

CS-GY 9223 D: Lecture 8

Acceleration, preconditioning, coordinate methods

NYU Tandon School of Engineering, Prof. Christopher Musco

IMPROVING GRADIENT DESCENT

We now have a good understanding of gradient descent.

Number of iterations for ϵ error:

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

\downarrow
 $K = \frac{\beta}{\alpha}$

How do we use this understanding to design faster algorithms?

ACCELERATION

ACCELERATED GRADIENT DESCENT

Nesterov's accelerated gradient descent:

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

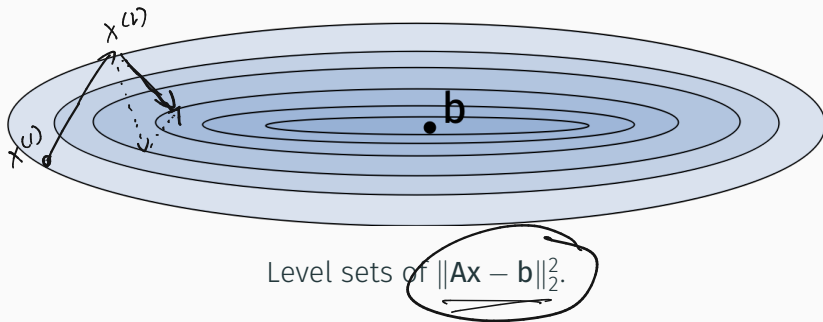
Theorem (AGD for β -smooth, α -strongly convex.)

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

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

Corollary: If $T = O(\sqrt{\kappa} \log(\kappa/\epsilon))$ achieve error ϵ .

INTUITION BEHIND ACCELERATION



Other terms for similar ideas:

- Momentum
- Heavy-ball methods

What if we look back beyond two iterates?

PRECONDITIONING

PRECONDITIONING

Main idea: Instead of minimizing $f(\mathbf{x})$, find another function $g(\mathbf{x})$ 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 $g(\mathbf{x}) = f(h(\mathbf{x}))$. Then

$$\min_{\mathbf{x}} f(\mathbf{x}) = \min_{\mathbf{x}} g(\mathbf{x}) \quad \text{and} \quad \arg \min_{\mathbf{x}} f(\mathbf{x}) = h \left(\arg \min_{\mathbf{x}} g(\mathbf{x}) \right).$$
$$\begin{aligned} \min_{\mathbf{x}} f(\mathbf{x}) &\leq \min_{\mathbf{x}} g(\mathbf{x}) & \mathbf{x}^* &= \arg \min_{\mathbf{x}} g(\mathbf{x}) & \min_{\mathbf{x}} f(\mathbf{x}) &\leq f(h(\mathbf{x}^*)) = \underline{g(\mathbf{x}^*)} \\ \min_{\mathbf{x}} g(\mathbf{x}) &\leq \min_{\mathbf{x}} f(\mathbf{x}) & \mathbf{x}^* &= \arg \min_{\mathbf{x}} f(\mathbf{x}) & \min_{\mathbf{x}} g(\mathbf{x}) &\leq \underline{g(h^{-1}(\mathbf{x}^*))} \\ & & & & &= f(h(h^{-1}(\mathbf{x}^*))) \\ & & & & &= f(\mathbf{x}^*) \end{aligned}$$

PRECONDITIONING

First Goal: We need $g(\mathbf{x})$ to still be convex.

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

$g(\mathbf{x})$ is always convex.

$$h(\mathbf{x}) = \mathbf{P}\mathbf{x}$$

\rightarrow h convex, so $f \circ h$ is convex

$$\lambda f(\mathbf{p}_x) + (1-\lambda)f(\mathbf{p}_y) \geq f(\lambda\mathbf{p}_x + (1-\lambda)\mathbf{p}_y)$$

$$\begin{aligned} \lambda g(\mathbf{x}) + (1-\lambda)g(\mathbf{y}) &\geq f(\mathbf{P}[\lambda\mathbf{x} + (1-\lambda)\mathbf{y}]) \\ &= g(\lambda\mathbf{x} + (1-\lambda)\mathbf{y}) \end{aligned}$$

If \mathbf{y}^* = $\arg \min g(\mathbf{y})$, then $\mathbf{x}^* = \mathbf{P}\mathbf{y}^*$ minimizes $f(\mathbf{x})$.

PRECONDITIONING

Second Goal: $\|A x - b\|^2$ $\|A P x - b\|^2$
 $A^T (A x - b)$ $P^T A^T (A P x - b)$
 $g(x)$ should have better condition number κ than $f(x)$. $P^T \nabla f(Px)$

High dimensional chain rule:

If $g(x) = f(Px)$, $\nabla^2 g(x) = \cancel{P^T \nabla^2 f(Px) P}$
 $P^T (\nabla^2 f(Px)) P$

Recall that the condition number is equal to:

$$\max_x \frac{\lambda_{\max}(\nabla^2 g(x))}{\lambda_{\min}(\nabla^2 g(x))} \quad \frac{d^2}{dx^2} f(Px) = P^2 f''(Px)$$

want to be small

$n \times n$

$$H = \nabla^2 f(x)$$

$$\lambda_{\min}(H) \geq \alpha$$

$$\lambda_{\max}(H) \leq \beta$$

$$\alpha I \preceq H \preceq \beta I$$

PRECONDITIONING

Example: $f(x) = g(x)$

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

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

Ideal preconditioner: Choose P so that $\underline{P^T A^T AP = I}$. For example, could set $P = \sqrt{(A^T A)^{-1}}$. But obviously this is too expensive to compute.

$$(A^T A)^{-1} \rightarrow O(nd^2 + d^3) \text{ time}$$

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.

Example: Diagonal preconditioner for least squares problems.

- Let $\mathbf{D} = \text{diag}(\mathbf{A}^T \mathbf{A})$



$O(nd)$

- Want $\mathbf{P}^T \mathbf{A}^T \mathbf{A} \mathbf{P}$ to be close to identity \mathbf{I} .

- Let $\mathbf{P} = \sqrt{\mathbf{D}^{-1}}$

$$\mathbf{P}^T \mathbf{A}^T \mathbf{A} \mathbf{P} = \mathbf{I} \rightarrow k=1 \text{ when } \mathbf{A}^T \mathbf{A} \text{ is diagonal}$$

\mathbf{P} is often called a **Jacobi preconditioner**. Often works very well in practice!

$$\mathbf{D} \approx \mathbf{A}^T \mathbf{A}$$

$\mathbf{D} = \text{diag min}$
diagonal matrix \mathbf{S}

$$\|\mathbf{A}^T \mathbf{A} - \mathbf{S}\|_F$$

$$\sqrt{\mathbf{D}^{-1}} \approx \sqrt{(\mathbf{A}^T \mathbf{A})^{-1}}$$

DIAGONAL PRECONDITIONER

A =

-734	1	33	9111	0
-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))));
```

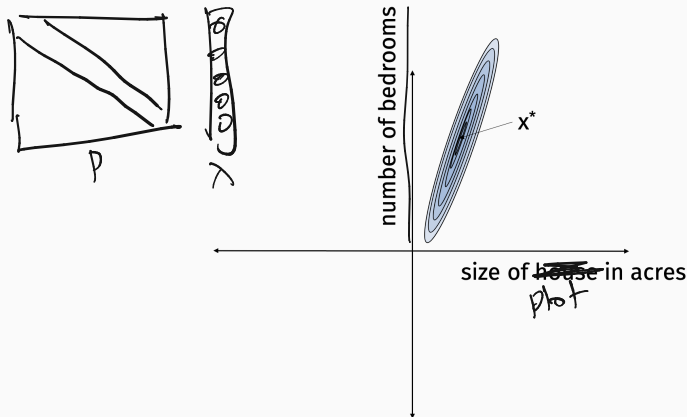
```
>> cond(P*A'*A*P)
```

ans =

10.3878

DIAGONAL PRECONDITIONER INTUITION

$g(\mathbf{x}) = f(\|\mathbf{A}P\mathbf{x} - \mathbf{b}\|_2^2)$ is the same least squares problem as $f(\mathbf{x}) \equiv \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$, but with each feature (column of \mathbf{A}) scaled differently. The i^{th} column is scaled by P_{ii} .

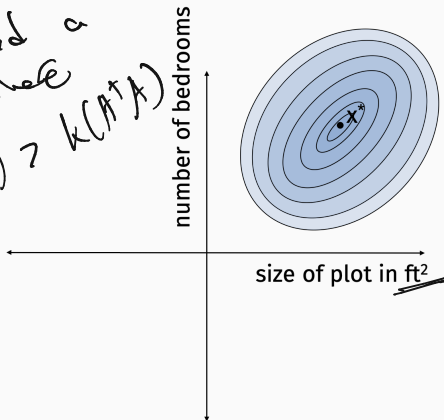


Feature scaling can have a huge impact on conditioning.

DIAGONAL PRECONDITIONER INTUITION

$g(\mathbf{x}) = f(\|\mathbf{A}\mathbf{P}\mathbf{x} - \mathbf{b}\|_2^2)$ is the same least squares problem as $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$, but with each feature (column of \mathbf{A}) scaled differently. The i^{th} column is scaled by P_{ii} .

Exercise: Find a matrix \mathbf{P} where $\kappa(\mathbf{P}^T \mathbf{A}^T \mathbf{A} \mathbf{P}) \approx \kappa(\mathbf{A}^T \mathbf{A})$



Feature scaling can have a huge impact on conditioning.

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.

$$x^* = P y^*$$
$$y^* = \arg \min_y g(y)$$

Gradient descent on g :

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

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

Gradient descent on g :

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

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

$$y^{(i)} = P^{(i)}$$

→ vector

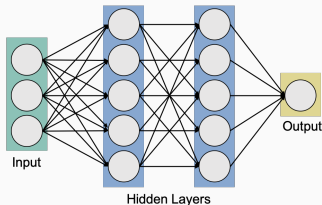
When P is diagonal, this is just gradient descent with a different step size for each parameter!

ADAPTIVE STEPSIZES

Less clear how to set P for general optimization problems where the Hessian is changing, but lots of heuristic algorithms based on this idea:

- AdaGrad, AdaDelta
- 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 .

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:

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

1/d cheaper (maybe)

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

Update: $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \eta \nabla_{if}(\mathbf{x}^{(t)})$.

COORDINATE DESCENT

$$A^T A \sim \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \end{bmatrix} \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \end{bmatrix} \rightarrow O(nd)$$

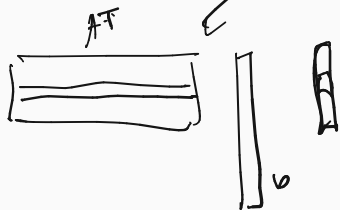
When \mathbf{x} has d parameters, computing $\nabla_i f(\mathbf{x})$ sometimes costs just a $1/d$ fraction of what it costs to compute $\nabla f(\mathbf{x})$

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

$$\cdot \nabla f(\mathbf{x}) = 2\mathbf{A}^T(\mathbf{Ax}) - 2\mathbf{A}^T\mathbf{b} \quad O(nd)$$

$$\cdot \nabla_i f(\mathbf{x}) = 2 \underbrace{[\mathbf{A}^T \mathbf{Ax}]_i}_{O(n)} - 2 \underbrace{[\mathbf{A}^T \mathbf{b}]_i}_{O(d)} \quad nd + nd + nd$$

Computing full gradient takes $O(nd)$ time. Can we do better here?



$$O(n)$$

$$\mathbf{A}^T (\mathbf{Ax}) \quad \underbrace{\hspace{2cm}}_{O(nd)}$$

$$nd + n + n$$

COORDINATE DESCENT

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

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

- $\nabla f(\mathbf{x}) = 2\mathbf{A}^T \mathbf{Ax} - 2\mathbf{A}^T \mathbf{b}$.
- $\nabla_i f(\mathbf{x}) = 2 [\mathbf{A}^T \mathbf{Ax}]_i - 2 [\mathbf{A}^T \mathbf{b}]_i$.

$$\begin{aligned} & \mathbf{Ax}^{(t+1)} = \mathbf{Ax}^{(t)} + c \mathbf{e}_i \\ & 2 [\mathbf{A}^T (\mathbf{Ax}^{(t+1)} - \mathbf{b})]_i \end{aligned}$$

$$\mathbf{Ax}^{(t)} = \begin{bmatrix} b_1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\mathbf{A} \mathbf{e}_i = \begin{bmatrix} a_{1i} \\ a_{2i} \\ a_{3i} \\ a_{4i} \end{bmatrix}$$

$$\mathbf{Ax}^{(t+1)} = \mathbf{Ax}^{(t)} + c \mathbf{A} \mathbf{e}_i$$

$O(n)$ time

$O(n)$ time

$$c a_i$$

Stochastic Coordinate Descent:

- Choose number of steps T and step size η .
- For $t = 1, \dots, T$:
 - ~~set~~ $\mathbf{x}^{(0)} = \mathbf{0}$
 - Pick random $\underline{j} \in 1, \dots, d$ uniformly at random.
 - $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \underline{\nabla_{j,j} f(\mathbf{x}^{(t)})}$
- Return $\hat{\mathbf{x}} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}^{(t)}$.

STOCHASTIC COORDINATE DESCENT

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:

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

$$T = O\left(\frac{G^2 R^2}{\epsilon^2}\right) \cdot d \quad \text{gives error } \epsilon$$

$$O\left(\frac{G^2 \cdot R^2}{\epsilon^2}\right)$$

IMPORTANCE SAMPLING

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

$$A = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & -5.5 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & -2 & 0 & 0 & 0 \end{bmatrix} \quad b = \begin{bmatrix} 10 \\ 42 \\ -11 \\ -51 \\ 34 \\ -22 \end{bmatrix}$$

Select indices i proportional to $\|a_i\|_2^2$:

$$\Pr[\text{select index } i \text{ to update}] = \frac{\|a_i\|_2^2}{\sum_{j=1}^d \|a_j\|_2^2} = \frac{\|a_i\|_2^2}{\|A\|_F^2}$$

Let's analyze this approach.

$$\|A\|_F^2 = \sum_{i,j} A_{ij}^2$$

STOCHASTIC COORDINATE DESCENT

Specialization of SCD to $\|Ax - b\|_2^2$:

Randomized Coordinate Descent (Strohmer, Vershynin 2007 /
Leventhal, Lewis 2018)

- For iterate $x^{(t)}$, let $r^{(t)}$ be the residual:

$$\underline{r^{(t)}} = \underline{Ax^{(t)} - b}$$

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

- $r^{(t+1)} = r^{(t)} - ca_j$. Here a_j is the j^{th} column of A .

$$Ax^{(t+1)} - b = A(x^{(t)} - ce_j) - b = Ax^{(t)} - b - cAe_j$$

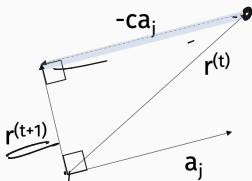
Handwritten notes: $r^{(t)}$ above $r^{(t+1)}$, and $\rightarrow = ca_j$ pointing to $-ca_j$.

Typically c depends on fixed learning rate. Here we will choose it optimally – similar idea to gradient descent with line search.

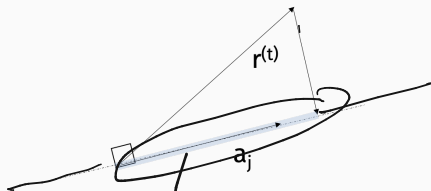
STOCHASTIC COORDINATE DESCENT

What choice for c minimizes $\|r^{(t+1)}\|_2^2$?

- $\|r^{(t+1)}\|_2^2 = \|r^{(t)} - ca_j\|_2^2$ \rightarrow small α \rightarrow $\frac{1}{2} \rightarrow 5, 1, 10$
- Requires projecting $r^{(t)}$ onto perpendicular of a_j .



$$c = \frac{a_j^T r^{(t)}}{\|a_j\|_2^2}$$



$$\frac{a_j^T r^{(t)}}{\|a_j\|_2^2} \cdot a_j$$

Note that $\|r^{(t+1)}\|_2^2 = \|r^{(t)}\|_2^2 - \|ca_j\|_2^2 = \|r^{(t)}\|_2^2 - \frac{(a_j^T r^{(t)})^2}{\|a_j\|_2^2}$

$$\frac{(a_j^T r^{(t)})^2}{\|a_j\|_2^2} \cdot \frac{1}{\|a_j\|_2^2}$$

Specialization of SCD to $\|\mathbf{Ax} - \mathbf{b}\|_2^2$:

Randomized Coordinate Descent

- Choose number of steps T .
- Let $\mathbf{x}^{(1)} = \mathbf{0}$ and $\mathbf{r}^{(1)} = \mathbf{b}$.
- For $t = 1, \dots, T$:
 - Pick random $j \in 1, \dots, d$. Index j is selected with probability proportional to $\|\mathbf{a}_j\|_2^2 / \|\mathbf{A}\|_F^2$.
 - Set $c = \mathbf{a}_j^T \mathbf{r}^{(t)} / \|\mathbf{a}_j\|_2^2$
 - $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - c \mathbf{e}_j$
 - $\mathbf{r}^{(t+1)} = \mathbf{r}^{(t)} - c \mathbf{a}_j$
- Return $\mathbf{x}^{(T)}$.

CONVERGENCE

Claim

$$\left(I - \frac{\lambda_{\max}(A^T A)}{\|A\|_F^2} \cdot I \right)$$

$$\rightarrow \lambda_{\max}(A^T A) \cdot \|r^{(t)}\|_2^2 - \frac{\lambda_{\max}(A^T A)}{\|A\|_F^2} \cdot \|r^{(t)}\|_2^2$$

$$\mathbb{E} \|r^{(t+1)}\|_2^2 = \|r^{(t)}\|_2^2 - \frac{1}{\|A\|_F^2} \|A^T r^{(t)}\|_2^2$$

$$\mathbb{E} \|r^{(t+1)}\|_2^2 = \mathbb{E} \left[\|r^{(t)}\|_2^2 - \frac{(a_j^T r^{(t)})^2}{\|a_j\|_2^2} \right]$$

$$= \sum_{j=1}^d \frac{\|a_j\|_2^2}{\|A\|_F^2} \cdot \left(\|r^{(t)}\|_2^2 - \frac{(a_j^T r^{(t)})^2}{\|a_j\|_2^2} \right)$$

$$= \|r^{(t)}\|_2^2 - \sum_{j=1}^d \frac{(a_j^T r^{(t)})^2}{\|A\|_F^2} = \|A^T r^{(t)}\|_2^2$$

$$= \|r^{(t)}\|_2^2 - \frac{1}{\|A\|_F^2} \sum_{j=1}^d (a_j^T r^{(t)})^2$$

CONVERGENCE

$$r^* = Ax^* - b \quad r^*$$

Any residual \mathbf{r} can be written as $\mathbf{r} = \underline{\mathbf{r}^*} + \bar{\mathbf{r}}$ where $\mathbf{r}^* = \mathbf{A}\mathbf{x}^* - \mathbf{b}$ and $\bar{\mathbf{r}} = \mathbf{A}(\mathbf{x}^t - \mathbf{x}^*)$. Note that $\mathbf{A}^T \mathbf{r}^* = 0$ and $\bar{\mathbf{r}} \perp \mathbf{r}^*$.

Claim $\bar{\mathbf{r}}^{(t)} = \mathbf{r}^{(t)} - \mathbf{r}^*$

$$\underline{\mathbb{E} \|\bar{\mathbf{r}}^{(t+1)}\|_2^2} \leq \underline{\|\bar{\mathbf{r}}^{(t)}\|_2^2} - \frac{\lambda_{\min}(\mathbf{A}^T \mathbf{A})}{\|\mathbf{A}\|_F^2}$$

$$\mathbb{E} \|\bar{\mathbf{r}}^{(t+1)}\|_2^2 + \|\mathbf{r}^*\|_2^2 \leq \|\bar{\mathbf{r}}^{(t)}\|_2^2 + \|\mathbf{r}^*\|_2^2 - \frac{1}{\|\mathbf{A}\|_F^2} \|\mathbf{A}^T \bar{\mathbf{r}}^{(t)}\|_2^2 \|\bar{\mathbf{r}}^{(t)}\|_2^2$$

Exercise: Because $\bar{\mathbf{r}}$ is in the column span of \mathbf{A} ,

$$\|\mathbf{A}^T \bar{\mathbf{r}}^{(t)}\|_2^2 \geq \lambda_{\min}(\mathbf{A}^T \mathbf{A}) \|\bar{\mathbf{r}}^{(t)}\|_2^2$$

CONVERGENCE

Theorem (Randomized Coordinate Descent convergence)

After T steps of RCD with importance sampling run on

$f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2$, we have:

$$\approx 1/e \quad \log(1/\epsilon)$$

$$\mathbb{E}[f(\mathbf{x}^{(t)}) - f(\mathbf{x}^*)] \leq \left(1 - \frac{\lambda_{\min}(\mathbf{A}^T \mathbf{A})}{\|\mathbf{A}\|_F^2}\right)^t [f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*)]$$

Corollary: After $T = O\left(\underbrace{\frac{\|\mathbf{A}\|_F^2}{\lambda_{\min}(\mathbf{A}^T \mathbf{A})} \log \frac{1}{\epsilon}}_{f \approx \frac{\|\mathbf{A}\|_F^2}{\lambda_{\min}(\mathbf{A}^T \mathbf{A})}}\right)$ we obtain error $\epsilon \|\mathbf{b}\|_2^2$.

Is this more or less iterations than the $T = O\left(\frac{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}{\lambda_{\min}(\mathbf{A}^T \mathbf{A})} \log \frac{1}{\epsilon}\right)$ required for gradient descent to converge?

COMPARISON

Recall useful linear algebraic fact:

$$\|A\|_F^2 = \text{tr}(A^T A) = \sum_{i=1}^d \lambda_i(A^T A) \leq d \cdot \lambda_{\max}(A^T A)$$

Handwritten notes: $\geq \lambda_{\max}(A^T A)$ (with an arrow pointing to the sum) and $\leq d \cdot \lambda_{\max}(A^T A)$ (with an arrow pointing to the final term).

$$\lambda_{\max}(A^T A) \leq \|A\|_F^2 \leq d \cdot \lambda_{\max}(A^T A)$$

For solving $\|Ax - b\|_2^2$,

$$(\# \text{ GD Iterations}) \leq (\# \text{ RCD Iterations}) \leq d \cdot (\# \text{ GD Iterations})$$

But RCD iterations are cheaper by a factor of d .

COMPARISON

When does $\|A\|_F^2 = \text{tr}(A^T A) = d \cdot \lambda_{\max}(A^T A)$?

$$A = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}$$

$$\|A\|_F^2 = d$$

$$A^T A = I$$

↓
eigenvalues all 1
 $\lambda_{\max} = 1$

$$\downarrow \cdot \lambda_{\max}(A^T A) = d$$

When does $\|A\|_F^2 = \text{tr}(A^T A) = \underline{1 \cdot \lambda_{\max}(A^T A)}$?

$$\begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & \ddots & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix}$$

$$\|A\|_F^2 \leq \underline{\text{rank}(A)} \cdot \lambda_{\max}(A^T A)$$



Roughly:

Stochastic Gradient Descent performs well when data points (rows) are repetitive.

Stochastic Coordinate Descent performs well when data features (columns) are repetitive.

NON-CONVEX OPTIMIZATION

STATIONARY POINTS

We understand much less about optimizing non-convex functions in comparison to convex functions, but not nothing. In many cases, we're still figuring out the right questions to ask

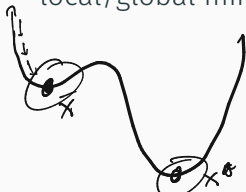
Definition (Stationary point)

For a differentiable function f , a stationary point is any x with:

$$\nabla f(x) = 0$$

$$\| \nabla f(x) \|_2 = 0$$

local/global minima - local/global maxima - saddle points



Reasonable goal: Find an approximate stationary point $\hat{\mathbf{x}}$ with

$$\|\nabla f(\hat{\mathbf{x}})\|_2^2 \leq \epsilon.$$

(Handwritten note: $\epsilon \rightarrow 0$)

Definition

A differentiable (potentially non-convex) function f is β smooth if for all \mathbf{x}, \mathbf{y} ,

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq \beta \|\mathbf{x} - \mathbf{y}\|_2$$

Corollary: For all \mathbf{x}, \mathbf{y}

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

Theorem

If GD is run with step size $\eta = \frac{1}{\beta}$ on a differentiable function f with global minimum \mathbf{x}^* then after $T = O(\frac{\beta[f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*)]}{\epsilon})$ we will find an ϵ -approximate stationary point $\hat{\mathbf{x}}$.

- $\nabla f(\mathbf{x}^{(t)})^T (\mathbf{x}^{(t)} - \mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)}) + f(\mathbf{x}^{(t+1)}) \leq \frac{\beta}{2} \|\mathbf{x}^{(t)} - \mathbf{x}^{(t+1)}\|_2^2.$
- $f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)}) \leq \frac{\beta}{2} \eta^2 \|\nabla f(\mathbf{x}^{(t)})\|_2^2 - \eta \|\nabla f(\mathbf{x}^{(t)})\|_2^2$
- $f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)}) \leq \frac{-\eta}{2} \|\nabla f(\mathbf{x}^{(t)})\|_2^2$
- $\frac{1}{T} \sum_{t=1}^T \frac{\eta}{2} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{1}{T} \sum_{t=1}^T f(\mathbf{x}^{(t)}) - f(\mathbf{x}^{(t+1)})$
- $\frac{\eta}{2} \min_t \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{1}{T} [f(\mathbf{x}^{(1)}) - f(\mathbf{x}^{(T)})]$

If GD can find a stationary point, are there algorithms which find a stationary point faster using preconditioning, acceleration, stochastic methods, etc.?

QUESTIONS IN NON-CONVEX OPTIMIZATION

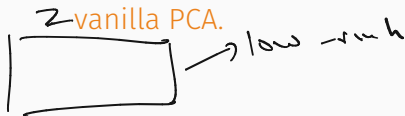
What if my function only has global minima and ~~stationary~~ ^{saddle} points? Randomized methods (SGD, perturbed gradient methods, etc.) can “escape” stationary points under some minor assumptions.

Example: $\min_x \frac{-x^T A^T A x}{x^T x}$



- **Global minimum:** Top eigenvector of $A^T A$ (i.e., top principal component of A).
- **Stationary points:** All other eigenvectors of A .

Useful for lots of other matrix factorization problems beyond



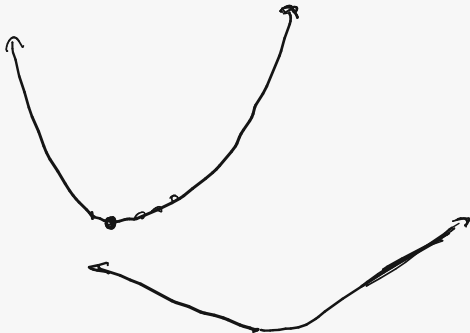
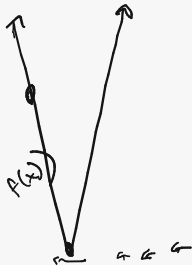
$$\max_x \| \nabla f(x) \|_2 = G$$

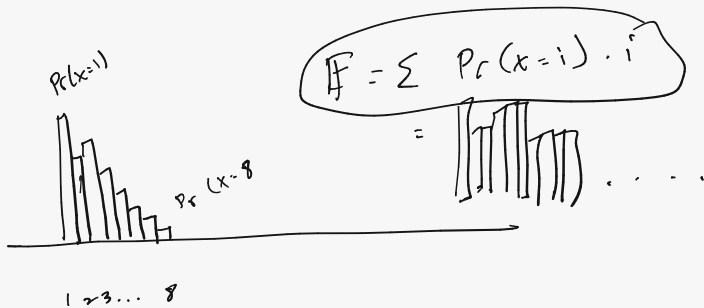
$$\frac{1}{G}$$

$$\min \nabla f(x)$$

$$G = \max$$

$$\max (\nabla^2 P(x))$$





$$\sum_{i=1}^n Pr\{X \geq i\}$$