

Randomized Sketching for Convex and Non-Convex Optimization

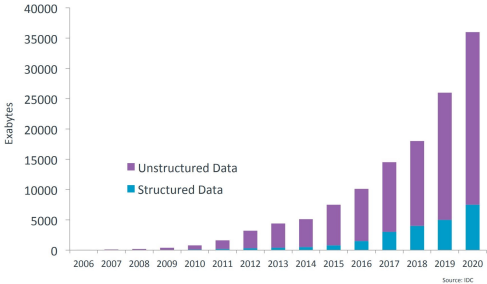
Mert Pilanci

Department of Electrical Engineering

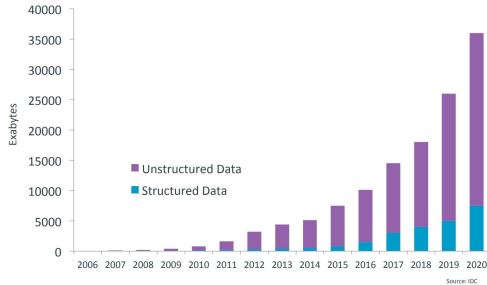
Stanford University

December 6, 2018

Scale of data

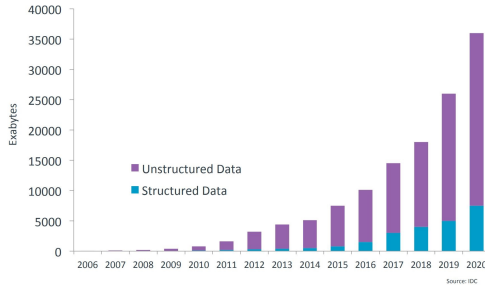


Scale of data



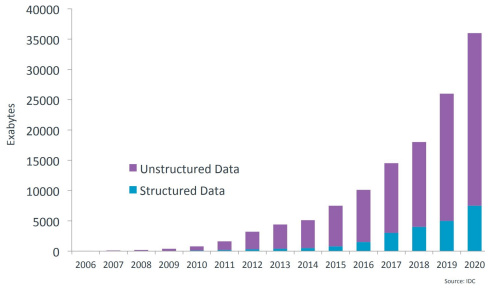
- Every day, we create 2.5 billion gigabytes of data

Scale of data



- Every day, we create 2.5 billion gigabytes of data
- Data stored grows 4x faster than world economy (Mayer-Schonberger)

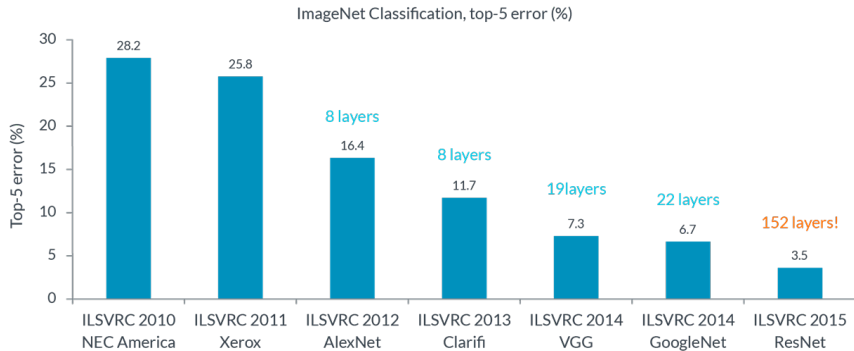
Scale of data




- Every day, we create 2.5 billion gigabytes of data
- Data stored grows 4x faster than world economy (Mayer-Schonberger)



Deep learning revolution



- Machine learning
- Statistical estimation and data analysis
- Signal processing and control theory
- Computational imaging
- Design and manufacturing
- Decision making



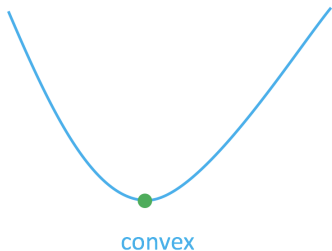
minimize $g(x)$ subject to constraints

Machine learning and statistics

- More data points reduce sampling error, higher significance
→ Large scale optimization problems

Machine learning and statistics

- More data points reduce sampling error, higher significance
→ Large scale optimization problems
- Complex models can improve accuracy
→ Non-convex optimization problems harder to solve as dimensions grow



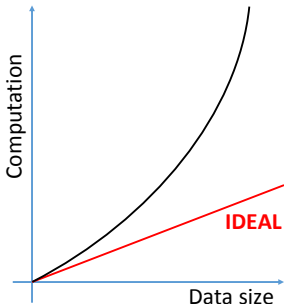
Machine learning and statistics

- More data points reduce sampling error, higher significance
→ Large scale optimization problems
- Complex models can improve accuracy
→ Non-convex optimization problems harder to solve as dimensions grow

What if we could reduce the **data volume without losing any significant information ?**

Machine learning and statistics

- More data points reduce sampling error, higher significance
→ Large scale optimization problems
- Complex models can improve accuracy
→ Non-convex optimization problems harder to solve as dimensions grow



- Goals:**
1. Find optimal trade-offs between **computation** and **accuracy**
 2. Leverage **distributed computation**
 3. Tackle **non-convexity**

Optimization and Big Data

Sketching

Distributed Sketching

Non-convex Problems

Ongoing work

minimize $f(Ax)$ subject to $x \in \mathcal{C}$

- Data matrix $A \in \mathbb{R}^{n \times d}$ is extremely large

$$\text{minimize } f(Ax) \quad \text{subject to } x \in \mathcal{C}$$

- Data matrix $A \in \mathbb{R}^{n \times d}$ is extremely large

Examples:

- Airline dataset (120GB) $n = 120 \times 10^6$, $d = 28$
Flight arrival and departure details from 1987 to 2008

$$\text{minimize } f(Ax) \quad \text{subject to } x \in \mathcal{C}$$

- Data matrix $A \in \mathbb{R}^{n \times d}$ is extremely large

Examples:

- Airline dataset (120GB) $n = 120 \times 10^6$, $d = 28$
Flight arrival and departure details from 1987 to 2008
- Imagenet dataset (1.31TB) $n = 14 \times 10^6$, $d = 2 \times 10^5$
14 Million images for visual recognition

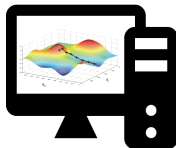
[US Department of Transportation]

[Deng et al. 2009]

DATA



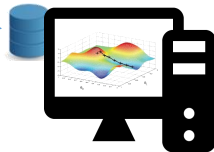
OPTIMIZER



DATA



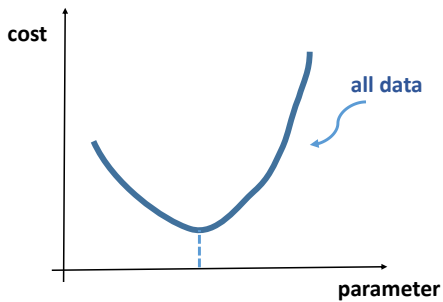
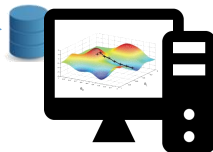
OPTIMIZER



DATA



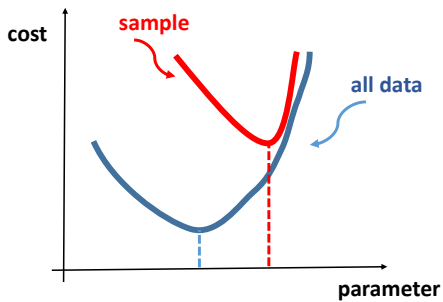
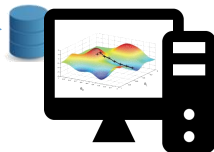
OPTIMIZER



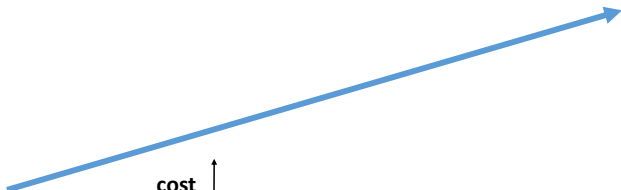
DATA



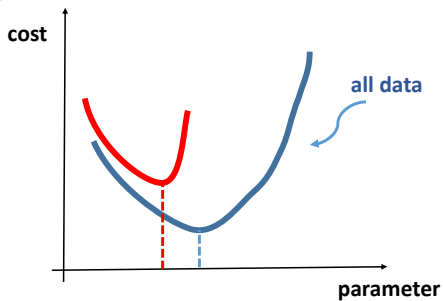
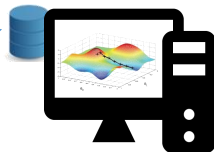
OPTIMIZER



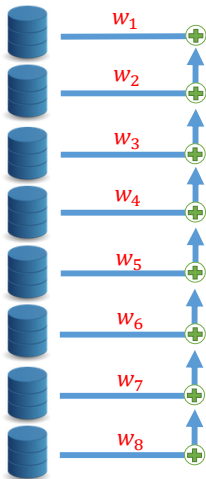
DATA



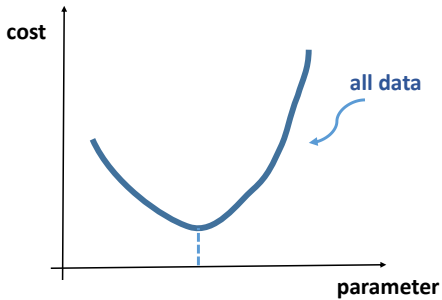
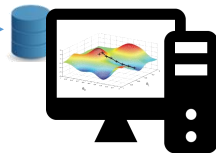
OPTIMIZER



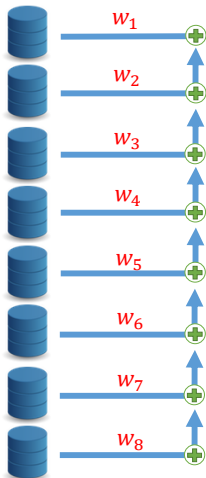
DATA



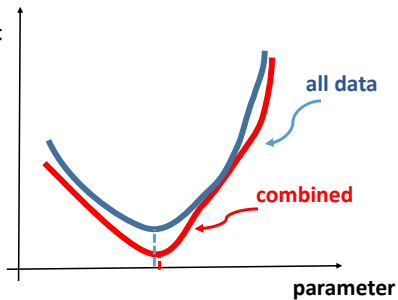
OPTIMIZER



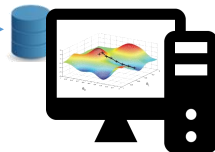
DATA



cost



OPTIMIZER



Example : Least squares prediction

- $A : n \times d$ feature matrix, and $y : n \times 1$ response vector
- Original $\mathbf{OPT} = \min_{x \in \mathcal{C}} \underbrace{\|Ax - y\|^2}_{f(Ax)}$

Example : Least squares prediction

○ $A : n \times d$ feature matrix, and $y : n \times 1$ response vector

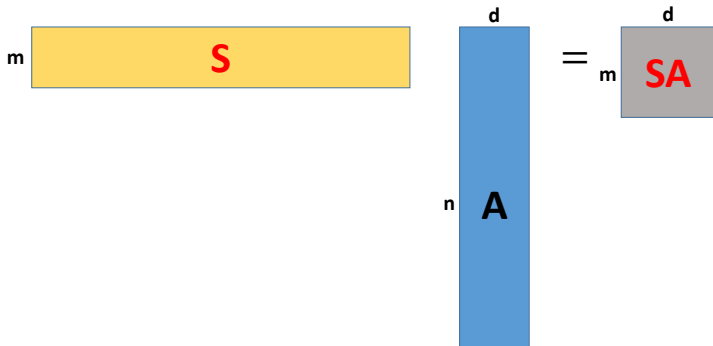
○ Original $\text{OPT} = \min_{x \in \mathcal{C}} \underbrace{\|Ax - y\|^2}_{f(Ax)}$



Example : Least squares prediction

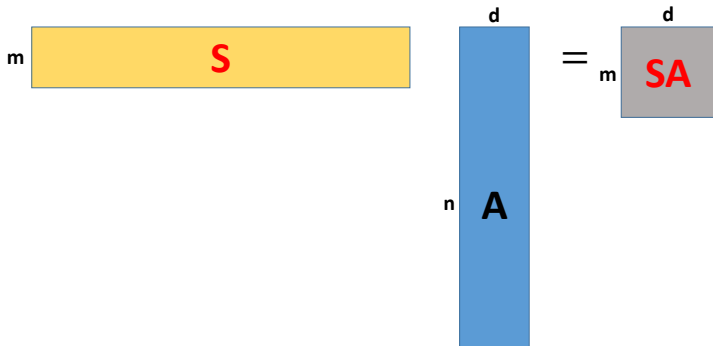
○ $A : n \times d$ feature matrix, and $y : n \times 1$ response vector

○ Original $\text{OPT} = \min_{x \in \mathcal{C}} \underbrace{\|Ax - y\|^2}_{f(Ax)}$



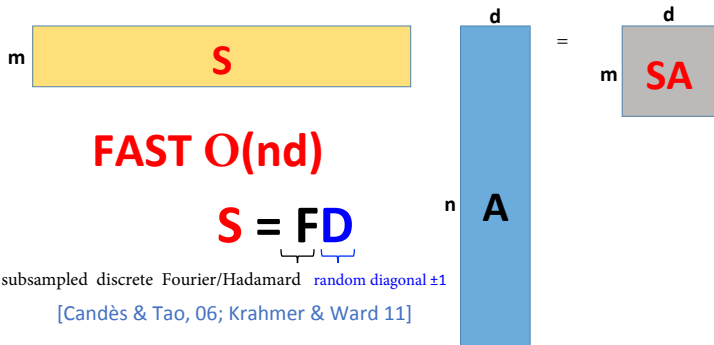
Example : Least squares prediction

- $A : n \times d$ feature matrix, and $y : n \times 1$ response vector
- Original $\text{OPT} = \min_{x \in \mathcal{C}} \underbrace{\|Ax - y\|^2}_{f(Ax)}$
- Approximate $\hat{x} = \arg \min_{x \in \mathcal{C}} \|S(Ax - y)\|^2$
- $S : m \times n$ *sketching* matrix (e.g., i.i.d. ± 1 random matrix)



Example : Least squares prediction

- $A : n \times d$ feature matrix, and $y : n \times 1$ response vector
- Original $\text{OPT} = \min_{x \in \mathcal{C}} \underbrace{\|Ax - y\|^2}_{f(Ax)}$
- Approximate $\hat{x} = \arg \min_{x \in \mathcal{C}} \|S(Ax - y)\|^2$
- $S : m \times n$ *sketching* matrix (e.g., i.i.d. ± 1 random matrix)



Example : Least squares prediction

- $A : n \times d$ feature matrix, and $y : n \times 1$ response vector
- Original $\text{OPT} = \min_{x \in \mathcal{C}} \underbrace{\|Ax - y\|^2}_{f(Ax)}$
- Approximate $\hat{x} = \arg \min_{x \in \mathcal{C}} \|S(Ax - y)\|^2$
- $S : m \times n$ *sketching* matrix (e.g., i.i.d. ± 1 random matrix)

Theorem : Cost approximation

If $m \geq 2 \text{rank}^*(A)/\epsilon$, then

$$\text{OPT} \leq f(A\hat{x}) \leq (1 + \epsilon)\text{OPT}$$

with high probability

Example : Least squares prediction

- $A : n \times d$ feature matrix, and $y : n \times 1$ response vector
- Original $\mathbf{OPT} = \min_{x \in \mathcal{C}} \underbrace{\|Ax - y\|^2}_{f(Ax)}$
- Approximate $\hat{x} = \arg \min_{x \in \mathcal{C}} \|S(Ax - y)\|^2$
- $S : m \times n$ *sketching* matrix (e.g., i.i.d. ± 1 random matrix)

Theorem : Converse

If $m \leq c_0 \text{rank}^*(A)/\epsilon$, then
 $f(A\hat{x}) \geq (1 + \epsilon)\mathbf{OPT}$
with probability $> \frac{1}{2}$

Example : Least squares prediction

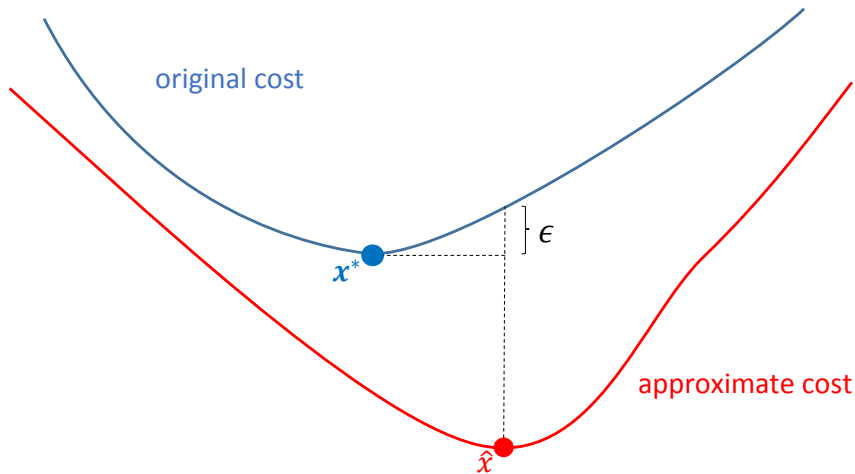
Practical use

Airline dataset $n = 120,000,000$, $d = 28$

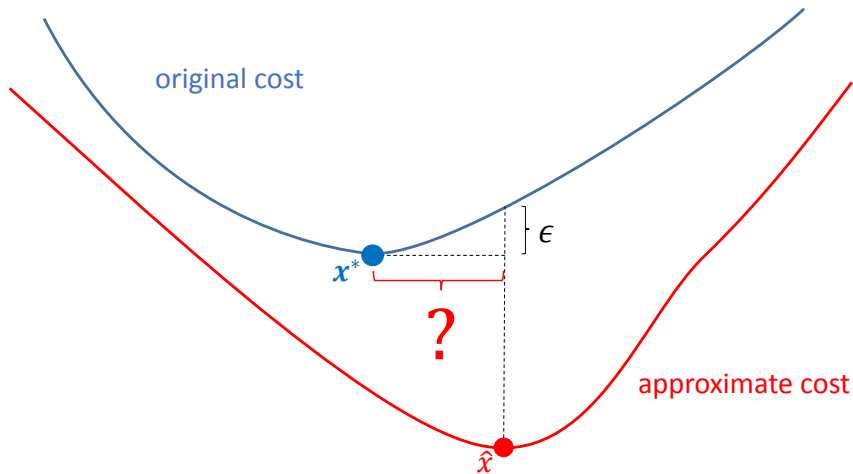
$m = 500$ gives 1.1-approximation

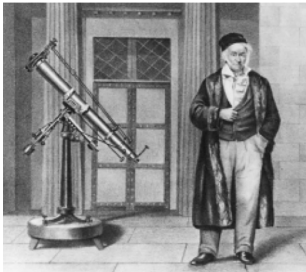
$m = 5000$ gives 1.01-approximation

$$\min_x \|Ax - y\|_2^2 \quad \text{and} \quad \min_x \|S(Ax - y)\|_2^2$$



$$\min_x \|Ax - y\|_2^2 \quad \text{and} \quad \min_x \|S(Ax - y)\|_2^2$$

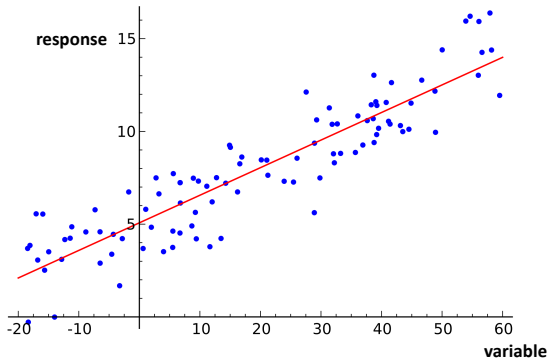




Least squares

$$\min_x \|Ax - y\|^2$$

[Gauss, 1795]



- Consider the noisy observation model

$$y = Ax^* + w, \quad \text{where } w_1, w_2 \dots \sim N(0, \sigma^2)$$

- Consider the noisy observation model

$$y = Ax^* + w, \quad \text{where } w_1, w_2, \dots \sim N(0, \sigma^2)$$

- estimation error:

$$\mathbb{E} \|x_{LS} - x^*\|_2^2 = O\left(\frac{d}{n}\sigma^2\right)$$

Statistical error

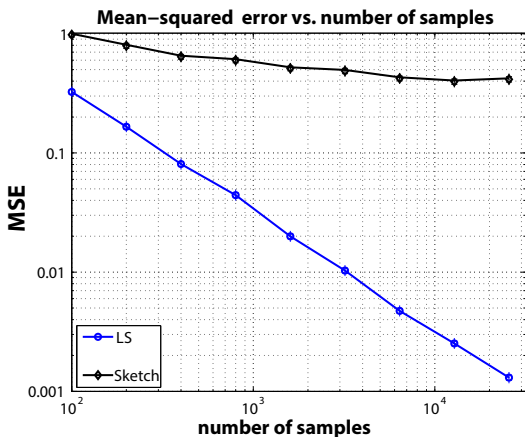
- Observation model $y = Ax^* + w$, where $w \sim N(0, \sigma^2 I_n)$
- Is the sketched solution \hat{x} statistically optimal?

$$\hat{x} = \arg \min_x \|SAx - Sy\|^2 \quad \text{where } S \in \mathbb{R}^{m \times n}$$

Statistical error

- Observation model $y = Ax^* + w$, where $w \sim N(0, \sigma^2 I_n)$
- Is the sketched solution \hat{x} statistically optimal?

$$\hat{x} = \arg \min_x \|SAx - Sy\|^2 \quad \text{where } S \in \mathbb{R}^{m \times n}$$



Suboptimality of the Classical Sketch

- Is the sketched solution \hat{x} statistically optimal?

Suboptimality of the Classical Sketch

- Is the sketched solution \hat{x} statistically optimal?
- **No**, information-theoretically impossible!

Suboptimality of the Classical Sketch

- Is the sketched solution \hat{x} statistically optimal?
- **No**, information-theoretically impossible!
- For Gaussian, i.i.d. ± 1 or DFT/Hadamard based S and $m = \text{constant} \times \text{rank}(A)$

Theorem

Any estimator that is a function of (SA, Sy) obeys

$$\mathbb{E}_{S,w} [\|\hat{x} - x^*\|^2] \gtrsim \sigma^2$$

Suboptimality of the Classical Sketch

- Is the sketched solution \hat{x} statistically optimal?
- **No**, information-theoretically impossible!
- For Gaussian, i.i.d. ± 1 or DFT/Hadamard based S and $m = \text{constant} \times \text{rank}(A)$

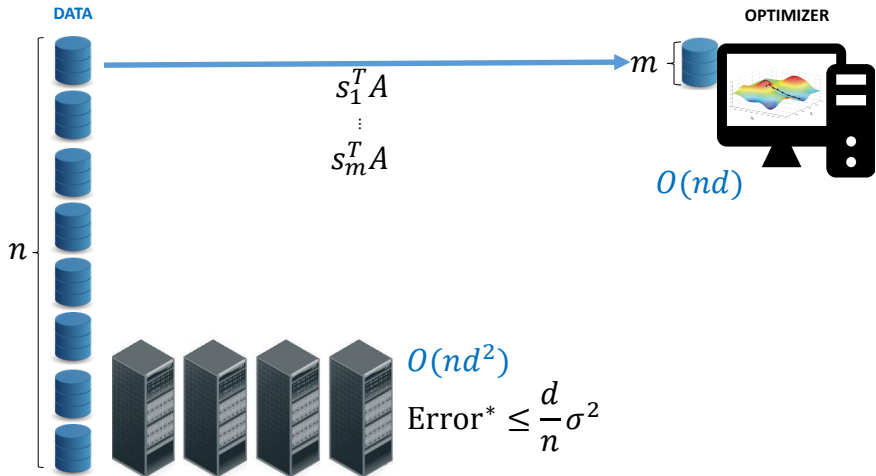
Theorem

Any estimator that is a function of (SA, Sy) obeys

$$\mathbb{E}_{S,w} [\|\hat{x} - x^*\|^2] \gtrsim \sigma^2$$

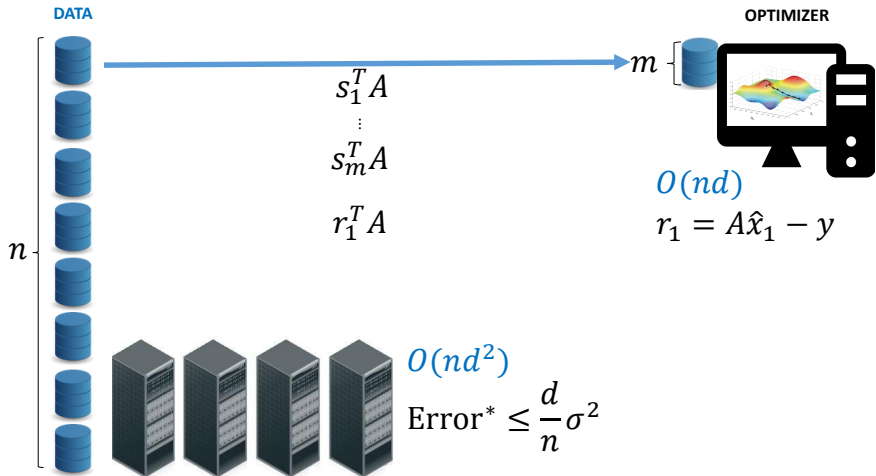
	statistical error	computation
sketch	suboptimal	$O(nd)$
original	optimal	$O(nd^2)$

Is there an **optimal** algorithm with complexity $O(nd)$?



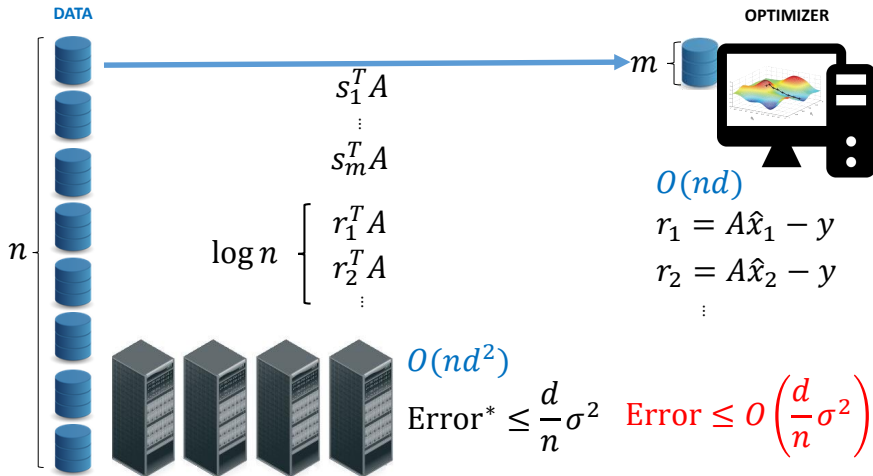
Is there an **optimal** algorithm with complexity $O(nd)$?

YES!

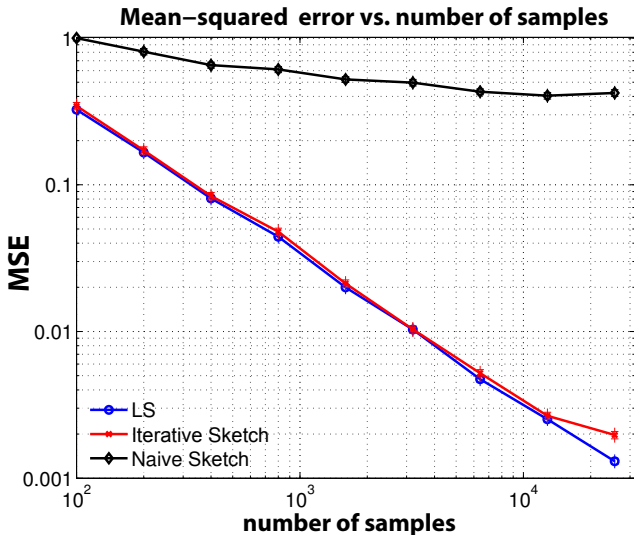


Is there an **optimal** algorithm with complexity $O(nd)$?

YES!



Numerical Simulation: Estimation Error



Iterative Sketch

$$\begin{array}{c} SA \\ r_1^T A \\ \vdots \\ r_t^T A \end{array}$$

Iterative Sketch

$$\begin{array}{c} SA \\ r_1^T A \\ \vdots \\ r_t^T A \end{array}$$

- Statistical model : $y = Ax^* + w$ where $x^* \in \mathcal{C}$

Theorem (Optimality)

Iterative Sketch achieves optimal prediction error with $\log(n/d)$ iterations for any convex set \mathcal{C}

	statistical error	computation
sketch	suboptimal	$O(nd)$
original	optimal	$O(nd^2)$
iterative sketch	optimal	$O(nd)$

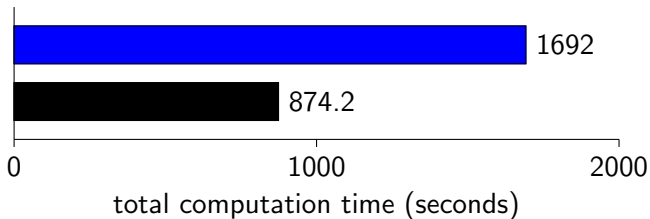
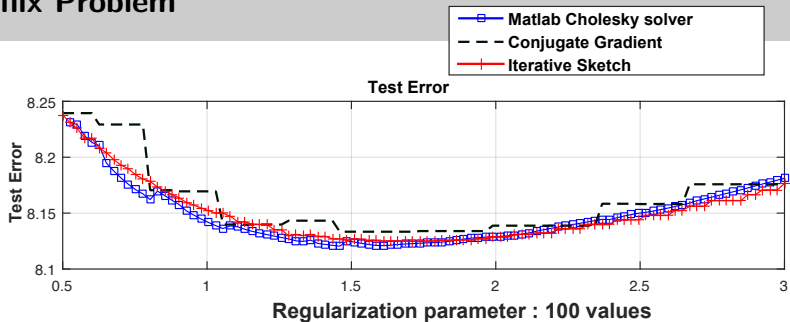
Netflix Problem

- 500000 × 17000 matrix A of ratings (users × movies)
- Predict the ratings of a particular movie
- Least-squares regression with ℓ_2 regularization

$$\min_x \|Ax - y\|^2 + \lambda \|x\|_2^2$$

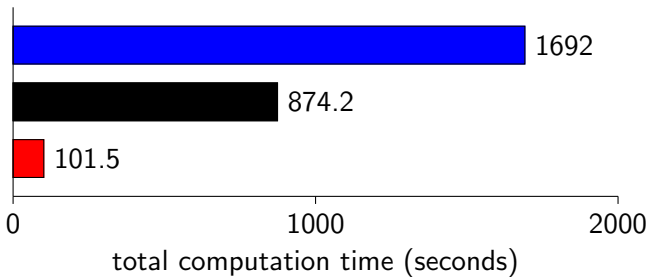
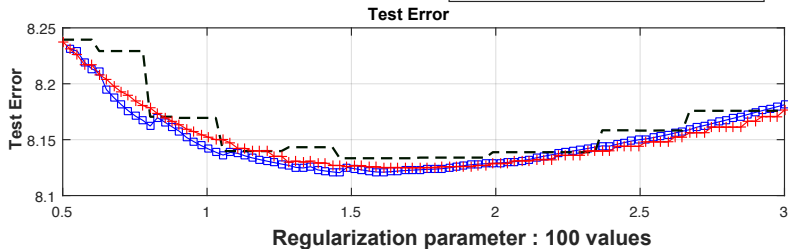
- Partition into test and training sets, solve for all values of $\lambda \in \{1, 2, \dots, 100\}$.

Netflix Problem



Netflix Problem

—○— Matlab Cholesky solver
- - - Conjugate Gradient
—+— Iterative Sketch

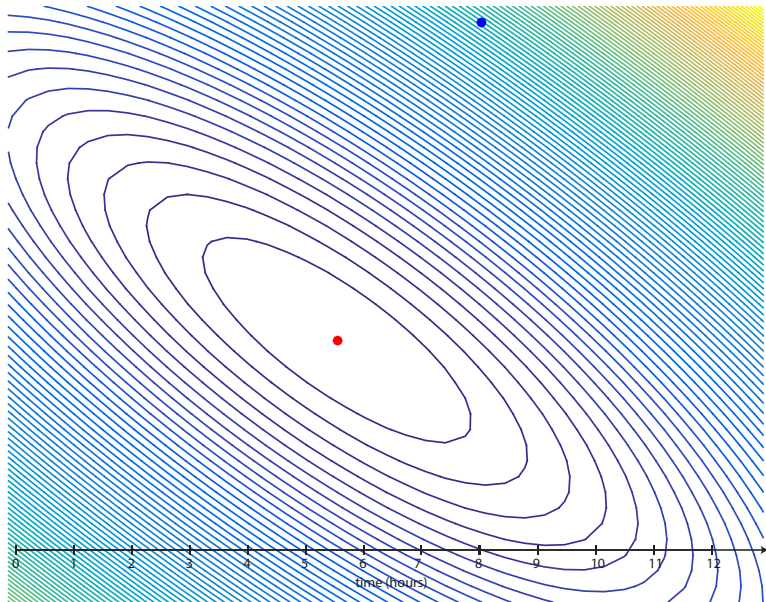


How to generalize to arbitrary functions ?

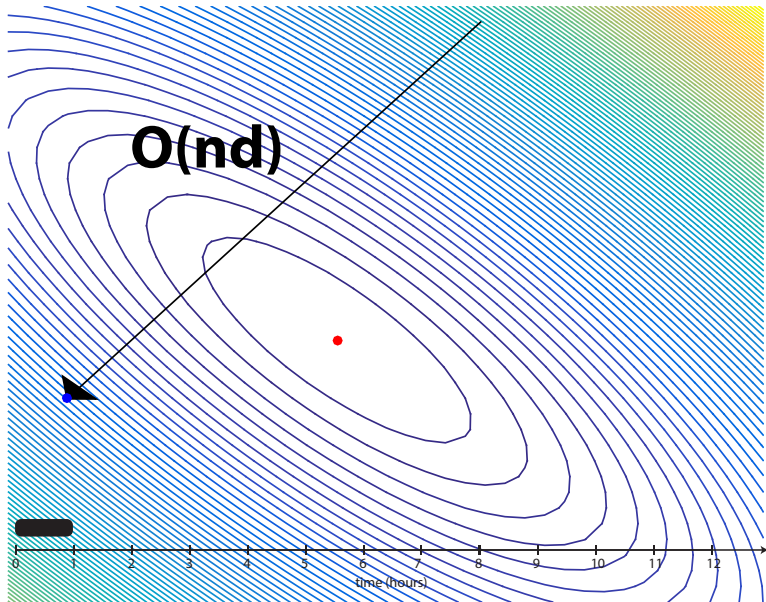
Convex objective, where $A \in \mathbb{R}^{n \times d}$ is a large data matrix

$$x^* = \arg \min_{x \in \mathcal{C}} f(Ax)$$

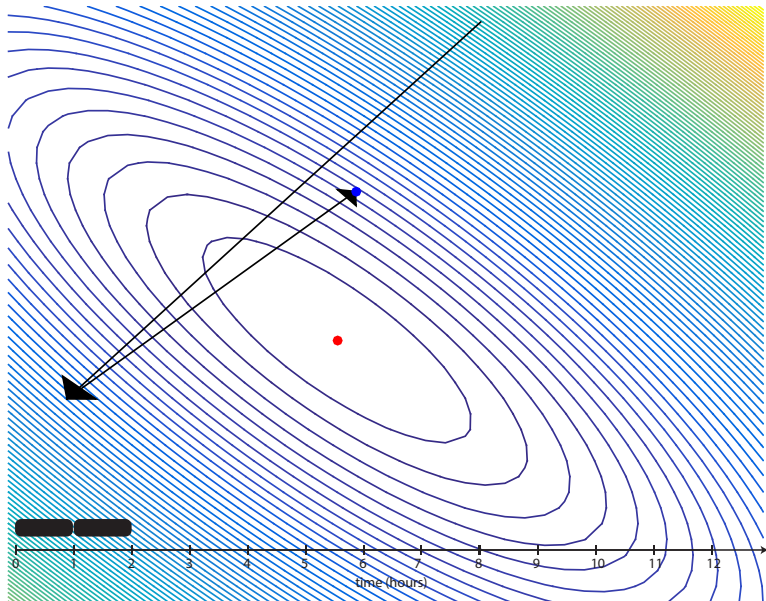
Gradient Descent



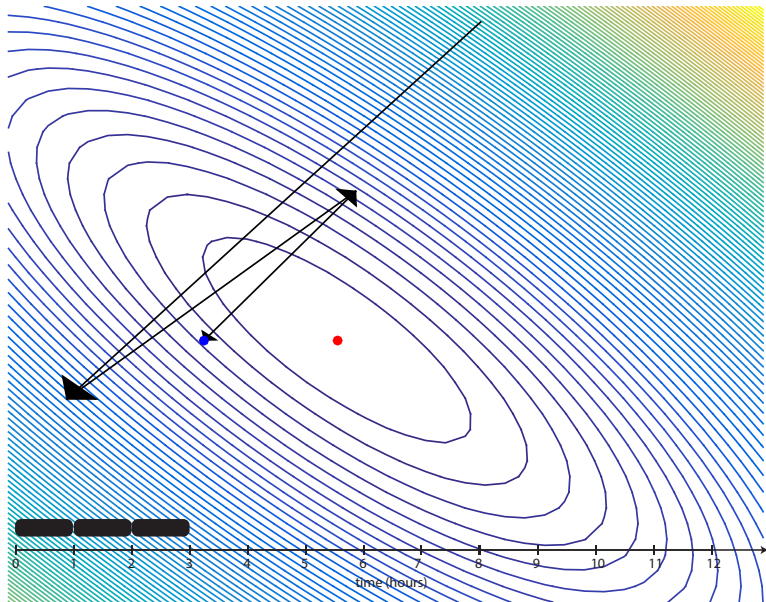
Gradient Descent



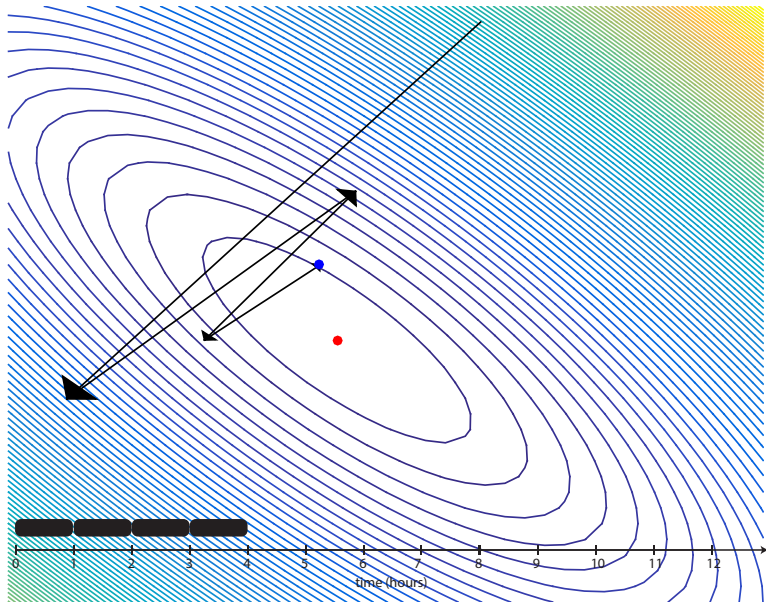
Gradient Descent



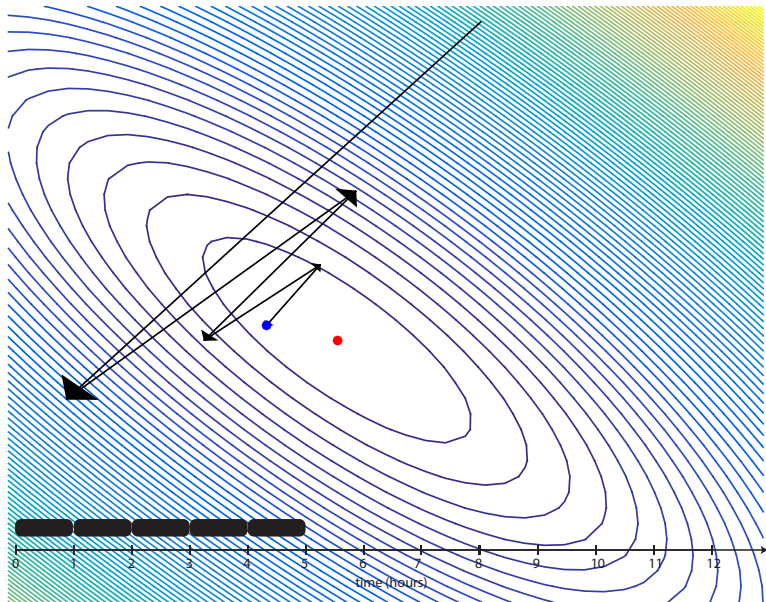
Gradient Descent



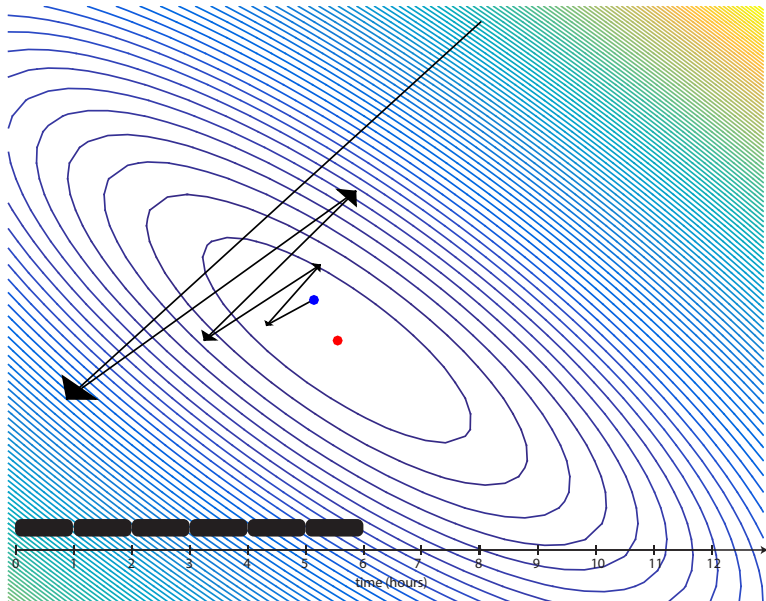
Gradient Descent



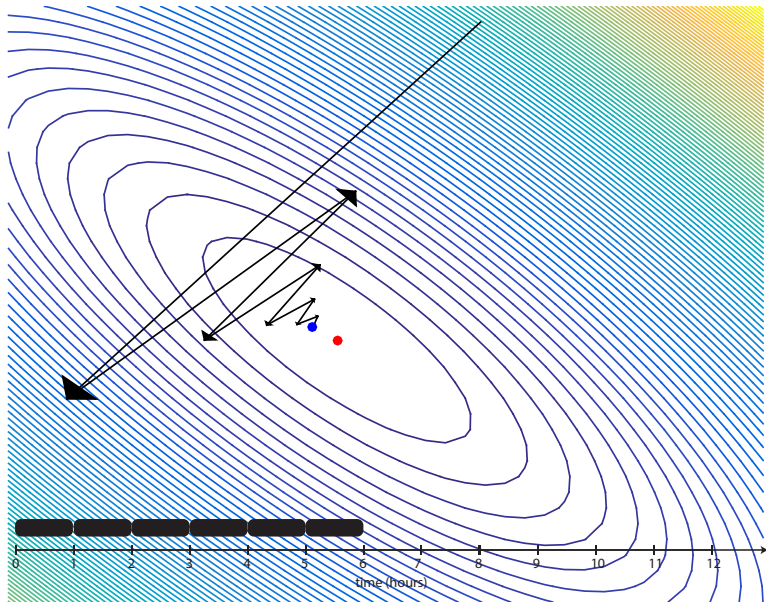
Gradient Descent



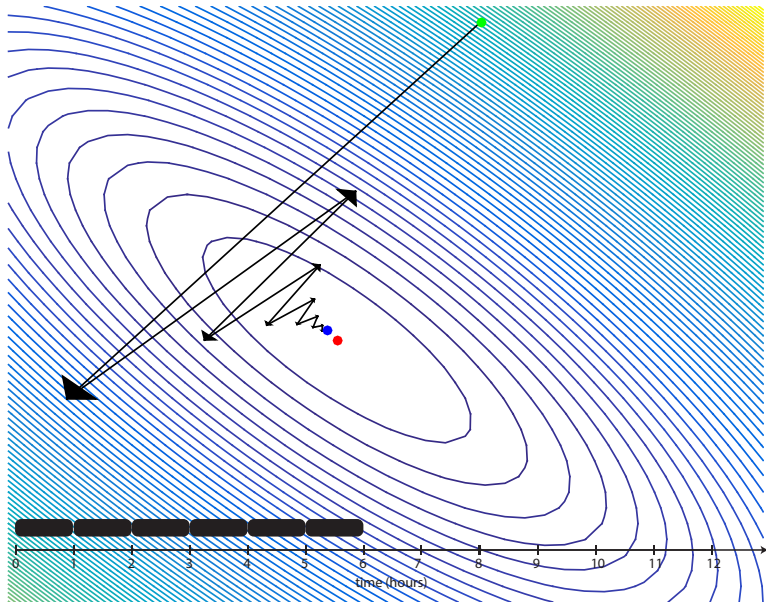
Gradient Descent



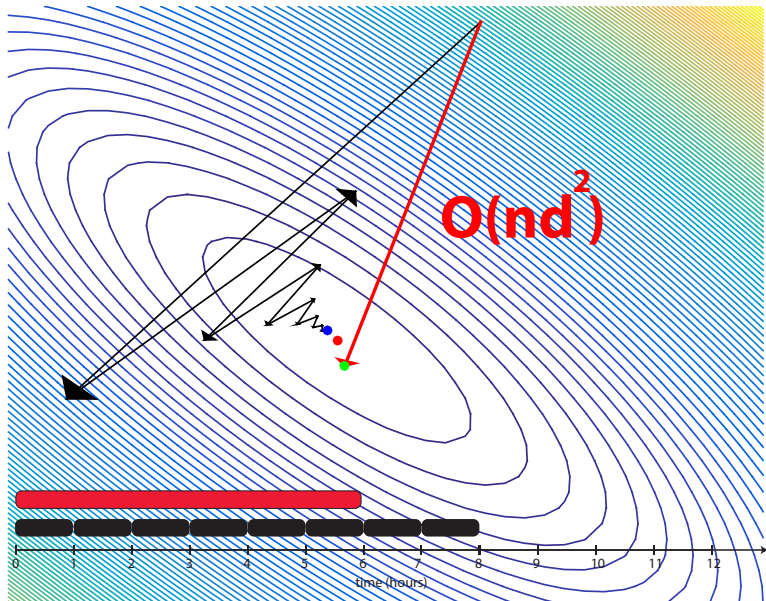
Gradient Descent vs



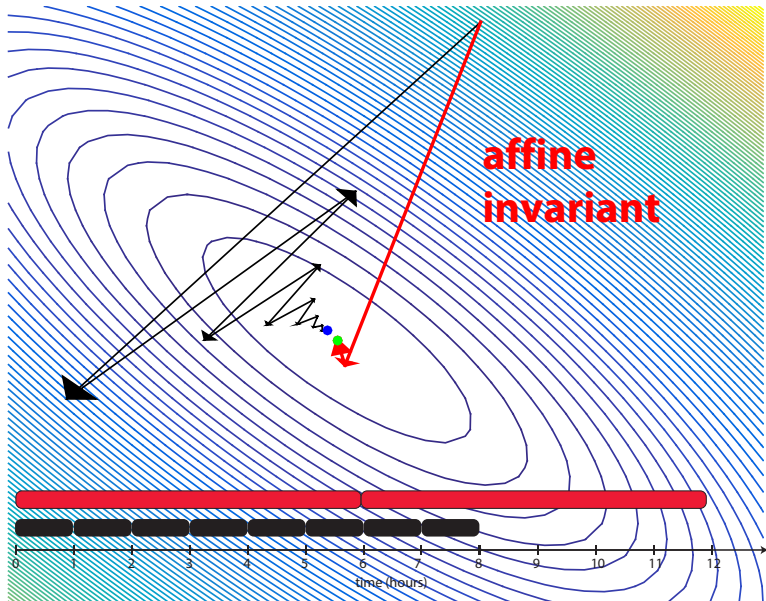
Gradient Descent vs Newton's Method



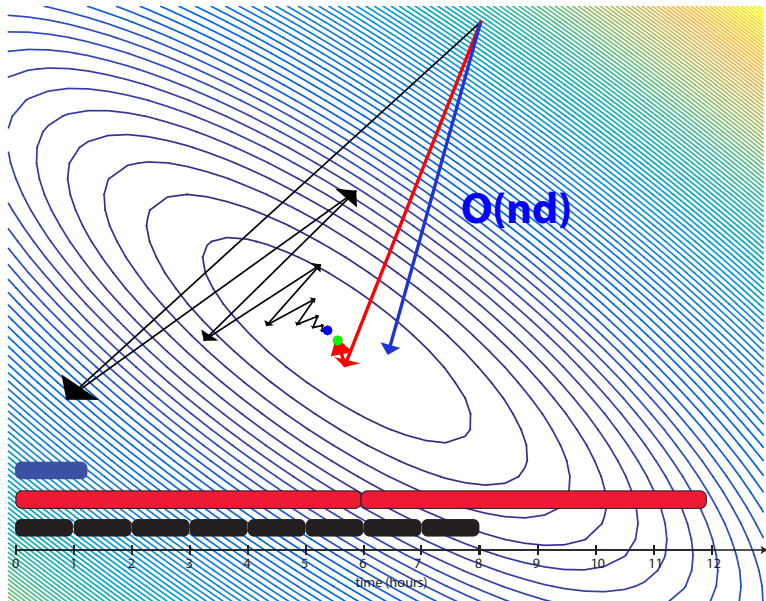
Gradient Descent vs Newton's Method



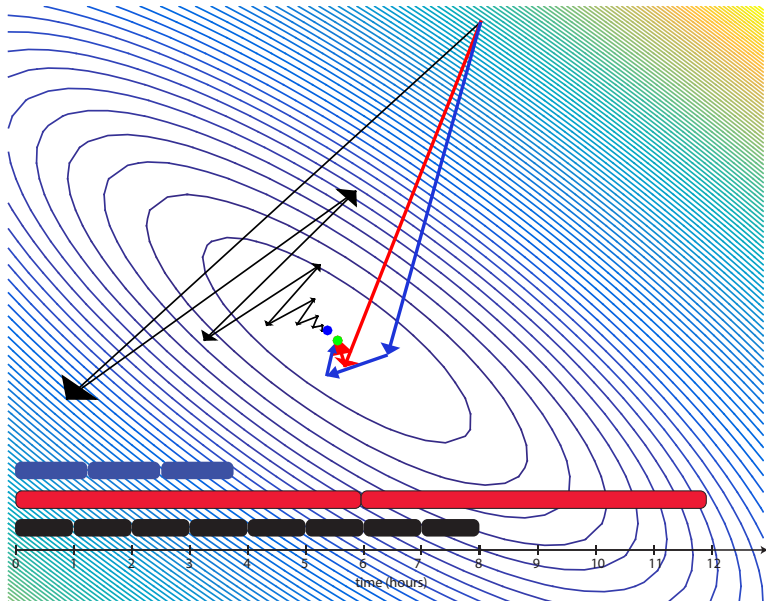
Gradient Descent vs Newton's Method



Gradient Descent vs Newton's Method



Gradient Descent vs Newton's Method



- Newton's Method

$$x^{t+1} = \arg \min_{x \in \mathcal{C}} \langle \nabla g(x^t), x - x^t \rangle + \frac{1}{2} \|\nabla^2 g(x^t)^{1/2} (x - x^t)\|^2$$

Introducing Newton Sketch

- Newton's Method

$$x^{t+1} = \arg \min_{x \in \mathcal{C}} \langle \nabla g(x^t), x - x^t \rangle + \frac{1}{2} \|\nabla^2 g(x^t)^{1/2} (x - x^t)\|^2$$

Definition (Newton Sketch)

$$x^{t+1} = \arg \min_{x \in \mathcal{C}} \langle \nabla g(x^t), x - x^t \rangle + \frac{1}{2} \|S^t \nabla^2 g(x^t)^{1/2} (x - x^t)\|^2$$

Introducing Newton Sketch

- Newton's Method

$$x^{t+1} = \arg \min_{x \in \mathcal{C}} \langle \nabla g(x^t), x - x^t \rangle + \frac{1}{2} \|\nabla^2 g(x^t)^{1/2} (x - x^t)\|^2$$

Definition (Newton Sketch)

$$x^{t+1} = \arg \min_{x \in \mathcal{C}} \langle \nabla g(x^t), x - x^t \rangle + \frac{1}{2} \|S^t \nabla^2 g(x^t)^{1/2} (x - x^t)\|^2$$

- Iterative Sketch is a special case $g(x) = \|Ax - y\|^2$

Convergence of Newton Sketch

Theorem

*Newton Sketch is **affine invariant** in distribution and the number of iterations for ϵ accuracy is less than*

$$C \log(1/\epsilon)$$

*C is a constant **independent** of f & data (same assumptions).*

[P. and Wainwright. SIAM Journal on Optimization, 2017]

Convergence of Newton Sketch

Theorem

Newton Sketch is **affine invariant** in distribution and the number of iterations for ϵ accuracy is less than

$$C \log(1/\epsilon)$$

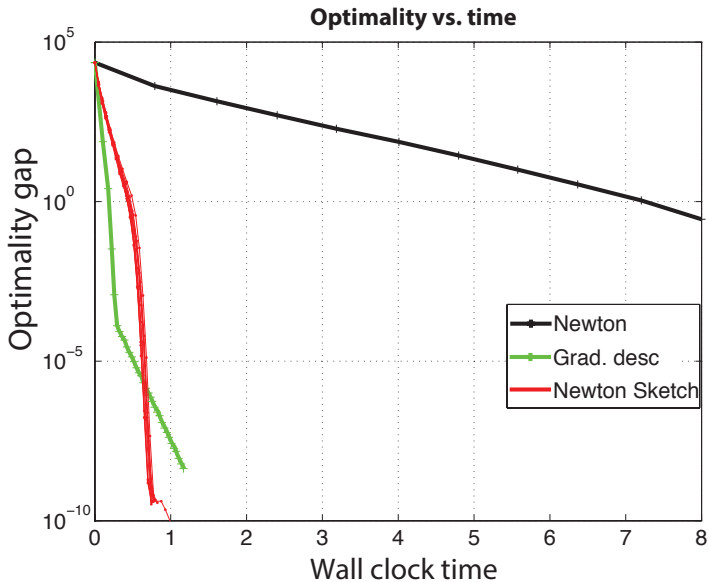
C is a constant **independent** of f & data (same assumptions).

[P. and Wainwright. SIAM Journal on Optimization, 2017]

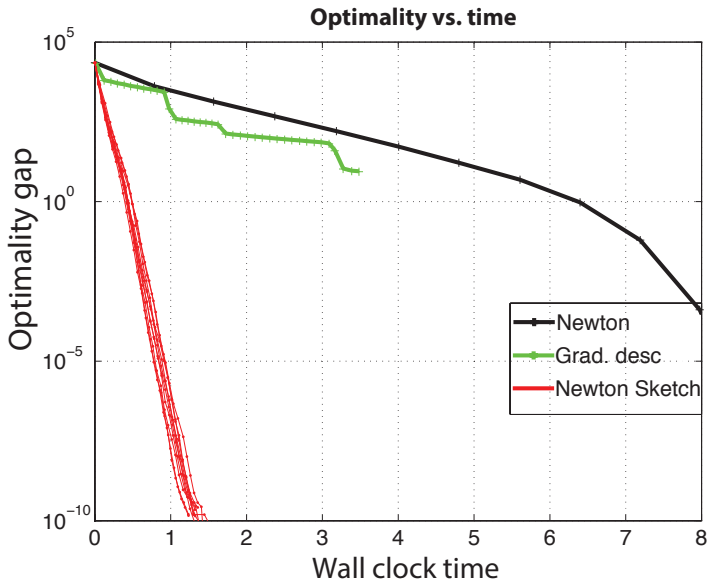
Gradient Descent	computation $O(\kappa nd \log(1/\epsilon))$
Newton's Method	$O(nd^2 \log \log(1/\epsilon))$
Newton Sketch	$O(nd \log(1/\epsilon))$

Dependence on curvature κ is **unavoidable** among first order methods [Nesterov, 04]

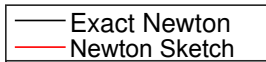
Logistic Regression ($n = 500,000$, $d = 5,000$ uncorrelated)



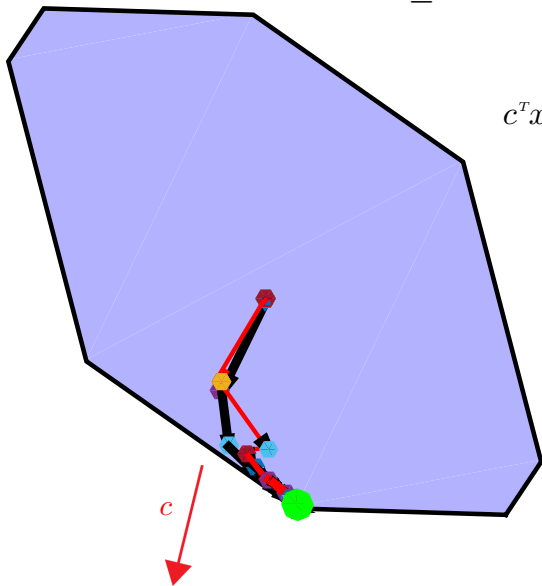
Logistic Regression ($n = 500,000$, $d = 5,000$ correlation 0.1)



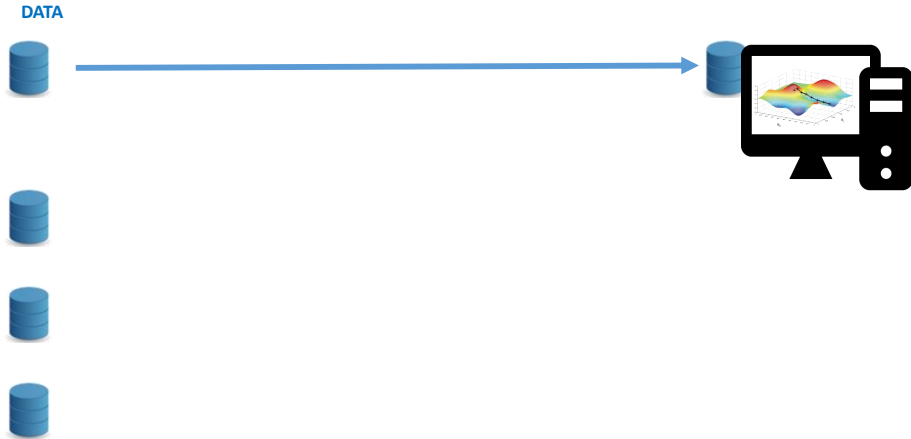
$$\min c^T x \\ Ax \leq b$$



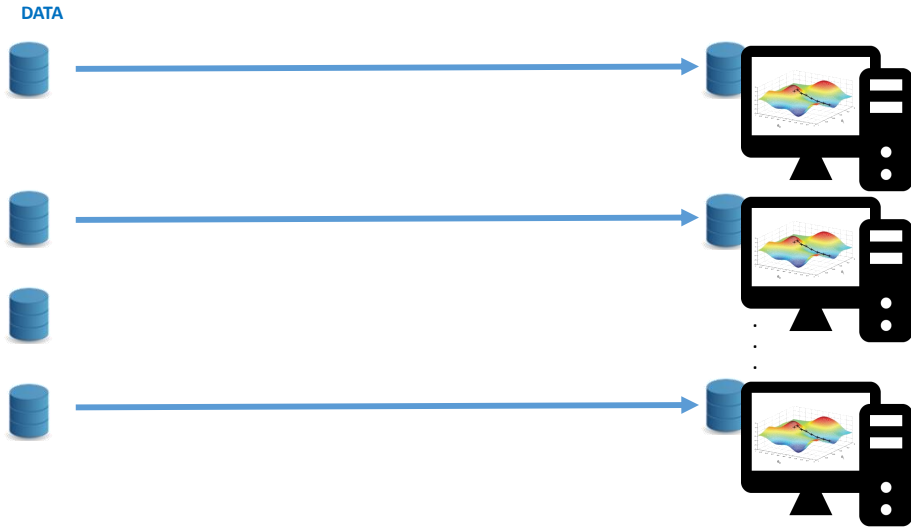
$$c^T x - \mu \sum_{i=1}^n \log(b_i - a_i^T x)$$



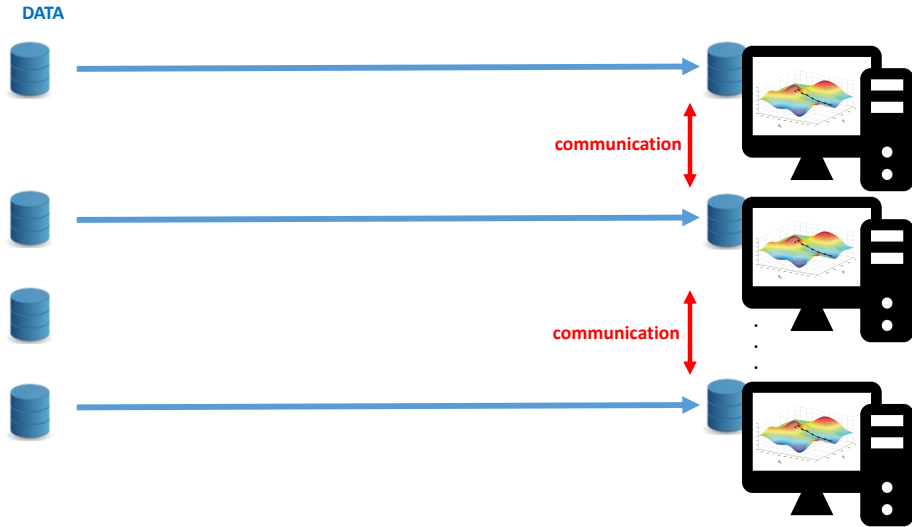
Distributed Optimization



Distributed Optimization



Distributed Optimization



Distributed optimization

- A : $n \times d$ feature matrix, and y : $n \times 1$ response vector
- Partition data $A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$, $y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$
- Least squares cost

$$\min_x \|Ax - y\|^2 = \min_x \|A_1x - y_1\|^2 + \|A_2x - y_2\|^2$$

Distributed optimization

- $A : n \times d$ feature matrix, and $y : n \times 1$ response vector
- Partition data $A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$, $y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$
- Least squares cost

$$\min_x \|Ax - y\|^2 = \min_x \|A_1x - y_1\|^2 + \|A_2x - y_2\|^2$$

Alternating Directions Method of Multipliers (ADMM)

(Hestenes, Powell 1969, Gabay et al. 1976, Boyd et al. 2011)

$$\min_x \|Ax - y\|^2 = \min_{x_1=x_2} \|A_1x_1 - y_1\|^2 + \|A_2x_2 - y_2\|^2$$

Distributed optimization

- $A : n \times d$ feature matrix, and $y : n \times 1$ response vector
- Partition data $A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$, $y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$
- Least squares cost

$$\min_x \|Ax - y\|^2 = \min_x \|A_1x - y_1\|^2 + \|A_2x - y_2\|^2$$

Alternating Directions Method of Multipliers (ADMM)

(Hestenes, Powell 1969, Gabay et al. 1976, Boyd et al. 2011)

$$\begin{aligned} \min_x \|Ax - y\|^2 &= \min_{x_1=x_2} \|A_1x_1 - y_1\|^2 + \|A_2x_2 - y_2\|^2 \\ &= \min_{x_1, x_2} \|A_1x_1 - y_1\|^2 + \|A_2x_2 - y_2\|^2 + \lambda^T(x_1 - x_2) \end{aligned}$$

Alternating Directions Method of Multipliers

→ x_1 update on machine 1

$$\min_{x_1} \|A_1x_1 - y_1\|^2 + \underbrace{\|A_2x_2 - y_2\|^2}_{\text{constant}} + \lambda^T(x_1 - x_2) + \rho\|x_1 - x_2\|^2$$

involves only A_1, y_1, x_2

Alternating Directions Method of Multipliers

→ x_1 update on machine 1

$$\min_{x_1} \|A_1x_1 - y_1\|^2 + \underbrace{\|A_2x_2 - y_2\|^2}_{\text{constant}} + \lambda^T(x_1 - x_2) + \rho\|x_1 - x_2\|^2$$

involves only A_1, y_1, x_2

→ x_2 update on machine 2

$$\min_{x_2} \underbrace{\|A_1x_1 - y_1\|^2}_{\text{constant}} + \|A_2x_2 - y_2\|^2 + \lambda^T(x_1 - x_2) + \rho\|x_1 - x_2\|^2$$

involves only A_2, y_2, x_1

Alternating Directions Method of Multipliers

→ x_1 update on machine 1

$$\min_{x_1} \|A_1x_1 - y_1\|^2 + \underbrace{\|A_2x_2 - y_2\|^2}_{\text{constant}} + \lambda^T(x_1 - x_2) + \rho\|x_1 - x_2\|^2$$

involves only A_1, y_1, x_2

→ x_2 update on machine 2

$$\min_{x_2} \underbrace{\|A_1x_1 - y_1\|^2}_{\text{constant}} + \|A_2x_2 - y_2\|^2 + \lambda^T(x_1 - x_2) + \rho\|x_1 - x_2\|^2$$

involves only A_2, y_2, x_1

→ communicate $x_1 \iff x_2$, update $\lambda \leftarrow \lambda + \rho(x_1 - x_2)$

Alternating Directions Method of Multipliers

- (Informal) Under some assumptions, ADMM converges in $O(\kappa_A \log(1/\epsilon))$ iterations, where κ_A is a conditioning parameter
- # iterations = rounds of communication

Distributed Sketching and ADMM

A_1 and A_2 are $n \times d$

$$A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$$

Define multiple sketching matrices

$$S_1 = \begin{bmatrix} I_{n \times n} & 0_{n \times n} \end{bmatrix}, \quad S_2 = \begin{bmatrix} 0_{n \times n} & I_{n \times n} \end{bmatrix}$$

Distributed Sketching and ADMM

A_1 and A_2 are $n \times d$

$$A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$$

Define multiple sketching matrices

$$S_1 = \begin{bmatrix} I_{n \times n} & 0_{n \times n} \end{bmatrix}, \quad S_2 = \begin{bmatrix} 0_{n \times n} & I_{n \times n} \end{bmatrix}$$

$$S_1 A = A_1$$

$$S_2 A = A_2$$

Distributed Sketching and ADMM

A_1 and A_2 are $n \times d$

$$A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$$

Define multiple sketching matrices

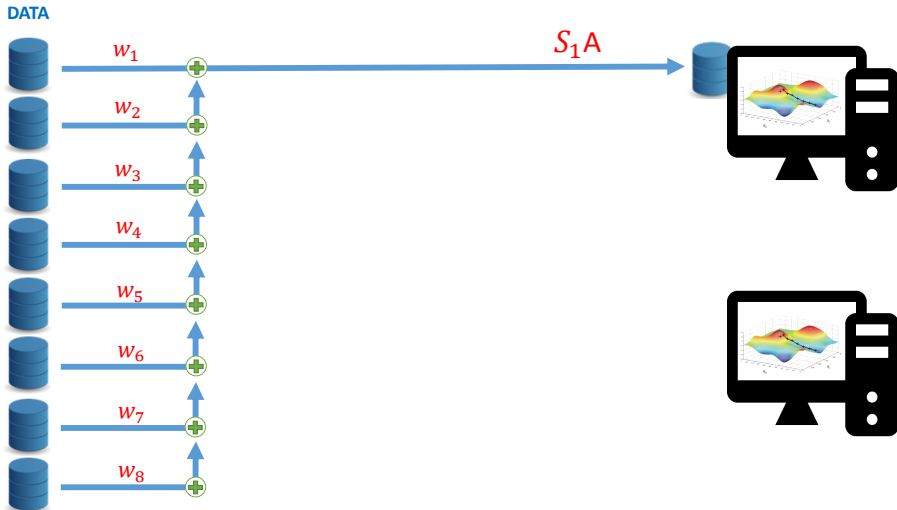
$$S_1 = \begin{bmatrix} I_{n \times n} & 0_{n \times n} \end{bmatrix}, \quad S_2 = \begin{bmatrix} 0_{n \times n} & I_{n \times n} \end{bmatrix}$$

$$S_1 A = A_1$$

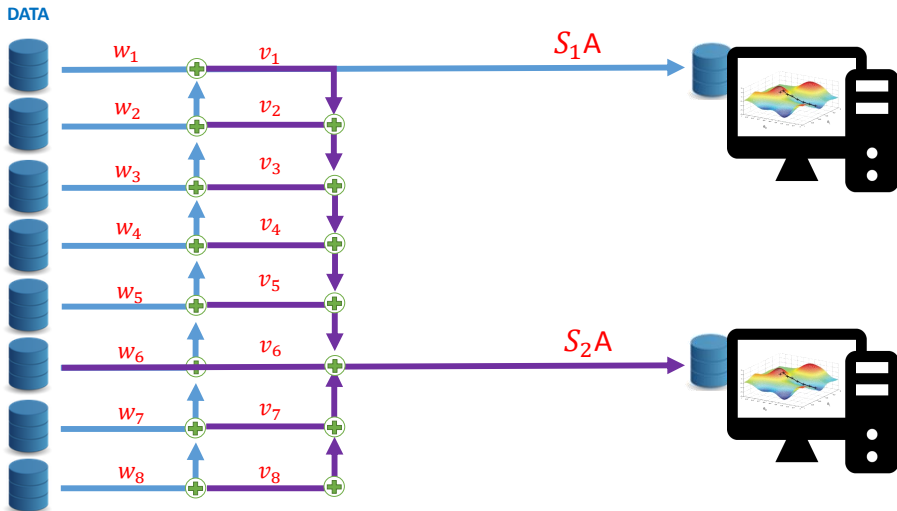
$$S_2 A = A_2$$

ADMM is operating on (naive) sketches!

Distributed Sketching



Distributed Sketching



Randomized Direction Method of Multipliers

Let S be a random orthogonal matrix, $S = \begin{bmatrix} S_1 \\ S_2 \end{bmatrix}$, $S^T S = I$
e.g., DFT matrix with randomly permuted rows

Randomized Direction Method of Multipliers

Let S be a random orthogonal matrix, $S = \begin{bmatrix} S_1 \\ S_2 \end{bmatrix}$, $S^T S = I$

e.g., DFT matrix with randomly permuted rows

→ x_1 update on machine 1

$$\min_{x_1} \|S_1(Ax_1 - y)\|^2 + \underbrace{\|S_2(Ax_2 - y)\|^2}_{\text{constant}} + \lambda^\top (x_1 - x_2) + \frac{\rho}{2} \|A(x_1 - x_2)\|^2$$

→ x_2 update on machine 2

$$\min_{x_2} \underbrace{\|S_1(Ax_1 - y)\|^2}_{\text{constant}} + \|S_2(Ax_2 - y)\|^2 + \lambda^\top (x_1 - x_2) + \frac{\rho}{2} \|A(x_1 - x_2)\|^2$$

→ communicate $x_1 \iff x_2$ and update λ

Randomized Direction Method of Multipliers

Let S be a random orthogonal matrix, $S = \begin{bmatrix} S_1 \\ S_2 \end{bmatrix}$, $S^T S = I$

e.g., DFT matrix with randomly permuted rows

→ x_1 update on machine 1

$$\min_{x_1} \|S_1(Ax_1 - y)\|^2 + \underbrace{\|S_2(Ax_2 - y)\|^2}_{\text{constant}} + \lambda^\top (x_1 - x_2) + \frac{1}{2} \|A(x_1 - x_2)\|^2$$

→ x_2 update on machine 2

$$\min_{x_2} \underbrace{\|S_2(Ax_2 - y)\|^2}_{\text{constant}} + \|S_2(Ax_2 - y)\|^2 + \lambda^\top (x_1 - x_2) + \frac{1}{2} \|A(x_1 - x_2)\|^2$$

→ communicate $x_1 \iff x_2$ and update λ

Theorem

Local solutions converge, and number of iterations for ϵ accuracy is less than

$$C \log(1/\epsilon),$$

*where C is a constant **independent** of data.*

[P. and Candès, 2018]

Theorem

Local solutions converge, and number of iterations for ϵ accuracy is less than

$$C \log(1/\epsilon),$$

*where C is a constant **independent** of data.*

[P. and Candès, 2018]

- $O(\log(1/\epsilon))$ rounds of communication. No condition number dependency.

Theorem

Local solutions converge, and number of iterations for ϵ accuracy is less than

$$C \log(1/\epsilon),$$

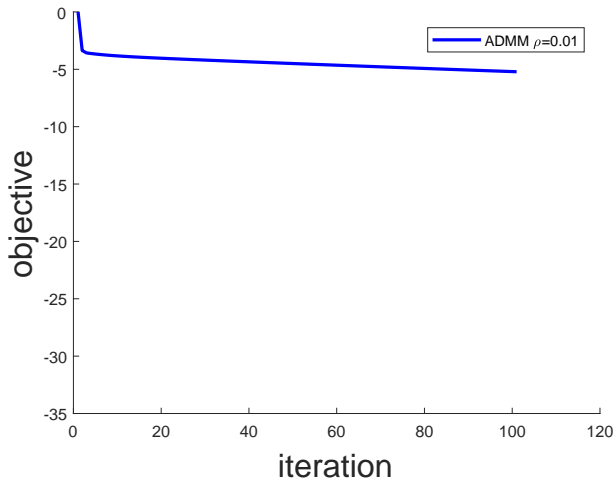
*where C is a constant **independent** of data.*

[P. and Candès, 2018]

- $O(\log(1/\epsilon))$ rounds of communication. No condition number dependency.
- Exchanging $O(d \log(1/\epsilon))$ bits to communicate x_1, x_2, \dots, x_M is information theoretically **optimal**

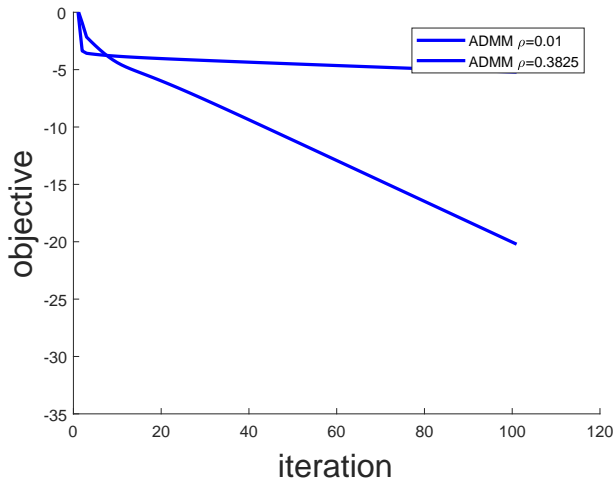
Randomized Direction Method of Multipliers

Random i.i.d. heavy tailed data



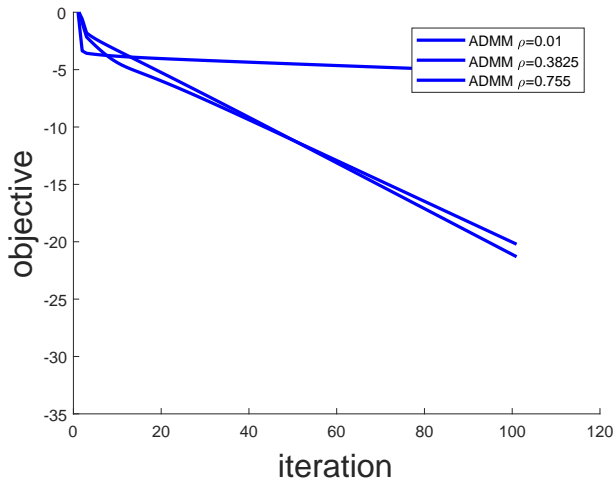
Randomized Direction Method of Multipliers

Random i.i.d. data



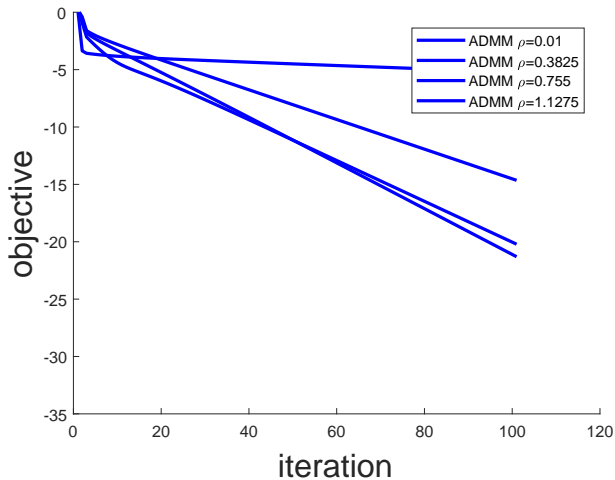
Randomized Direction Method of Multipliers

Random i.i.d. data



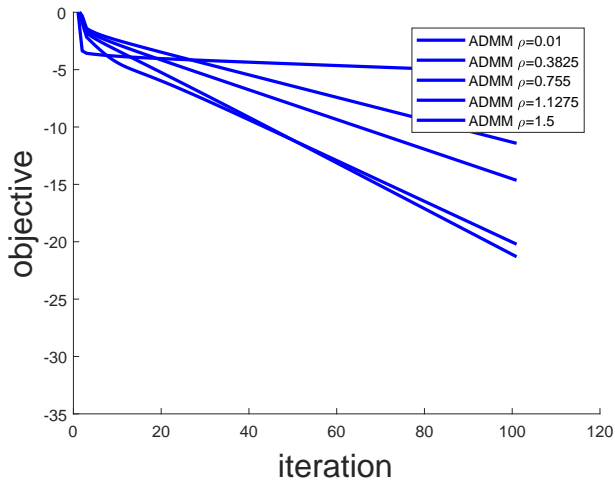
Randomized Direction Method of Multipliers

Random i.i.d. data



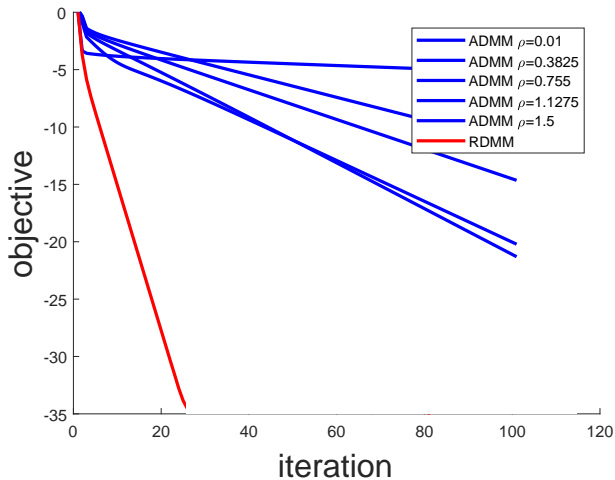
Randomized Direction Method of Multipliers

Random i.i.d. data



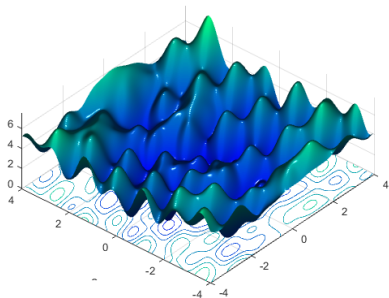
Randomized Direction Method of Multipliers

Random i.i.d. data



Non-convex Optimization Problems

- In general, very difficult to solve globally
- Need to make further assumptions



Non-convex Optimization Problems

$$\min_x \sum_{i=1}^n (f_x(a_i) - y_i)^2$$

Non-convex Optimization Problems

$$\min_x \sum_{i=1}^n (f_x(a_i) - y_i)^2$$

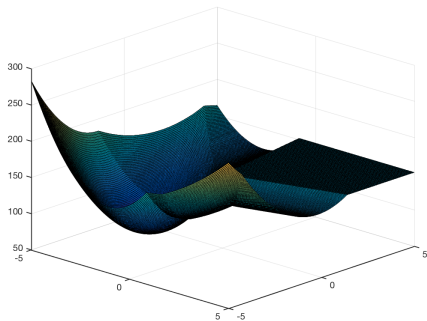
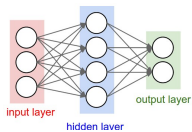
→ Heuristic: Gauss-Newton method

$$x_{t+1} = \arg \min_x \underbrace{\| f_{x_t}(A) + J_t x - y \|_2^2}_{\text{Taylor's approx for } f_x}$$

where $(J_t)_{ij} = \frac{\partial}{\partial x_j} f_x(a_i)$ is the Jacobian matrix

Non-convex Optimization Problems

Deep learning, nonlinear least squares...



2 layer ReLU neural network training loss

→ **Randomized** Gauss-Newton method

$$x_{t+1} = \arg \min_x \|S_t(f_{x_t}(A) + J_t x - y)\|_2^2.$$

- $S_t J_t$ backpropagation (Pearlmutter, 1994)

→ **Randomized** Gauss-Newton method

$$x_{t+1} = \arg \min_x \|S_t(f_{x_t}(A) + J_t x - y)\|_2^2.$$

- $S_t J_t$ backpropagation (Pearlmutter, 1994)

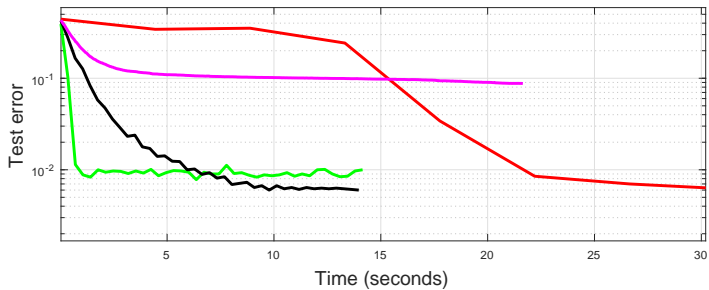
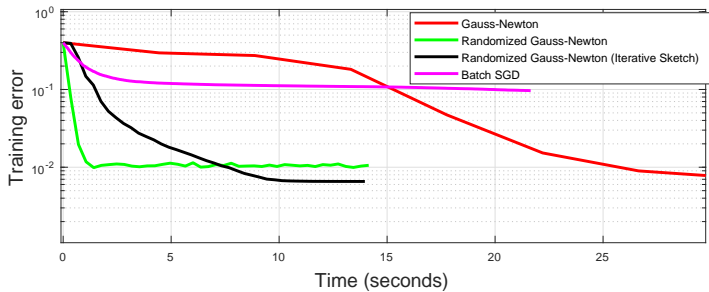
Theorem

(informal) Consider a **single hidden layer** neural network, **Gaussian** input data A . Randomized Gauss-Newton method converges to a **global minimum** in

$$C \log(1/\epsilon),$$

iterations.

[P. and Candès, 2018]



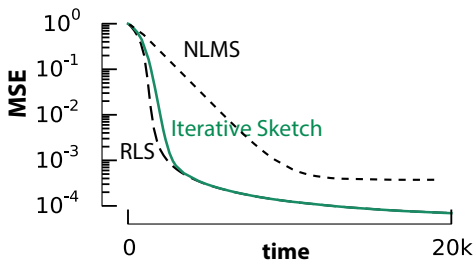
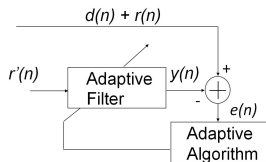
Summary

1. Information theoretic lower bounds for sketching
2. Iterative sketching with statistical optimality
3. Distributed sketching
4. Non-convex problems

Extensions: Streaming optimization

- Sketch can be updated by $[S \ s_+] \begin{bmatrix} A \\ a_+^T \end{bmatrix} = SA + s_+ a_+^T$

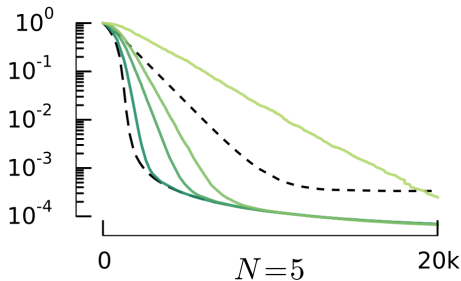
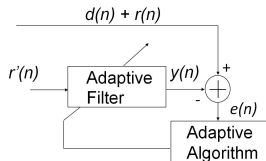
Example: Adaptive filtering



Extensions: Streaming optimization

- Sketch can be updated by $[S \ s_+] \begin{bmatrix} A \\ a_+^T \end{bmatrix} = SA + s_+ a_+^T$

Example: Adaptive filtering



Extensions

- Privacy preserving optimization



- Privacy preserving optimization



- SA provides privacy
- Mutual information constraint $I(SA; A) \leq \epsilon$

Yang, P., Wainwright, Annals of Statistics, 2015, P. (book chapter) 2018,

- Distributed and fault tolerant computing



- $S_1A, S_2A \dots, S_mA$ can be lost due to point failures

- Distributed and fault tolerant computing



- $S_1A, S_2A \dots, S_mA$ can be lost due to point failures
- Generate another sketch $S_{m+1}A$ i.i.d.

Thank you!
Questions ?