

# **EE270**

# **Large scale matrix computation, optimization and learning**

Instructor : Mert Pilanci

Stanford University

# Randomized Linear Algebra and Optimization

## Lecture 13: Gradient Descent with Momentum and Preconditioning

$$A = U \Sigma V^T \quad A^T = V \Sigma^T U^T$$

## Optimizing convex least squares cost



$$x_{is} = A^T b$$

$$\begin{aligned} A x_{is} &= A \cdot A^T b \\ &= P_{\text{Range}(A)} \cdot b \\ &= U \Sigma^T b \end{aligned}$$

- ▶ Consider

$$\min_x \underbrace{\frac{1}{2} \|Ax - b\|_2^2}_{f(x)} = (\bar{A}^T A) \cdot x - \bar{A}^T b$$

- ▶ gradient  $\nabla f(x) = A^T(\underbrace{Ax - b}_{f(x)})$   $O(nd)$  per iter

- ▶ Gradient Descent:

$$x_{t+1} = x_t - \mu A^T(Ax_t - b) \quad t=1, \dots, T$$

- ▶ fixed step size  $\mu_t = \mu$

Compute once and store the  $d \times d$  matrix

$$\bar{A}^T A: \boxed{O(nd^2)} + d^2 \cdot T$$

Sparse A: can fit A but not  $\bar{A}^T A$

not compute  $A^T A$   
 $\nearrow$   
 $nd \cdot T$

if might not  
 be possible to  
 compute  $A^T A$

load A  $n \times d$   
 $\bar{A}^T A$   $d \times d$

# Optimizing convex least squares cost

- ▶ Basic (in)equality method

- (1)  $x^*$  minimizes  $f(x)$ , hence  $\nabla f(x^*) = A^T(Ax^* - b) = 0$
- (2)  $x_{t+1} = x_t - \mu A^T(Ax_t - b)$
- (3) define error  $\Delta_t = x_t - x^*$

# Optimizing convex least squares cost

- ▶ Basic (in)equality method

- (1)  $x^*$  minimizes  $f(x)$ , hence  $\nabla f(x^*) = A^T(Ax^* - b) = 0$

- (2)  $x_{t+1} = x_t - \mu A^T(Ax_t - b)$

- (3) define error  $\Delta_t = x_t - x^*$

- ▶  $\Delta_{t+1} = \Delta_t - \mu A^T A \Delta_t$

# Optimizing convex least squares cost



- ▶ run gradient descent  $M$  iterations, i.e.,  $t = 1, \dots, M$

$$\Delta_M = (I - \mu A^T A)^M \Delta_0$$

$$\|\Delta_M\|_2 \leq \sigma_{\max}((I - \mu A^T A)^M) \|\Delta_0\|_2$$

$$\sigma_{\max}(I - \mu A^T A)^M = \max_{i=1,\dots,d} |1 - \mu \lambda_i(A^T A)|^M$$

where  $\lambda_i$  is the  $i$ -th eigenvalue in decreasing order

- ▶ Define

$\lambda_-$  as the smallest eigenvalue of  $A^T A$

$\lambda_+$  as the largest eigenvalue of  $A^T A$

$x^*$  unknown

$$\max_{\|x\|_2 \leq R} \|x\|_2 = (\lambda_{\max}(A^T A))^{1/2}$$

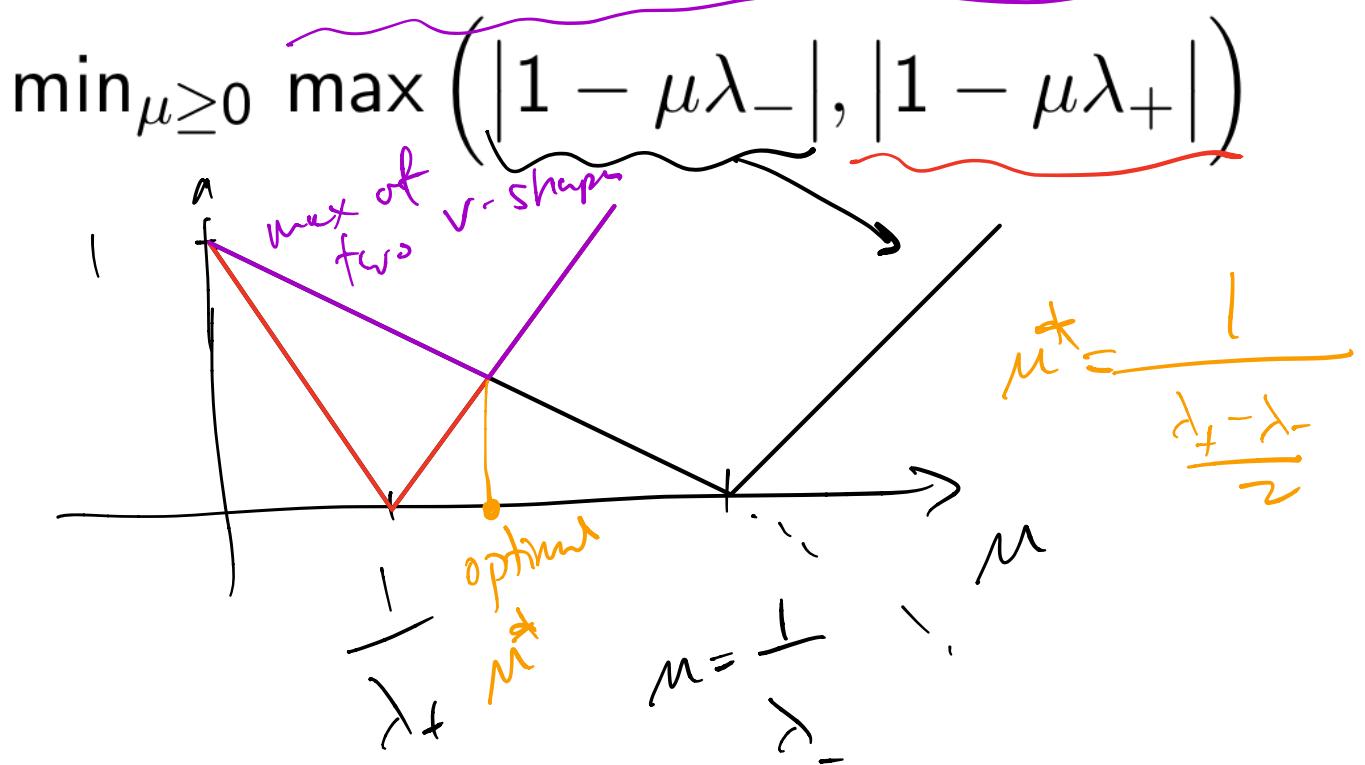
$$\max_{i=1,\dots,d} |1 - \mu \lambda_i(A^T A)| = \max(|1 - \mu \lambda_-|, |1 - \mu \lambda_+|)$$

- ▶ optimal step size that minimizes above

$$\min_{\mu \geq 0} \max(|1 - \mu \lambda_-|, |1 - \mu \lambda_+|)$$

- ▶ optimal  $\mu = \mu^*$  satisfies  $|1 - \mu^* \lambda_-| = |1 - \mu^* \lambda_+|$

$$\text{which implies } \mu^* = \frac{2}{\lambda_+ + \lambda_-}$$



$|1 - \mu \lambda_+|$  at optimum

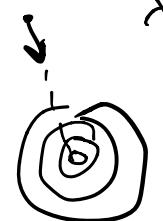
$$|1 - \mu \cdot \lambda_-| = |1 - \mu \lambda_+|$$

$$(1 - \mu \lambda_-) = (1 - \mu \lambda_+) \text{ inconsistent } \lambda_- \neq \lambda_+$$

$$1 - \mu \lambda_- = \mu \lambda_+ - 1 \Rightarrow \mu (\lambda_+ + \lambda_-) = 2 \Rightarrow \mu = \frac{2}{\lambda_+ + \lambda_-}$$

# Optimizing convex least squares cost

$$1 - \frac{2\cdot\lambda_-}{\lambda_+ + \lambda_-} = \frac{\lambda_+ - \lambda_-}{\lambda_+ + \lambda_-}$$

$$= \frac{\lambda_+ - \lambda_-}{\lambda_+ + \lambda_-}$$


- ▶ Convergence rate using  $\mu^* = \frac{2}{\lambda_+ + \lambda_-}$
- ▶  $\max(|1 - \mu^* \lambda_-|, |1 - \mu^* \lambda_+|) = \frac{\lambda_+ - \lambda_-}{\lambda_+ + \lambda_-}$
- ▶  $\|x_M - x^*\|_2 \leq \left(\frac{\lambda_+ - \lambda_-}{\lambda_+ + \lambda_-}\right)^M \|x_0 - x^*\|_2$

convergence depends on the eigenvalues of  $A^T A$

Two extremes:

- ▶ Identical eigenvalues (extremely well conditioned)  $\lambda_- = \lambda_+$ , i.e.,  $\lambda_1 = \lambda_2 = \dots = \lambda_d \implies$  convergence in one iteration
- ▶ Distant eigenvalues (poorly conditioned)  $\lambda_+ \gg \lambda_- \implies \frac{\lambda_+ - \lambda_-}{\lambda_+ + \lambda_-} \approx 1$  leads to slow convergence
- ▶ Condition number  $\kappa := \frac{\lambda_+}{\lambda_-}$   $\left(\frac{\lambda_+ - \lambda_-}{\lambda_+ + \lambda_-}\right)^M = \left(\frac{\frac{\lambda_+}{\lambda_-} - 1}{\frac{\lambda_+}{\lambda_-} + 1}\right)^M$
- ▶  $\|x_M - x^*\|_2 \leq \left(\frac{\kappa - 1}{\kappa + 1}\right)^M \|x_0 - x^*\|_2$

# Computational complexity

$$\left(\frac{\kappa-1}{\kappa+1}\right)^M \|x_0 - x^*\|_2 = \epsilon$$

$$m \cdot \log(\quad) = \log(e)$$

$$\|x_M - x^*\|_2 \leq \left(\frac{\kappa-1}{\kappa+1}\right)^M \|x_0 - x^*\|_2$$

- ▶ Initialize at  $x_0 = 0$
  - ▶ For  $\epsilon$  accuracy, i.e.,  $\|x_M - x^*\|_2 \leq \epsilon$
  - ▶ We need to set the number of iterations

$$- M \log \left( \frac{\kappa - 1}{\kappa + 1} \right) + \log \|x^*\|_2 \geq \overbrace{- \log(\epsilon)}^{\log(1/\epsilon)}$$

- $M = O\left(\frac{\log(\frac{1}{\epsilon})}{\log\left(\frac{\kappa+1}{\kappa-1}\right)}\right)$   $\log(1+x) \leq x$   $\nabla f(x_t) = \tilde{A}(A x_t - b)$
  - $\log\left(\frac{\kappa+1}{\kappa-1}\right) \approx \frac{2}{\kappa-1}$  for large  $\kappa$   $O(n \cdot d)$  per it
  - $M = O\left(\frac{\log(\frac{1}{\epsilon})}{\log\left(\frac{\kappa+1}{\kappa-1}\right)}\right) = O\left(\kappa \log\left(\frac{1}{\epsilon}\right)\right)$  for large  $\kappa$   $\xrightarrow{\text{FDT} \& G}$
  - Total computational cost  $\kappa n d \log\left(\frac{1}{\epsilon}\right)$  for  $\epsilon$  accuracy  $\xrightarrow{n d \log(1/\epsilon)}$   $\xrightarrow{n d \log n}$

O(n.d) per iter

FJLT & GD

nd log(1/e)  
+ nd log n

$$\text{nd logn} + \frac{d^2}{c^2}$$

15

GD

Kadlog(1/t)

## Improving condition number dependence: momentum

- ▶  $\min_x f(x)$
- ▶ Gradient Descent with Momentum

$$x_{t+1} = x_t - \mu_t \nabla f(x_t) + \beta_t(x_t - x_{t-1})$$

- ▶ the term  $\beta_t(x_t - x_{t-1})$  is referred to as **momentum**

# Momentum

$$\frac{d \dot{x}(t)}{dt} = \ddot{x} = f(x)$$
$$\lim_{\Delta \rightarrow 0} \frac{x_{t+\Delta} - x_t}{\Delta} = f(x)$$

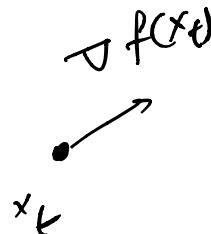
- ▶ Gradient Descent with Momentum

$$x_{t+1} = x_t - \underbrace{\mu_t \nabla f(x_t)} + \beta_t (x_t - x_{t-1})$$

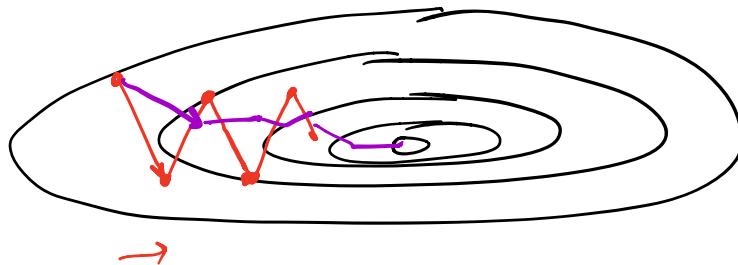
- ▶ related to a discretization of the second order ordinary differential equation

$$\ddot{x} + a\dot{x} + b\nabla f(x)$$

- ▶ which models the motion of a body in a potential field given by  $f$



# Momentum



- ▶ also called accelerated gradient descent, or heavy-ball method
- ▶ can be re-written as

$$p_t = \beta_t p_{t-1} - \nabla f(x_t)$$

$$x_{t+1} = x_t + \alpha_t p_t$$

- ▶  $p_t$  is the search direction
- ▶ there is a short-term memory
- ▶ typically we set  $p_0 = 0$

(we need to store another vector  
of size  $d$ )

$p_1 \Rightarrow$   
LS problems additional memory does not help.

Anderson Acceleration:

# Gradient Descent with Momentum for Least Squares Problems

$$\vec{A}^T (\vec{A}\vec{x}_t - \vec{b}) = \vec{A}^T (\vec{A}\vec{x}_t - \vec{A}\vec{x}^* - \vec{b}^*) = \vec{A}^T \vec{A} \cdot \vec{\Delta}_t$$

- $\min_x f(x)$  where  $f(x) = \|\vec{A}\vec{x} - \vec{b}\|_2^2$  Chebyshev Iteration

- Gradient Descent with momentum (Heavy Ball Method)

$$\underbrace{\vec{x}_{t+1}}_{\vec{\Delta}_{t+1}} = \underbrace{\vec{x}_t}_{\vec{\Delta}_t} - \mu_t \nabla f(\vec{x}_t) + \beta_t \underbrace{(\vec{x}_t - \vec{x}_{t-1})}_{\beta_t (\vec{\Delta}_t - \vec{\Delta}_{t-1})}$$

- Recall that when  $\beta = 0$  (Gradient Descent) we defined

$\vec{\Delta}_t := \vec{x}_t - \vec{x}^*$  where  $\boxed{\vec{x}^* = \vec{A}^\dagger \vec{b}}$  and established the recursion

$$\vec{\Delta}_{t+1} = \vec{\Delta}_t - \mu \cdot \vec{A}^T \vec{A} \cdot \vec{\Delta}_t + \beta_t \cdot (\vec{\Delta}_t - \vec{\Delta}_{t-1})$$

$$\boxed{\vec{\Delta}_{t+1} = (\vec{I} - \mu \vec{A}^T \vec{A}) \vec{\Delta}_t}$$

$$\|\vec{\Delta}_t\|_2 \leq v_t$$

$$v_t \rightarrow 0$$

- Since there is one time step memory, consider

$$\leq \boxed{V_t := \|\vec{\Delta}_{t+1}\|_2^2 + \|\vec{\Delta}_t\|_2^2}$$
 instead

- we can write  $V_t$  in terms of  $V_{t-1} = \|\vec{\Delta}_t\|_2^2 + \|\vec{\Delta}_{t-1}\|_2^2$

- Lyapunov analysis

$V_t$  is an energy function that decays to zero and upper-bounds error, i.e.,  $\|\vec{\Delta}_t\|_2^2 \leq V_t$

# Convergence analysis

- ▶  $\min_x f(x)$  where  $f(x) = \|Ax - b\|_2^2$
- ▶ Gradient Descent with momentum (Heavy Ball Method)

$$x_{t+1} = x_t - \mu_t \nabla f(x_t) + \beta_t (x_t - x_{t-1})$$

- ▶ let  $\Delta_t := x_t - x^*$  where  $x^* = A^\dagger b$

- ▶ note that  $b = Ax^* + b^\perp$  and  $\nabla f(x_t) = A^T A \Delta_t$

$$\tilde{\Delta}_{th} = \begin{bmatrix} \Delta_{t+1} \\ \Delta_t \end{bmatrix} = \begin{bmatrix} x_t - \overbrace{\alpha \nabla f(x_t)}^{A^T A \Delta_t} + \beta (\overbrace{x_t - x_{t-1}}{\Delta_t} - x^*) \\ \Delta_t \end{bmatrix}$$

*Augmented State*

$$= \begin{bmatrix} (1 + \beta)I - \alpha A^T A & \beta I \\ I & 0 \end{bmatrix} \begin{bmatrix} \Delta_t \\ \Delta_{t-1} \end{bmatrix}$$

$$\tilde{\Delta}_{th} = \begin{bmatrix} (1 + \beta)I - \alpha A^T A & \beta I \\ I & 0 \end{bmatrix} \cdot \tilde{\Delta}_t$$

# Convergence analysis

- iterating for  $t = 1, \dots, M$

$$\begin{bmatrix} \Delta_{M+1} \\ \Delta_M \end{bmatrix} = \begin{bmatrix} (1 + \beta)I - \alpha A^T A & \beta I \\ I & 0 \end{bmatrix}^M \begin{bmatrix} \Delta_1 \\ \Delta_0 \end{bmatrix}$$

- taking norms

$$\begin{aligned} \|\Delta_M\|_2 &\leq \left\| \begin{bmatrix} \Delta_{M+1} \\ \Delta_M \end{bmatrix} \right\|_2 = \left\| \begin{bmatrix} (1 + \beta)I - \alpha A^T A & \beta I \\ I & 0 \end{bmatrix}^M \begin{bmatrix} \Delta_t \\ \Delta_{t-1} \end{bmatrix} \right\|_2 \\ &\leq \sigma_{\max} \underbrace{\left( \begin{bmatrix} (1 + \beta)I - \alpha A^T A & \beta I \\ I & 0 \end{bmatrix}^M \right)}_{\gamma} \underbrace{\left\| \begin{bmatrix} \Delta_t \\ \Delta_{t-1} \end{bmatrix} \right\|_2}_{\sqrt{V_t}} \\ \sqrt{V_{M+1}} &\leq (\gamma) \cdot V_1 \end{aligned}$$

$$\begin{aligned} \mathcal{D} = \mathcal{D}^\top &\Rightarrow \mathcal{D} = \sum u_i u_i^\top \lambda_i = U \Lambda U^\top \\ \mathcal{D} \neq \mathcal{D}^\top &\Rightarrow \mathcal{D}^m = (U \Lambda U^\top)^m \\ \mathcal{D} = U \Sigma V & \\ \mathcal{D}^m = U \underbrace{\Sigma^m}_{\geq} V &= U \Lambda^m U^\top \end{aligned}$$

# Spectral Radius

$$M = M^\top = U \Lambda U^\top$$

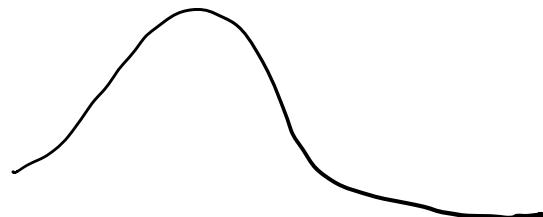
Symmetric case

$$M^k = U \Lambda^k U^\top$$
$$\sigma_{\max}(M^k) = (\sigma_{\max}(M))^k$$

- ▶ Let  $M$  be an  $d \times d$  matrix with eigenvalues  $\lambda_1, \dots, \lambda_d$
- ▶ spectral radius is defined as

$$\rho(M) := \max_{i=1, \dots, d} |\lambda_i(M)|$$

**Lemma** (Gelfand's formula)  $\lim_{k \rightarrow \infty} \sigma_{\max}(M^k)^{\frac{1}{k}} = \rho(M)$



$$\sigma_{\max}(M^k) \approx (\rho(M))^k$$

$k \rightarrow \infty$

asymptotic case

- ▶ Let  $\lambda_i$  denote the eigenvalues of  $A^T A$  for  $i = 1, \dots, d$
- ▶ **Lemma** The eigenvalues of

$$\begin{bmatrix} (1 + \beta)I - \alpha A^T A & \beta I \\ I & 0 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \lambda \cdot \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

are given by the eigenvalues of  $2 \times 2$  matrices

$$\begin{bmatrix} 1 + \beta - \alpha \lambda_i & -\beta \\ 1 & 0 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \lambda \cdot \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$$

- ▶ for  $i = 1, \dots, d$
  - ▶ These are given by the roots of  $\overbrace{u^2 - (1 + \beta - \alpha \lambda_i)u + \beta}^{= 0}$
  - ▶ setting  $\alpha = \frac{4}{\sqrt{\lambda_+} + \sqrt{\lambda_-}}$  and  $\beta = \frac{\sqrt{\lambda_+} - \sqrt{\lambda_-}}{\sqrt{\lambda_+} + \sqrt{\lambda_-}}$  yields (minimized the abs of the max. root)
  - ▶ spectral radius:  $\rho \left( \begin{bmatrix} (1 + \beta)I - \alpha A^T A & \beta I \\ I & 0 \end{bmatrix} \right) = \frac{\sqrt{\lambda_+} - \sqrt{\lambda_-}}{\sqrt{\lambda_+} + \sqrt{\lambda_-}}$
- (  $\frac{\sqrt{\lambda_+} - \sqrt{\lambda_-}}{\sqrt{\lambda_+} + \sqrt{\lambda_-}}$  )  $\downarrow$  as  $k \rightarrow \infty$
- Convergence  
rate

## Convergence result

$$\boxed{\text{GD } (\beta=0) \quad \|x_M - x^*\|_2 \leq \left( \frac{\lambda_+ - \lambda_-}{\lambda_+ + \lambda_-} \right)^M \|x_0 - x^*\|_2}$$

- ▶ setting  $\alpha = \frac{4}{\sqrt{\lambda_+} + \sqrt{\lambda_-}}$  and  $\beta = \frac{\sqrt{\lambda_+} - \sqrt{\lambda_-}}{\sqrt{\lambda_+} + \sqrt{\lambda_-}}$  yields

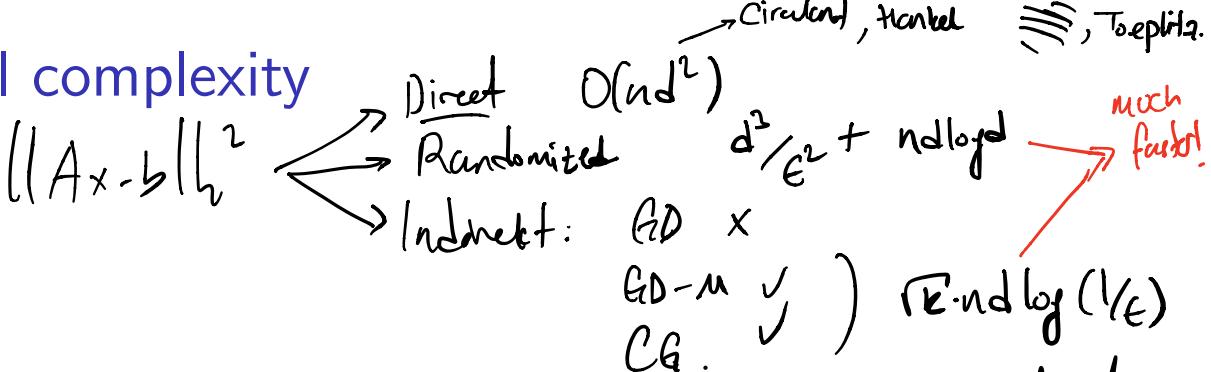
$$\|\Delta_m\|_2 \leq \left\| \begin{bmatrix} \Delta_{m+1} \\ \Delta_m \end{bmatrix} \right\|_2 \leq \dots \left( \frac{\sqrt{\lambda_+} - \sqrt{\lambda_-}}{\sqrt{\lambda_+} + \sqrt{\lambda_-}} \right)^M \left\| \begin{bmatrix} \Delta_1 \\ \Delta_{1-1} \end{bmatrix} \right\|_2$$

$$k = \frac{\lambda_+}{\lambda_-} \quad = \quad \left( \frac{\sqrt{k} - 1}{\sqrt{k} + 1} \right)^M \cdot \left\| \begin{bmatrix} \Delta_1 \\ \Delta_{1-1} \end{bmatrix} \right\|_2 = \epsilon$$

solve for  $M$ :  $M \cdot \log\left(\frac{\sqrt{k} + 1}{\sqrt{k} - 1}\right) = \log \frac{1}{\epsilon}$

$$M = \frac{\log(1/\epsilon)}{\log\left(\frac{\sqrt{k} + 1}{\sqrt{k} - 1}\right)} = \sqrt{k} \cdot \log(1/\epsilon)$$

# Computational complexity



- ▶ Gradient Descent ( $\beta = 0$ ) total computational cost  
 $\boxed{\kappa nd \log(\frac{1}{\epsilon})}$  for  $\epsilon$  accuracy
- ▶ Gradient Descent with Momentum total computational cost  
 $\boxed{\sqrt{\kappa} nd \log(\frac{1}{\epsilon})}$  for  $\epsilon$  accuracy
- ▶ we need to know eigenvalues of  $A^T A$  to find optimal step-sizes

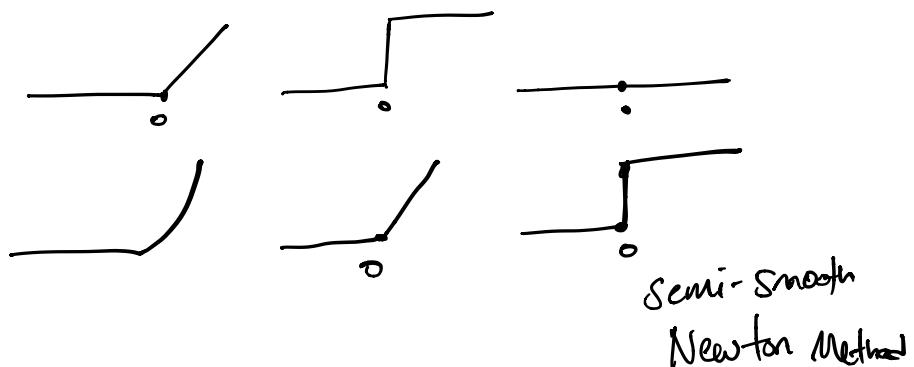
Conjugate gradient method (automatically) tunes the step size.

- ▶ Higher order momentum can not do better than  $\sqrt{\kappa} nd \log(\frac{1}{\epsilon})$ .
- ▶ Sketching can improve

# Computational complexity

- ▶ Gradient Descent ( $\beta = 0$ ) total computational cost  $\kappa nd \log(\frac{1}{\epsilon})$  for  $\epsilon$  accuracy
- ▶ Gradient Descent with Momentum total computational cost  $\sqrt{\kappa}nd \log(\frac{1}{\epsilon})$  for  $\epsilon$  accuracy
- ▶ we need to know eigenvalues of  $A^T A$  to find optimal step-sizes
- ▶ Conjugate Gradient doesn't require the eigenvalues explicitly and results in  $\sqrt{\kappa}nd \log(\frac{1}{\epsilon})$  operations

## Newton's Method



- ▶ Suppose  $f$  is twice differentiable, and consider a second order Taylor approximation at a point  $x_t$

$$f(y) \approx f(x_t) + \underbrace{\nabla f(x_t)^T (y - x_t)}_{\text{smooth part}} + \frac{1}{2}(y - x^t) \nabla^2 f(x^t) (y - x^t)$$

- ▶ and minimize the approximation

$$\nabla_y (\quad ) = 0$$

## Newton's Method

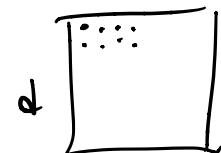
$$f(x) = \|Ax - b\|^2$$

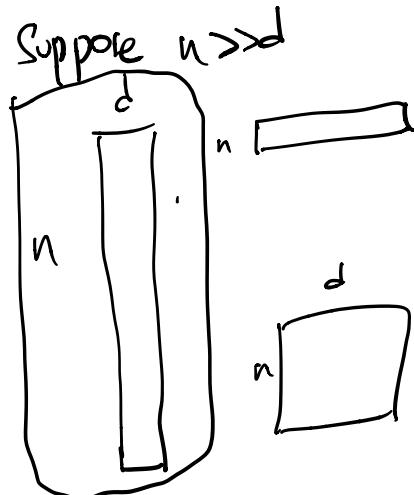
$$[\nabla^2 f(x)]_{ij} = [A^T A]_{ij} = A^{(i)\top} A^{(j)}$$

$A: n \times d$

►  $x_{t+1} = x_t - \mu_t (\nabla^2 f(x))^{-1} \nabla f(x)$

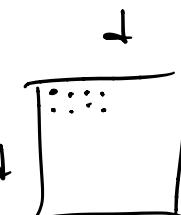
► complexity:

Compute  $\nabla^2 f(x)$    $O(n d^2)$

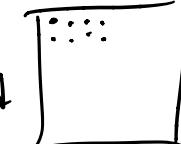


Invert  $(\nabla^2 f(x))^{-1}$

$O(d^3)$



$d$



$d$



$d$



$d$



$d$



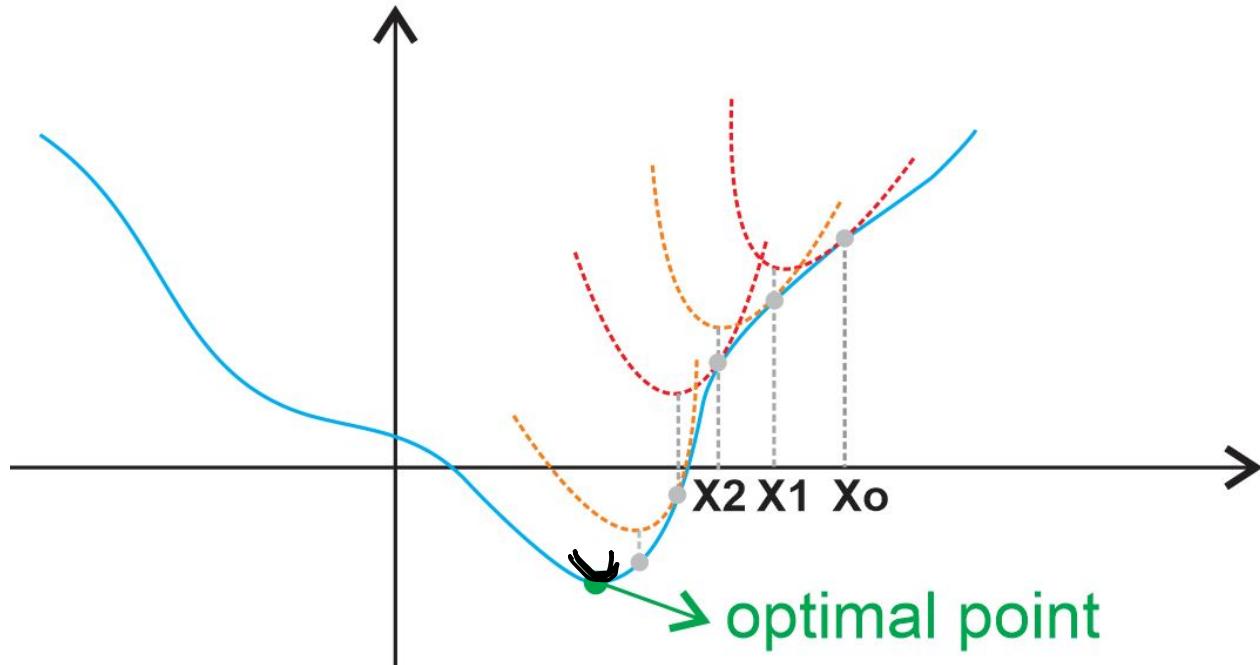
$d$



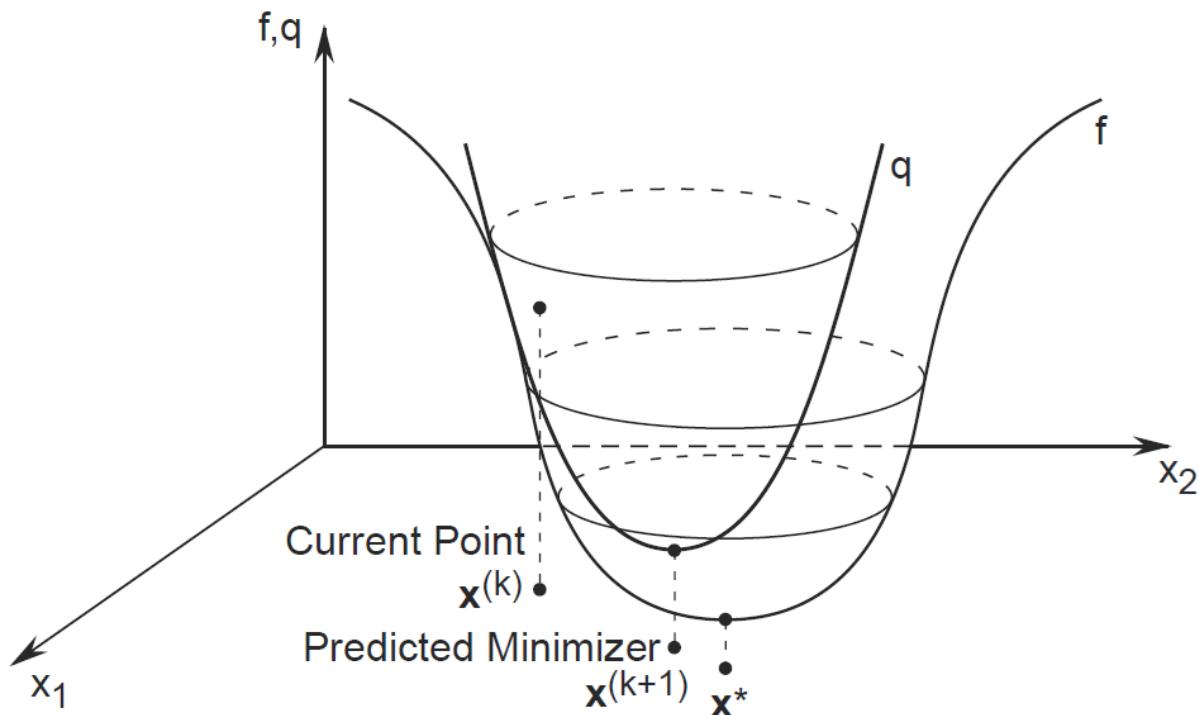
$d$



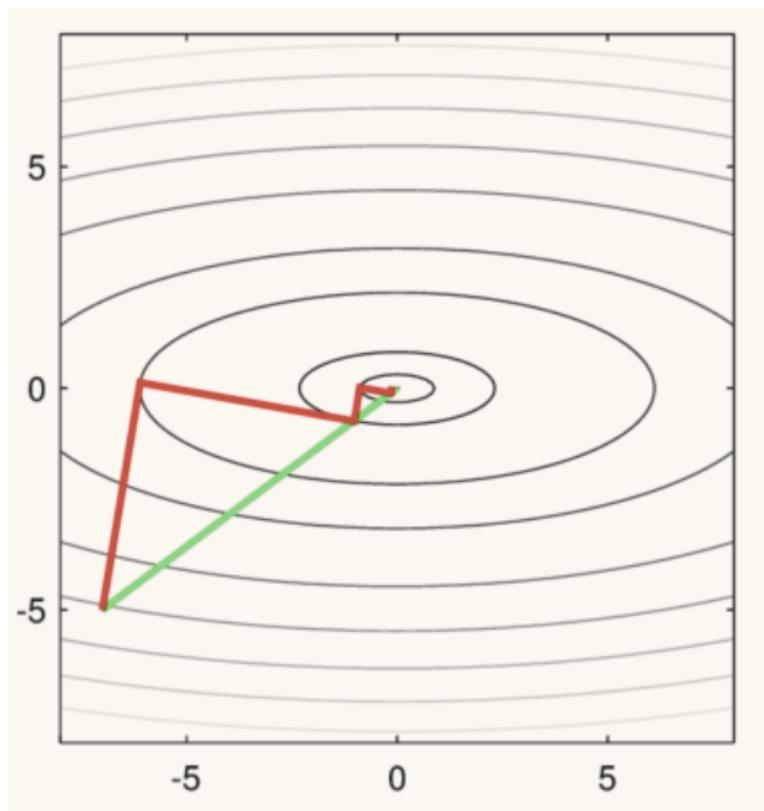
# Newton's Method in one dimension



# Newton's Method in higher dimensions



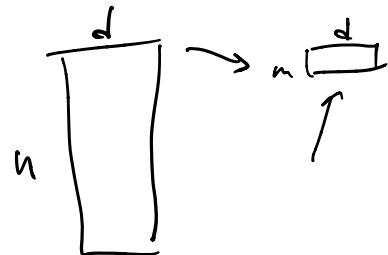
# Newton's Method vs Gradient Descent



# Newton's Method for least squares converges in one step

- ▶ Consider

$$\min_x \underbrace{\frac{1}{2} \|Ax - b\|_2^2}_{f(x)}$$



- ▶ gradient  $\nabla f(x) = A^T(Ax - b)$
- ▶ Hessian  $\nabla^2 f(x) = A^T A$
- ▶ Gradient Descent:

$$x_{t+1} = x_t - \mu A^T (Ax_t - b)$$

- ▶ Newton's Method:

$$(A^T A)^{-1} \cdot A^T (Ax_t - b) \quad RNM$$

$$x_{t+1} = x_t - \underbrace{\mu (A^T A)^{-1} A^T (Ax_t - b)}_{\text{fixed step size}}$$

- ▶ fixed step size  $\mu_t = \mu$

$$x_{th} = x_t - \mu x_t + \mu \underbrace{(A^T A)^{-1} A^T b}_{x^*} \quad \begin{array}{l} \text{if } \mu = 1 \\ \Rightarrow x_{th} = x^* \end{array}$$

# Questions?