ 0 & \text{with probability } 1 - 1/\kappa, \end{cases}$$

e.g., $\kappa = 3$ we only need to store/use (1/3) of the projected entries. It turns out that this will also satisfy the preceding theorems.⁷

⁷See Li et al., KDD 2006, for results on even more sparse (sparser?) projections.

Preserving inner products

Theorem

Let $\mathbf{x}, \mathbf{v} \in \mathbb{R}^p$ be such that $\|\mathbf{x}\| \leq 1$, and $\|\mathbf{v}\| \leq 1$. Define $f(\mathbf{x}) = k^{-1/2} \Omega \mathbf{x}$ where Ω is composed of i.i.d. standard normal entries. Then for any $\epsilon > 0$

$$P \{ |\mathbf{x}^T \mathbf{v} - f(\mathbf{x})^T f(\mathbf{v})| \geq \epsilon \} \leq 4 \exp \{ -\epsilon^2(1 - \epsilon)/4 \}.$$

Blank page for notes



Blank page for notes



Roadmap

- ▶ It's time to relax: running sums
- ▶ Random projections
- ▶ **Stochastic approximation**



On randomness

The world of science lives fairly comfortably with paradox. We know that light is a wave, and also that light is a particle. The discoveries made in the infinitely small world of particle physics indicate randomness and chance, and I do not find it any more difficult to live with the paradox of a universe of randomness and chance and a universe of pattern and purpose than I do with light as a wave and light as a particle...the distinction means little when that light is concentrated by a deep state space laser.

—Marjorie Taylor Greene

Warm-up: estimating a mean (5 min)

- Suppose we have data $\mathbf{X}_1, \dots, \mathbf{X}_n, \dots$ drawn i.i.d. from a distribution with mean μ . Our goal is to construct an estimator from streaming data. Show that for any n we have

$$\bar{\mathbf{X}}_n = \bar{\mathbf{X}}_{n-1} - \alpha_n (\bar{\mathbf{X}}_{n-1} - \mathbf{X}_n),$$

for some learning rate $\{\alpha_n\}_{n=1}^\infty$ such that $\alpha_n > 0$,
 $\sum_n \alpha_n = \infty$, and $\sum_{n \geq 1} \alpha_n^2 < \infty$.

Warm-up II: Jensen's inequality (5 min)

- Explain Jensen's inequality to your stats group. Use it to prove that

$$\left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_n) \right\}^2 \leq \mathbb{P}_n(Y - \mathbf{X}^\top \hat{\boldsymbol{\beta}}_n)^2.$$

Estimating a mean with more general step-sizes

- ▶ Suppose we have $X_1, X_2, \dots \sim_{i.i.d.} (\mu, \sigma^2)$
- ▶ Consider recursive estimator such that $\hat{\mu}_1 = \alpha_1 X_1$ and for $n \geq 2$

$$\hat{\mu}_n = \hat{\mu}_{n-1} - \alpha_n (\hat{\mu}_{n-1} - X_n)$$

our goal is to show that this converges to μ if the following two conditions hold

$$(C1) \sum_{n \geq 1} \alpha_n = \infty$$

$$(C2) \sum_{n \geq 1} \alpha_n^2 < \infty$$

What's the intuition behind these assumptions?

Blank page for notes



Blank page for notes



Blank page for notes



Blank page for notes



Stochastic gradient descent

- ▶ Suppose we are interested in estimating

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} F(\boldsymbol{\theta}),$$

where $F(\boldsymbol{\theta}) = Pf(\mathbf{Z}; \boldsymbol{\theta})$

- ▶ Canonical example $\mathbf{Z} = (\mathbf{X}, Y)$, $f(\mathbf{x}, y; \boldsymbol{\beta}) = (y - \mathbf{x}^\top \boldsymbol{\beta})^2$ and

$$F(\boldsymbol{\beta}) = P(Y - \mathbf{X}^\top \boldsymbol{\beta})^2$$

and $\boldsymbol{\beta}^*$ is the usual population-level least-squares soln

Vanilla SGD

- ▶ Many flavors/variants of SGD but most build on

$$\hat{\theta}_i = \hat{\theta}_{i-1} - \alpha_i \nabla_{\mathbf{z}} f(\mathbf{Z}_i; \hat{\theta}_{i-1})$$

which processes a single observation at time

- ▶ Average SGD averages iterates at the end of data-stream, i.e.,

$$\bar{\theta}_n = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_i$$

- ▶ In all cases, require (C1) and (C2)

Vanilla SGD with linear model

- ▶ Observed data $\{(\mathbf{X}_i Y_i)\}_{i=1}^n$ and

$$\hat{\boldsymbol{\beta}}_i = \hat{\boldsymbol{\beta}}_{i-1} + \alpha_i (Y_i - \mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_{i-1}) \mathbf{X}_i,$$

when $\alpha_i \equiv \alpha$, this is known as the Widrow-Hoff algorithm⁸

- ▶ To `widrow_hoff.R`

⁸After its inventors: Widrow Wilson and Hof Falgorithm.



Evaluating Widrow-Hoff

- ▶ WH offers computational/memory savings: at what cost?
- ▶ Baseline for comparison: using all the data at once
 - ▶ Full data: $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$
 - ▶ Given sequence of estimators $\{\hat{\boldsymbol{\beta}}_i\}_{i=1}^n$ define the cumulative squared-error loss

$$L_n\left(\{\hat{\boldsymbol{\beta}}_i\}_{i=1}^n\right) = \sum_{i=1}^n \left(Y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}_i\right)^2,$$

i.e., how good are our predictions if we use them online

- ▶ For any $\boldsymbol{\beta} \in \mathbb{R}^p$ define $L_n(\boldsymbol{\beta}) = L_n(\{\boldsymbol{\beta}, \boldsymbol{\beta}, \dots, \boldsymbol{\beta}\})$ and note that the OLS estimator is given by $\arg \min_{\boldsymbol{\beta}} L_n(\boldsymbol{\beta})$

Widrow-Hoff is pretty good

Theorem

Assume $\|\mathbf{X}\| \leq 1$ with probability one and let $\hat{\beta}_i^{\text{WH}}$ denote the Widrow-Hoff estimator at iteration $i = 1, \dots, n$. If we initialize the algorithm at $\hat{\beta}_1 \equiv 0$ then

$$L_n \left(\left\{ \hat{\beta}_i^{\text{WH}} \right\}_{i=1}^n \right) \leq \min_{\beta} \left(\frac{L_n(\beta)}{1 - \alpha} + \frac{\|\beta\|^2}{\alpha} \right)$$

Preamble to proof

- ▶ For $\beta \in \mathbb{R}^p$ arbitrary define
 - ▶ $e_i(\beta) = \|\hat{\beta}_i^{\text{WH}} - \beta\|^2$
 - ▶ $g_i(\beta) = Y_i - \mathbf{X}_i^\top \beta$ (Residual at β at round i)
 - ▶ $r_i = g_i(\hat{\beta}_i^{\text{WH}}) = Y_i - \mathbf{X}_i^\top \hat{\beta}_i^{\text{WH}}$ (Residual of WH at round i)
 - ▶ $\Delta_i = -\alpha(Y_i - \mathbf{X}_i^\top \hat{\beta}_i^{\text{WH}})\mathbf{X}_i = -\alpha r_i \mathbf{X}_i$ (update at round i)

Lemma

For any $\beta \in \mathbb{R}^p$ and $i = 1, \dots, n$

$$e_{i+1}(\beta) - e_i(\beta) \leq \frac{\alpha}{1 - \alpha} g_i^2(\beta) - \alpha r_i^2.$$

Blank page for lemma proof



Blank page for lemma proof



Blank page for WH error bound proof



Blank page WH error bound proof



WH with averaging

- For each $i = 1, \dots, n$ define the average WH estimator

$$\bar{\beta}_i^{\text{WH}} = \frac{1}{i} \sum_{j=1}^i \hat{\beta}_j^{\text{WH}}$$

Theorem

For any n it follows that

$$\mathbb{E} \left(Y - \mathbf{X}^\top \bar{\beta}_n^{\text{WH}} \right)^2 \leq \min_{\beta \in \mathbb{R}^p} \left\{ \frac{P(Y - \mathbf{X}^\top \beta)^2}{1 - \alpha} + \frac{\|\beta\|^2}{\alpha n} \right\}$$

Blank page for ave WH error bound proof



Blank page ave WH error bound proof



SGD and statistical inference

- Suppose that β^* satisfies

$$\arg \min_{\beta} Pf(\mathbf{Z}; \beta),$$

e.g., $\mathbf{Z} = (\mathbf{X}, Y)$ and $f(\mathbf{z}; \beta) = (y - \mathbf{x}^\top \beta)^2$ etc.

- Consider ASGD $\hat{\beta}_i = \hat{\beta}_{i-1} - \alpha_i \nabla_{\beta} f(\mathbf{Z}_i; \hat{\beta}_{i-1})$ and

$$\bar{\beta}_n = \frac{1}{n} \sum_{i=1}^n \hat{\beta}_i$$

SGD and statistical inference cont'd

- Define

$$\Sigma \triangleq \nabla_{\beta}^2 P f(\mathbf{Z}; \beta^*)$$

$$\Omega \triangleq P \nabla_{\beta} f(\mathbf{Z}; \beta^*) \nabla_{\beta} f(\mathbf{Z}; \beta^*)^{\top}$$

- Polyak and Juditsky (1992) show that if $\beta \mapsto P f(\mathbf{Z}; \beta)$ is strongly convex with a smooth gradient and conditions (C1) and (C2) on $\{\alpha_i\}_{i \geq 1}$ then

$$\sqrt{n}(\bar{\beta}_n - \beta^*) \rightsquigarrow \text{Normal}(0, \Sigma^{-1} \Omega \Sigma^{-1}).$$

*Proof is involved so we'll be skipping it (paper linked on website). How can you convince yourself this is true without any math?

SGD and statistical inference: IRL

- ▶ How to estimate Ω and Σ ?
 - ▶ Scandalous approach: sweep through the data twice (gasp!) and compute

$$\hat{\Sigma}_n = \mathbb{P}_n \nabla_{\beta}^2 f(\mathbf{Z}; \bar{\beta}_n)$$

$$\hat{\Omega}_n = \mathbb{P}_n \nabla_{\beta} f(\mathbf{Z}; \bar{\beta}_n) \nabla_{\beta} f(\mathbf{Z}; \bar{\beta}_n)^{\top}$$

- ▶ One-sweep approach: use running estimator of β^*

$$\tilde{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n \nabla_{\beta}^2 f(\mathbf{Z}_i; \hat{\beta}_i)$$

$$\tilde{\Omega}_n = \frac{1}{n} \sum_{i=1}^n \nabla_{\beta} f(\mathbf{Z}_i; \hat{\beta}_i) \nabla_{\beta} f(\mathbf{Z}_i; \hat{\beta}_i)^{\top}$$

Note: one may need to use a generalized inverse etc. when inverting $\tilde{\Sigma}_n$

Back to `widrow_hoff.R`



Thank you.

`eric.laber@duke.edu`

`laber-labs.com`

