

CS-GY 9223 I: Lecture 7

Preconditioning, acceleration, coordinate decent, etc.

NYU Tandon School of Engineering, Prof. Christopher Musco

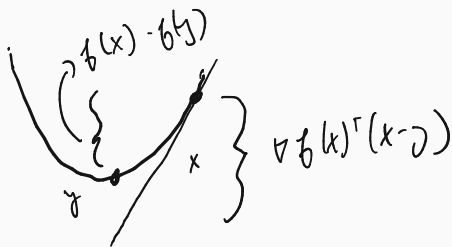
SMOOTH AND STRONGLY CONVEX

Recall from last lecture: a convex function f is β -smooth and α -strongly convex if, for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$,

$$\frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \leq \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{y}) - [f(\mathbf{x}) - f(\mathbf{y})] \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2.$$

$$\nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{y}) - [f(\mathbf{x}) - f(\mathbf{y})] \geq 0$$

when f is convex.



CONVERGENCE GUARANTEE

Theorem (GD for β -smooth, α -strongly convex.)

Let f be a β -smooth and α -strongly convex function. If we run GD for T steps (with step size $\eta = \frac{1}{\beta}$) we have:

$$\|x^{(t)} - x^*\|_2^2 \leq e^{-\frac{(t-1)\alpha}{\beta}} \|x^{(1)} - x^*\|_2^2$$

$\kappa = \frac{\beta}{\alpha}$ is called the "condition number" of f .

$$\log(1/\epsilon^2) = O(\log(1/\epsilon))$$

Corollary: If $T = O(\kappa \log(1/\epsilon))$ we have:

$$\|x^{(T)} - x^*\|_2 \leq \epsilon \|x^{(1)} - x^*\|_2.$$

ϵ^2

$$e^{-O(\kappa \log(1/\epsilon)) / \kappa} = e^{-O(\log(1/\epsilon))} = e^{O(\log(\epsilon))} \leq \epsilon.$$

$$e^{-t/\kappa}$$

$$\leq \epsilon$$

$$= e^{O(\log(\epsilon))} \leq \epsilon.$$

FROM LAST CLASS

Let $f(x) = \|Dx - b\|_2^2$ where D is a diagonal matrix.

- $\beta = 2 \max(D)^2$

- $\alpha = 2 \min(D)^2$

$$K = \frac{\max(D)^2}{\min(D)^2}$$

Gradient descent on f :

- $x^{(1)} = 0$

- For $t = 1, \dots, T$

- $\underline{x^{(t+1)}} = \underline{x^{(t)}} - \underline{\frac{1}{\beta}(2D(Dx^{(t)} - b))}$

$$\phi(x) = x^T D^2 x - 2x^T D b + b^T b$$

$$\nabla \phi(x) = 2D(Dx - b)$$

IN-CLASS EXERCISE

Theorem (GD for β -smooth, α -strongly convex.)

Let f be a β -smooth and α -strongly convex function. If we run GD for T steps (with step size $\eta = \frac{1}{\beta}$) we have:

$$\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 \leq e^{-t/\kappa} \|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2$$

Prove for $f(\mathbf{x}) = \|\mathbf{D}\mathbf{x} - \mathbf{b}\|_2^2$. You may assume that $\underline{\min}(\mathbf{D})^2 > 0$.

$$\mu = \frac{\max(\mathbf{D})^2}{\min(\mathbf{D})^2} \rightarrow \infty$$

IN-CLASS EXERCISE

Alternate view: $\nabla [\|Dx - b\|_2^2] = 2D(Dx - b)$
 $2D(Dx^* - b) = 0$

$$D^2 x^* = D^2 b$$

$$(x^{(t+1)} - x^*) = \left(1 - \frac{2}{\beta} D^2\right) (x^{(t)} - x^*)$$

Original gradient step:

$$x^{(t+1)} = x^{(t)} - \frac{1}{\beta} 2D(Dx^{(t)} - b)$$

$$x^{(t+1)} - x^* = x^{(t)} - x^* - \frac{2}{\beta} D^2 x^{(t)} + \frac{2}{\beta} \underbrace{Db}_{D^2 x^*}$$

$$x^{(t+1)} - x^* = \mathbb{I}(x^{(t)} - x^*) - \frac{2}{\beta} D^2 (x^{(t)} - x^*)$$

IN-CLASS EXERCISE

$$(x^{(t+1)} - x^*) = \left(I - \frac{2}{\beta} D^2 \right) (x^{(t)} - x^*)$$

$$(x^{(t+1)} - x^*) = \left(I - \frac{2}{\beta} D^2 \right)^t (x^{(1)} - x^*)$$

$$\beta = 2 \max(D^2)$$

$0 \leq \text{entries} \leq 1$

max entry ≤ 1

$$0 \leq \text{entries} \leq 1 - \frac{\min(D^2)}{\max(D^2)}$$

$$\|x^{(t+1)} - x^*\|_2 \leq$$

$$= \left\| \left(I - \frac{2}{\beta} D^2 \right)^t (x^{(1)} - x^*) \right\|_2$$

$$\boxed{1 - 1/k}$$

$$V = \begin{bmatrix} v_1 & & 0 \\ & \ddots & \\ 0 & & v_b \end{bmatrix}$$

$$V^t = \begin{bmatrix} v_1^t & & 0 \\ & \ddots & \\ 0 & & v_b^t \end{bmatrix}$$

small

$$= \|V^t (x^{(1)} - x^*)\|_2 \leq e^{-t/k} \|x^{(1)} - x^*\|_2$$

$$0 \leq \text{entries} \leq (1 - 1/k)^t = (1 - 1/k)^{k(t/k)} \leq \boxed{e^{-t/k}}$$

$$(1 - 1/k)^k \leq e^{-1}$$

GENERAL LINEAR REGRESSION

This same analysis holds for any linear system minimized via gradient descent:

$$\| \int_A^D x - \int v \|_2^2$$

$$\min_x \|Ax - b\|_2^2 = \min_x x^T A^T A x - x^T A^T b$$

$$\boxed{\nabla f(x) = 2A^T(Ax - b)} \quad \boxed{A^T A x^* = A^T b}$$

Unrolling gradient descent updates leads to:

$$(x^{(t+1)} - x^*) = \underbrace{(I - \eta A^T A)^t}_{\eta = 2/b} (x^{(1)} - x^*).$$

$$x^{(t+1)} = x^{(t)} - \eta A^T A x^{(t)} + \underbrace{\eta A^T b}_{= \eta A^T A x^*}$$

$$x^{(t+1)} - x^* = x^{(t)} - x^* - \eta A^T A (x^{(t)} - x^*)$$

GENERAL LINEAR REGRESSION

Quick linear algebra review:

$$A^T A = \boxed{\quad} = \boxed{\quad}^d \boxed{\Lambda} \boxed{U^T}$$

• $A^T A$ is symmetric so has an orthogonal eigendecomposition: $U \Lambda U^T$.

• $U^T U = U U^T = I$. U is orthogonal

• Λ is diagonal with entries $\lambda_1 \geq \lambda_2 \geq \dots, \lambda_d$.

Claim: $\lambda_d \geq 0$ (i.e., $A^T A$ is positive semidefinite).

Defn. $A^T A$ is positive semidefinite if
for all x , $x^T A^T A x \geq 0$.

$$y^T y = \|y\|_2^2 \geq 0$$

GENERAL LINEAR REGRESSION

Verify outside of class: $\rightarrow u = \frac{\lambda_1}{\lambda_d} = \frac{\max(\lambda)}{\min(\lambda)}$

$f(x) = \|Ax - b\|_2^2$ is $2\lambda_1$ smooth and $2\lambda_d$ strongly convex. So we

have: $\kappa = \frac{\lambda_1}{\lambda_d}$

$$I = UV^T \quad I = UIV^T$$

$$A = U\Lambda U^T$$

$$(x^{(t+1)} - x^*) = (I - \eta A^T A)^t (x^{(1)} - x^*)$$

$$(I - \eta A^T A)^t = (U \underbrace{(I - \eta \Lambda)}_r U^T)^t = \underline{U(I - \eta \Lambda)^t U^T}$$

$$I - \eta A^T A \quad (UVU^T)^t$$

$$= UVU^T UVU^T \dots UVU^T$$

$$UV^T U^T$$

$$= e^{-t/\kappa} \|x^{(1)} - x^*\|_2$$

has all entries $\in [0, 1 - 1/\kappa]$

$$\|x^{(t+1)} - x^*\|_2 = \|U(I - \eta \Lambda)^t U^T (x^{(1)} - x^*)\|_2$$

$$= \|(I - \eta \Lambda)^t U^T (x^{(1)} - x^*)\|_2 \leq e^{-t/\kappa} \|U^T (x^{(1)} - x^*)\|_2 \leq 10$$

We now have a really good understanding of gradient descent.

Number of iterations for ϵ error:

	G -Lipschitz	β -smooth
<u>R bounded start</u>	$O\left(\frac{G^2 R^2}{\epsilon^2}\right)$	$O\left(\frac{\beta R^2}{\epsilon}\right)$
α -strong convex	$O\left(\frac{G^2}{\alpha \epsilon}\right)$	$O\left(\frac{\beta}{\alpha} \log(1/\epsilon)\right)$

How do we use this understanding to design faster algorithms?

ACCELERATION

LINEAR REGRESSION RUNTIME

Total runtime for solving linear regression via GD: $d \leq n$

(time per iteration) \times (number of iterations)

$\kappa \approx \|A^T A x - b\|_2^2$
 $2A^T(Ax - b)$
 $O(\kappa \log(1/\epsilon))$
 $\kappa = \frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)}$
 $O(nd \cdot \kappa \log(1/\epsilon))$
 for $A \in \mathbb{R}^{n \times d}$, $x \in \mathbb{R}^d$, $b \in \mathbb{R}^n$.

$$x^* = (A^T A)^{-1} A^T b$$

$(d \times d) (n \times d)$
 $O(nd)$

$$O(d^3)$$

$$MV(A, v) \rightarrow Av$$

Others can be implemented in $\ll O(nd)$

$$\rightarrow O(nd^2)$$

Theorem (Accelerated Iterative Regression)

Let $\mathbf{x}^* = \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2^2$. There is an algorithm which finds $\tilde{\mathbf{x}}$ with $\|\tilde{\mathbf{x}} - \mathbf{x}^*\|_2 \leq \epsilon \|\mathbf{x}^*\|_2$ in time:

$$O(nd \cdot \underbrace{\sqrt{\kappa \log(1/\epsilon)}})$$

THE POLYNOMIAL VIEW

→ degree q

Claim: For any η , polynomial $\underline{p(z)} = c_1z + c_2z^2 + \dots + c_qz^q$ with $\underline{p(1)} = \sum_{j=1}^q c_j = 1$, there is an algorithm running in $\underline{O(ndq)}$ time which outputs $\underline{\tilde{x}}$ satisfying:

$$\|Ax - b\|_2^2$$

$$\tilde{x} - \tilde{x}^* = p(I - \eta A^T A)x^*$$

$$x^* - \tilde{x}^2 = p(I - \eta A^T A)x^*$$

For standard gradient descent, $\underline{p(z)} = z^q$.

$$x^{(1)} = 0$$

Gradient descent gives:

$$z^q$$

$$x^* - \tilde{x} = (I - \eta A^T A)^q (x^* - x^{(1)})$$

$$x^* - \tilde{x} = (I - \eta A^T A)^q x^*$$

$$O(ndq)$$

THE POLYNOMIAL VIEW

Claim: For any η , polynomial $p(z) = c_1z + c_2z^2 + \dots + c_qz^q$ with $p(1) = \sum_{j=1}^q c_j = 1$, there is an algorithm running in $O(ndq)$ time which outputs \tilde{x} satisfying:

$$\tilde{x} - x^* = p(I - \eta A^T A) x^* \\ \tilde{x} = \underbrace{c_1 \cdot (I - \eta A^T A) x^* + c_2 \cdot (I - \eta A^T A)^2 x^* + \dots + c_q \cdot (I - \eta A^T A)^q x^*}$$

Claim: $c_j \cdot (I - \eta A^T A)^j x^* = \underline{c_j \cdot x^*} + \underline{p'_j(I - \eta A^T A) A^T A} x^*$ where p'_j is a polynomial with degree $j - 1$.

$$\begin{aligned} c_j (I - \eta A^T A)^j x^* &= c_j (I - \eta_j A^T A + \dots) x^* \\ &= c_j x^* - (\quad) \\ &= \underline{c_j x^*} - (\quad) \underline{A^T A x^*} \end{aligned}$$

THE POLYNOMIAL VIEW

Claim: For any η , polynomial $p(z) = c_1z + c_2z^2 + \dots + c_qz^q$ with $p(1) = \sum_{j=1}^q c_j = 1$, there is an algorithm running in $O(ndq)$ time which outputs \tilde{x} satisfying:

$$x^* - \tilde{x} = \underbrace{(c_1 + c_2 + \dots + c_q)}_{\downarrow 1} \cdot x^* + \underbrace{p'(I - \eta A^T A) A^T A x^*}_{p' = p'_1 + p'_2 + \dots + p'_q}$$

$\tilde{x} = p'(I - \eta A^T A) A^T b$ where p' is a polynomial with degree $q - 1$.

$$\cancel{x^*} - \tilde{x} = \cancel{x^*} + p'(I - \eta A^T A) A^T b \quad \begin{matrix} (I - \eta A^T A) y \\ y = \eta A^T A x \\ \underbrace{\eta A^T A x}_{O(nd)} \end{matrix}$$

$$\tilde{x} = \underbrace{-p'(I - \eta A^T A) A^T b}_{\substack{\downarrow \\ \text{depends on } p}} \quad \text{can compute } \tilde{x} \text{ in } O(ndq) \text{ time.} \quad \begin{matrix} \downarrow \\ O(nd) \end{matrix}$$

p' has degree $q - 1$

THE POLYNOMIAL VIEW

$$p'(z) = c_1' z + c_2' z^2 + \dots + c_g' z^g$$

Time to compute

$$p'(V) A^T b = c_1' V A^T b + c_2' V^2 A^T b + \dots + V^g y$$

is αudg

$$\tilde{x}^* - x^* = p(I - \eta A^T A) x^*$$

$$p(I - \eta A^T A) = U p(I - \eta \Lambda) U^T$$

$$U = (I - \eta A^T A)$$

compute
 \tilde{x} in $O(\text{udg})$
 time

eigenvalues

$$(I - \eta \Lambda^1) (I - \eta \Lambda^2) \dots (I - \eta \Lambda^g)$$

$$\|\tilde{x} - x^*\|_2 = \|U p(I - \eta \Lambda) U^T x^*\|_2$$

$$= \|p(I - \eta \Lambda) U^T x^*\|_2$$

$$p(1) = 1$$

As long as $\max [p(I - \eta \Lambda)] \leq \epsilon$,

$$\leq \epsilon \|U^T x^*\|_2$$

$$M = \frac{1}{\max(\Lambda)}$$

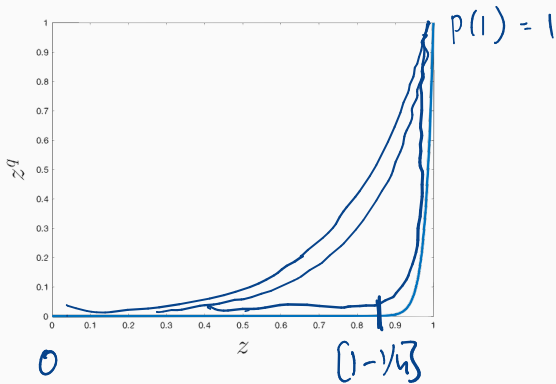
$$\|\tilde{x} - x^*\|_2 \leq \epsilon \|x^*\|_2$$

$$= \epsilon \|x^*\|_2$$

All entries in $I - \eta \Lambda$ are between $[0, 1 - 1/M]$

CONSTRUCTING A JUMP POLYNOMIAL

Goal: Find polynomial p such that $p(1) = 1$ and $p(z) \leq \epsilon$ for $z \in [0, 1 - \frac{1}{\kappa}]$.



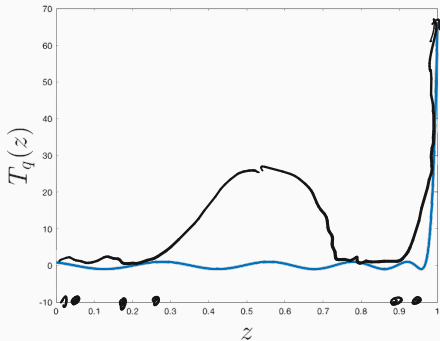
Gradient descent uses $p(z) = z^{O(\kappa \log(1/\epsilon))}$.

$p(z) \leq \epsilon$ for $z \in [0, 1 - \frac{1}{\kappa}]$

A BETTER JUMP POLYNOMIAL

Goal: Find polynomial p such that $p(1) = 1$ and $p(z) \leq \epsilon$ for $z \in [0, 1 - \frac{1}{\kappa}]$.

Gradient Descent for regression: Richardson Iteration



Accelerated:
"Nelder-Mead
Iteration"

"Conjugate
Gradient"

Can be done with degree $O(\sqrt{\kappa \log(1/\epsilon)})$ polynomial instead!

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

What are these polynomials?

Chebyshev polynomials of the first kind.

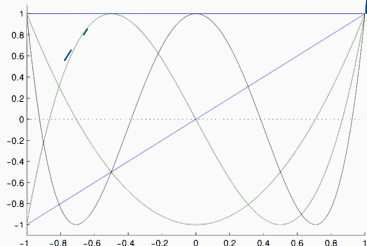
$$T_0(x) = 1$$

$$T_1(x) = x$$

$$T_2(x) = 2x^2 - 1$$

⋮

$$T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x)$$



“There’s only one bullet in the gun. It’s called the Chebyshev polynomial.” – Prof. Rocco Servedio

Nesterov's accelerated gradient descent:

- $\mathbf{x}^{(1)} = \mathbf{y}^{(1)} = \mathbf{z}^{(1)}$
- For $t = 1, \dots, T$
 - $\mathbf{y}^{(t+1)} = \mathbf{x}^{(t)} - \frac{1}{\beta} \nabla f(\mathbf{x}^{(t)})$
 - $\mathbf{x}^{(t+1)} = \left(1 + \frac{\sqrt{\kappa}-1}{\sqrt{\kappa+1}}\right) \mathbf{y}^{(t+1)} - \frac{\sqrt{\kappa}-1}{\sqrt{\kappa+1}} \mathbf{y}^{(t)}$

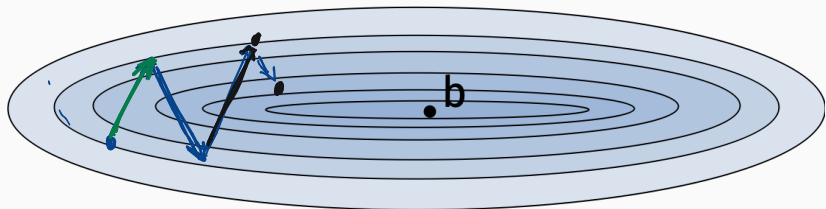
Theorem (AGD for β -smooth, α -strongly convex.)

Let f be a β -smooth and α -strongly convex function. If we run AGD for T steps we have:

$$f(\mathbf{x}^{(t)}) - f(\mathbf{x}^*) \leq \kappa e^{-(t-1)\sqrt{\kappa}} \left[f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*) \right]$$

Corollary: If $T = O(\underbrace{\sqrt{\kappa} \log(\kappa/\epsilon)}_{\text{achieve error } \epsilon})$

INTUITION BEHIND ACCELERATION



Level sets of $\|Ax - b\|_2^2$.

Other terms for similar ideas: GD greedy: $O(k)$

- Momentum

AGD (two iterates): $O(\sqrt{k})$

- Heavy-ball methods

look at d iterates: $O(k^{1/d})$

What if we look back beyond two iterates? $O(\sqrt{k})$

PRECONDITIONING

Preconditioning

Main idea: Instead of minimizing $f(\mathbf{x})$, find another function $\underline{g(\mathbf{x})}$ with the same minimum but which is better suited for first order optimization (e.g., has a smaller conditioner number).

Claim: Let $\underline{h(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^d}$ be an invertible function. Let $\underline{g(\mathbf{x}) = f(h(\mathbf{x}))}$. Then

$$\min_{\mathbf{x}} f(\mathbf{x}) = \min_{\mathbf{x}} g(\mathbf{x}) \quad \text{and} \quad \arg \min_{\mathbf{x}} f(\mathbf{x}) = h \left(\arg \min_{\mathbf{x}} g(\mathbf{x}) \right).$$

PRECONDITIONING

First Goal: We need $g(x)$ to still be convex.

$$f(h(x)) = g(x)$$

Claim: Let P be an invertible $d \times d$ matrix and let $g(x) = f(Px)$.

$$h(x) = Px$$

$g(x)$ is always convex.

$$\begin{matrix} \boxed{P} \\ \uparrow \downarrow \\ d \end{matrix} \quad \mathbb{R}^d$$

Need to prove: for all $x, y, \lambda \in [0, 1]$

$$\lambda g(x) + (1-\lambda)g(y) \geq g(\lambda x + (1-\lambda)y)$$

$$\lambda f(Px) + (1-\lambda)f(Py) \geq f(\lambda Px + (1-\lambda)Py)$$

$$\underbrace{\hspace{10em}} = f(P(\lambda x + (1-\lambda)y))$$

$$\lambda g(x) + (1-\lambda)g(y) \geq g(\lambda x + (1-\lambda)y)$$

PRECONDITIONING

Second Goal:

$$g(x) = Px$$

$g(x)$ should have better condition number κ than $f(x)$.

Example:

$$\cdot f(x) = \|Ax - b\|_2^2. \quad \kappa_f = \frac{\lambda_1(A^T A)}{\lambda_d(A^T A)}.$$

$$\cdot g(x) = \|APx - b\|_2^2. \quad \kappa_g = \frac{\lambda_1(P^T A^T A P)}{\lambda_d(P^T A^T A P)}. \quad \leftarrow \frac{\lambda_1(A^T A)}{\lambda_d(A^T A)}$$

Ideal preconditioner: Choose P so that $P^T A^T A P = I$. For example, could set $P = \sqrt{(A^T A)^{-1}}$. $\rightarrow U \sqrt{\kappa^{-1}} U^T$

$$A^T A = U \Lambda U^T$$

$$U \sqrt{\kappa^{-1}} U^T U \Lambda U^T U \sqrt{\kappa^{-1}} U^T = P^T A^T A P = I$$

What's the problem with this choice?

$$= \min_x \left(x^T P^T A^T A P x - x^T P^T A^T b \right)$$

$$\nabla g(x) = 2x - \underline{P^T A^T b}$$

DIAGONAL PRECONDITIONER

Third Goal: P should be easy to compute.

Many, many problem specific preconditioners are used in practice. Their design is usually a heuristic process.

$$(A^T A)_{ii} = a_i^T a_i$$

Example: Diagonal preconditioner.

- Let $D = \text{diag}(A^T A)$
- Intuitively, we roughly have that $D \approx A^T A$.
- Let $P = \sqrt{D^{-1}}$

does not take $O(n^2)$
takes $O(nd)$
 $\sqrt{(A^T A)^{-1}}$

P is often called a Jacobi preconditioner. Often works very well in practice!

DIAGONAL PRECONDITIONER

A =

-734	1	33	9111	0] β^x]
-31	-2	108	5946	-19	
232	-1	101	3502	10	
426	0	-65	12503	9	
-373	0	26	9298	0	
-236	-2	-94	2398	-1	
2024	0	-132	-6904	-25	
-2258	-1	92	-6516	6	
2229	0	0	11921	-22	
338	1	-5	-16118	-23	

```
>> cond(A'*A)
```

```
ans =
```

```
8.4145e+07
```

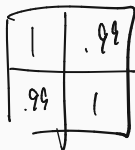
```
>> P = sqrt(inv(diag(diag(A'*A))));
```

```
>> cond(P*A'*A*P)
```

```
ans =
```

```
10.3878
```

Can you think of an example A where Jacobi preconditioning doesn't decrease a large κ ?



1	.99
.99	1

Can Jacobi preconditioning increase κ ?

ADAPTIVE STEPSIZES

Another view: If $g(x) = f(Px)$ then $\nabla g(x) = P^T \nabla f(Px)$.

$\nabla g(x) = P \nabla f(Px)$ when P is symmetric.

Gradient descent on g : $\nabla g(x^{(t)})$

- For $t = 1, \dots, T$,

$$\underline{P x^{(t+1)}} = \underline{P x^{(t)}} - \eta P^2 [\nabla f(Px^{(t)})]$$

$$x^{(t)} \approx x^* y$$

prev Gradient descent on f :

- For $t = 1, \dots, T$,

$$\underline{y^{(t+1)} = y^{(t)} - \eta P^2 [\nabla f(y^{(t)})]}$$

$$y^{(t)} = P x^{(t)} \approx x^* b$$

$$\eta P_{11}$$

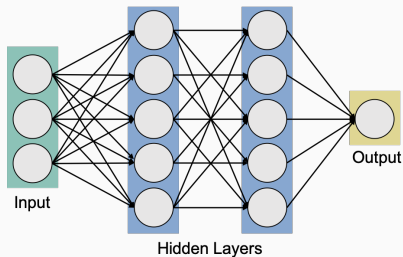
$$\eta P_{22}$$

$$\eta P_{dd}$$

When P is diagonal, this is just gradient descent with a different step size for each parameter!

Algorithms based on this idea:

- AdaGrad
- RMSprop
- Adam optimizer

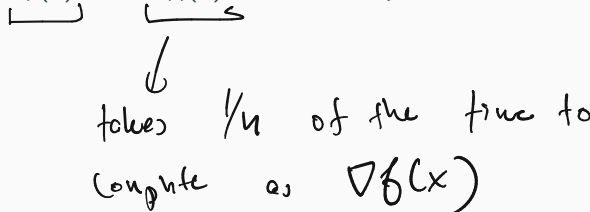


(Pretty much all of the most widely used optimization methods for training neural networks.)

COORDINATE DESCENT

Main idea: Trade slower convergence (more iterations) for cheaper iterations.

Stochastic Gradient Descent: When $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$, approximate $\nabla f(\mathbf{x})$ with $\nabla f_i(\mathbf{x})$ for randomly chosen i .


 takes $1/n$ of the time to
compute as $\nabla f(\mathbf{x})$

Main idea: Trade slower convergence (more iterations) for cheaper iterations.

Stochastic Coordinate Descent: Only compute a single random entry of $\nabla f(\mathbf{x})$ on each iteration:

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \frac{\partial f}{\partial x_2}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_d}(\mathbf{x}) \end{bmatrix}$$

$$\nabla_i f(\mathbf{x}) = \begin{bmatrix} 0 \\ \frac{\partial f}{\partial x_i}(\mathbf{x}) \\ \vdots \\ 0 \end{bmatrix}$$

Update: $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + \eta \nabla_i f(\mathbf{x}^{(t)})$

COORDINATE DESCENT

When x has d parameters, computing $\nabla f(x)$ often costs just a $1/d$ fraction of what it costs to compute $\nabla f(x)$

Example: $f(x) = \|Ax - b\|_2^2$ for $A \in \mathbb{R}^{n \times d}$, $x \in \mathbb{R}^d$, $b \in \mathbb{R}^n$.

• $\nabla f(x) = 2A^T Ax - 2A^T b$.

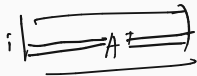
• $\nabla_i f(x) = 2 [A^T Ax]_i - 2 [A^T b]_i$

$O(d)$ times to compute $\nabla_i f(x)$.

$x^{(t)} = x^{(t-1)} + v_i$

where v_i is only non-zero in position i .

$O(d)$



$O(d)$

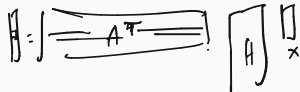
Have $Ax^{(t-1)} + Av_i$

$\neq a_i^T Ax$

can compute


$Ax^{(t)}$ in $O(d)$ time

a_i^T column $a_i^{(i)}$ row of A



Stochastic Coordinate Descent:

- Choose number of steps T and step size η .
- For $i = 1, \dots, T$:
 - Pick random $j_i \in 1, \dots, d$.
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \eta \nabla_{j_i} f(\mathbf{x}^{(i)})$
- Return $\hat{\mathbf{x}} = \frac{1}{T} \sum_{i=1}^T \mathbf{x}^{(i)}$.


$$\mathbb{E} \nabla_{j_i} f(x) = \nabla f(x)$$

Theorem (Stochastic Coordinate Descent convergence)

Given a G -Lipschitz function f with minimizer \mathbf{x}^* and initial point $\mathbf{x}^{(1)}$ with $\|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2 \leq R$, SCD with step size $\eta = \frac{1}{Rd}$ satisfies the guarantee:

$$\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \leq \frac{2GR}{\sqrt{T/d}}$$

\nearrow Lipschitz constant
 \rightarrow distance to optimal

γ_d

$$\frac{2GR}{\sqrt{T}}$$

$$T = \mathcal{O}\left(\frac{G^2 B^2}{\epsilon^2}\right)$$

Often it doesn't make sense to sample i uniformly at random:

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & -0.5 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & -2 & 0 & 0 & 0 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 10 \\ 42 \\ -11 \\ -51 \\ 34 \\ -22 \end{bmatrix}$$

Select indices i proportional to $\|\mathbf{a}_i\|_2^2$:

$$\Pr[\text{select index } i \text{ to update}] = \frac{\|\mathbf{a}_i\|_2^2}{\sum_{j=1}^d \|\mathbf{a}_j\|_2^2} = \frac{\|\mathbf{a}_i\|_2^2}{\|\mathbf{A}\|_2^2}$$