

# CS-GY 9223 I: Lecture 8

## Coordinate descent and non-convex models.

---

NYU Tandon School of Engineering, Prof. Christopher Musco

**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:**  $\underline{\mathbf{x}}^{(t+1)} \leftarrow \underline{\mathbf{x}}^{(t)} + \eta \nabla_{i} f(\mathbf{x}^{(t)})$ .

## COORDINATE DESCENT

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

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

- $\nabla f(x) = 2A^T Ax - 2A^T b$ .
- $\nabla_j f(x) = 2 [A^T Ax]_j - 2 [A^T b]_j$ .

Full gradient for:

$$f(x) = \|Ax - b\|_2^2$$

$\rightarrow O(nd)$   
time

- $Ax^{(t+1)} = A(x^{(t)} + c \cdot e_j)$
- $2 [A^T (Ax^{(t+1)} - b)]_j$

$O(n)$  time

$O(n)$  time

## Stochastic Coordinate Descent:

- Choose number of steps  $T$  and step size  $\eta$ .
- For  $t = 1, \dots, T$ :
  - Pick random  $j \in 1, \dots, d$ .
  - $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \nabla_j f(\mathbf{x}^{(t)})$
- 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:

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

$$\leq O\left(\frac{GR}{\sqrt{T}}\right) \text{ for full gradient descent}$$

How can we improve on this?

SCD takes  $d \times$  more iterations for same error.

## IMPORTANCE SAMPLING

Often it doesn't make sense to sample  $i$  uniformly at random:  $\|Ax - b\|_2^2$   $Ax \approx b$

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

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.  $\rightarrow$  with column of  $A$   
 $= \|A\|_F^2$

## STOCHASTIC COORDINATE DESCENT

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

Random Kaczmarz  
method

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

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

$$f(x^{(t)}) = \|r^{(t)}\|_2^2$$

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

- $x^{(t+1)} = x^{(t)} - c e_j$ . Here  $c$  is a scalar and  $e_j$  is a standard basis vector.
- $\underline{r^{(t+1)}} = r^{(t)} - ca_j$ . Here  $a_j$  is the  $i^{\text{th}}$  column of  $A$ .

$$\begin{aligned} Ax^{(t+1)} - b &= A(x^{(t)} - ce_j) - b \\ &= r^{(t)} - ca_j \end{aligned}$$

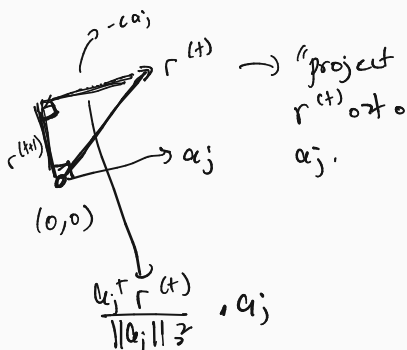
# STOCHASTIC COORDINATE DESCENT

Typically:  $\eta = \text{fixed}$ . "line search"  $\rightarrow f(x^{(t+1)})$

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

$$x^{(t+1)} = x^{(t)} - c e_j \quad \text{so that I minimize } f(x^{(t+1)})$$

$$\|r^{(t+1)}\|_2^2 = \|r^{(t)} - c a_j\|_2^2$$



$$g(c) = \|r^{(t)} - c a_j\|_2^2$$

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

$$g'(c) = -2a_j^T r^{(t)} + 2c \|a_j\|_2^2$$

$$= 0 \quad \text{when}$$

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



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

## Randomized Coordinate Descent

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



$$x^{(1)} \dots x^{(t)}$$

$$\|A\|_F^2 \quad \text{Frobenius}$$

$$= \sum_{ij} (A_{ij})^2$$

# CONVERGENCE

Claim

$$\|r^{(t+1)}\|_2^2 + \|c q_i\|_2^2 = \|r^{(t)}\|_2^2$$

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

$$\|r^{(t+1)}\|_2^2 = \|r^{(t)} - c q_i\|_2^2 = \|r^{(t)}\|_2^2 - c^2 \|q_i\|_2^2$$

$$= \|r^{(t)}\|_2^2 - \frac{(q_i^T r^{(t)})^2}{\|q_i\|_2^4} \|q_i\|_2^2$$

$$\mathbb{E} \|r^{(t+1)}\|_2^2 = \sum_{i=1}^d \frac{\|q_i\|_2^2}{\|A\|_F^2} \left( \|r^{(t)}\|_2^2 - \frac{(q_i^T r^{(t)})^2}{\|q_i\|_2^2} \right)$$

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

# CONVERGENCE

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

Claim

$$\mathbf{A}^T \mathbf{A} \mathbf{x}^* - \mathbf{A}^T \mathbf{b} = 0 \quad \forall \mathbf{b}(\mathbf{x}^*) = 0 \quad \bar{\mathbf{r}} \perp \mathbf{r}^*$$

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

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

$$= \left(1 - \frac{\lambda_{\min}(\mathbf{A}^T \mathbf{A})}{\|\mathbf{A}\|_F^2}\right) \|\bar{\mathbf{r}}^{(t)}\|_2^2$$

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

$$\|\bar{\mathbf{r}}^{(t)}\|_2^2 = \|\mathbf{r}^*\|_2^2 + \|\bar{\mathbf{r}}^{(t)}\|_2^2$$

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

$$\|\mathbf{A}^T \mathbf{r}^{(t)}\|_2^2 = \|\mathbf{A}^T \bar{\mathbf{r}}^{(t)}\|_2^2 + \mathbf{A}^T \mathbf{r}^* \rightarrow 0 = \|\mathbf{A}^T \bar{\mathbf{r}}^{(t)}\|_2^2$$

## 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:

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

*Handwritten notes:*  $\rightarrow \|\mathbf{r}^{(t)}\|_2^2 = \|\mathbf{r}^*\|_2^2 = \|\bar{\mathbf{r}}^{(t)}\|_2^2$

**Corollary:** After  $T = O\left(\frac{\|\mathbf{A}\|_F^2}{\lambda_{\min}(\mathbf{A}^T \mathbf{A})} \log \frac{1}{\epsilon}\right)$  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?

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

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

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

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

But RCD iterations are cheaper by a factor of  $d$ .

# COMPARISON

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

$$A^T A = \begin{matrix} \diagup & & \\ & \ddots & \\ & & \diagdown \end{matrix}$$

$$\begin{matrix} \longleftarrow & & \longrightarrow \\ & A^T & \\ \longrightarrow & & \longleftarrow \end{matrix}$$

orthonormal columns  
 For all  $i, j$   
 $a_j^T a_i = 0$   
 $a_j^T a_i = \text{small}$

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

$$\begin{matrix} \text{|||||} & & \\ \text{|||||} & & \\ \text{|||||} & & \\ \text{----} & & \end{matrix} = A^T A$$

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

$$a_i^T a_j = 1 \quad \text{for all } i, j$$

↓  
all ones matrix

↓  
similar vectors

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



## VISUALIZATION

Given  $f(x)$  which is potentially **non-convex**, find  $\hat{x}$  such that  $f(\hat{x}) \leq f(x^*) + \epsilon$ .



We understand very little 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.

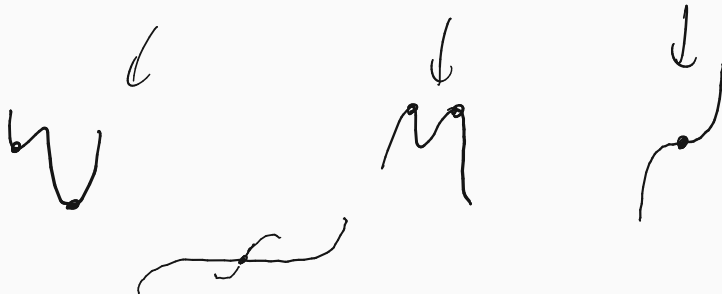
# STATIONARY POINTS

## Definition (Stationary point)

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

$$\nabla f(x) = 0 \quad \nabla f(x^*) = 0$$

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



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

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

## Definition

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

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

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

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



# GRADIENT DESCENT FINDS APPROXIMATE STATIONARY POINTS

## Theorem

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

$$1. \quad \underbrace{\nabla f(x^{(t)})^\top (x^{(t)} - x^{(t+1)})}_{= n \|\nabla f(x^{(t)})\|_2^2} - \underbrace{f(x^{(t)}) + f(x^{(t+1)})}_{= n \|\nabla f(x^{(t)})\|_2^2} \leq \frac{\beta}{2} \underbrace{\|x^{(t)} - x^{(t+1)}\|_2^2}_{= n \|\nabla f(x^{(t)})\|_2^2}$$

$$2. \quad f(x^{(t+1)}) - f(x^{(t)}) \leq \frac{\beta}{2} n^2 \|\nabla f(x^{(t)})\|_2^2 - n \|\nabla f(x^{(t)})\|_2^2 \\ = -\frac{n}{2} \|\nabla f(x^{(t)})\|_2^2 \quad \text{since } \beta = 1/n$$

$$3. \quad \frac{1}{T} \sum_{t=1}^T \frac{n}{2} \|\nabla f(x^{(t)})\|_2^2 \leq \frac{1}{T} \sum_{t=1}^T (f(x^{(t)}) - f(x^{(t+1)}))$$

$$4. \quad \min_t \|\nabla f(x^{(t)})\|_2^2 \cdot \frac{n}{2} \leq \frac{1}{T} [f(x^{(1)}) - \underbrace{f(x^{(T+1)})}_{\geq f(x^*)}] \quad \longrightarrow$$

## GRADIENT DESCENT FINDS APPROXIMATE STATIONARY POINTS

### Theorem

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

(cont.) Let  $\hat{\mathbf{x}} = \arg\min_{\mathbf{x}} \|\nabla f(\mathbf{x}^{(t)})\|_2$ .

$$\|\nabla f(\hat{\mathbf{x}})\|_2^2 \leq \frac{2}{T\eta} [f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*)]$$

$$\begin{aligned} &\hookrightarrow \\ &= \frac{2\beta}{T} \end{aligned}$$

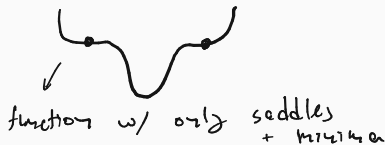
Setting  $T = \frac{2\beta}{\epsilon} \cdot [f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*)]$  gives the bound.

If GD can find a stationary point and that seems to work for your problem, 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 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 vanilla PCA.



## QUESTIONS IN NON-CONVEX OPTIMIZATION

- Can random or careful initialization lead to a good minima?
- Can we escape “shallow” local minima.
- Is a global minima even needed?