# EE270
# Large scale matrix computation, optimization and learning

Instructor : Mert Pilanci

Stanford University

Randomized Linear Algebra and Optimization
Lecture 12: Gradient Descent

# Summary of randomized least squares solvers

$$Sx: \ O(n\log n)$$

$$SA = [SA^{(1)} \cdots - SA^{(s)}] \quad O(d \cdot n\log n)$$
(near - linear time)

$S = PHD \quad (FJLT \text{ or } SRHT)$

subsampler)

$A = U\Sigma V^T$

AMM for

$U^T S^T S U \approx U^T U$

▶ Left Sketch

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

▶ $\min_x \|S(Ax - b)\|_2^2$

▶ Fast Johnson Lindenstrauss Transform (Randomized Hadamard Transform) / Sparse JL: [ ⠋⠉⠊ ]
SA and Sb can be computed in $O(nd\log n)$ time    $O(nd\log n)$

▶ Gaussian sketch / $\pm 1$
$= O(nd^2)$
SA and Sb can be computed in $O(ndm)$ time

▶ total complexity:

$m > d$

pick $m = \dfrac{d}{\epsilon^2}$ $\Rightarrow$ we'll have $\epsilon$-approximate LS solution

$\boxed{\dfrac{d^3}{\epsilon^2} + nd\log n}$ compare with $nd^2$

# Summary of randomized least squares solvers

*(handwritten annotation: $d$, a rectangle labeled $n$ on the left and $s$ on the right, with $n\,d\log d$ below)*

▶ Right Sketch

$$\min_{Ax=b} \|x\|_2^2$$

▶ $\min_{ASz=b} \|z\|_2^2$

▶ Fast Johnson Lindenstrauss Transform (Randomized Hadamard Transform)

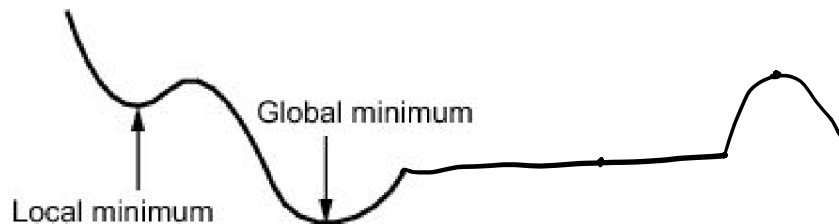  $AS$ can be computed in $O(ndlogd)$ time

▶ Gaussian sketch

  $AS$ can be computed in $O(ndm)$ time

▶ total complexity:   $ndlogd + \dfrac{n^2}{\epsilon^2}$

  $m = \dfrac{n}{\epsilon^2}.$

# Optimization: Gradient Descent



- ▶ Consider unconstrained minimization of $f : \mathbb{R}^d \to \mathbb{R}$, differentiable function
- ▶ we want to solve

$\mu_t$ : step size at iteration $t$

$$\min_{x \in \mathbb{R}^d} f(x)$$

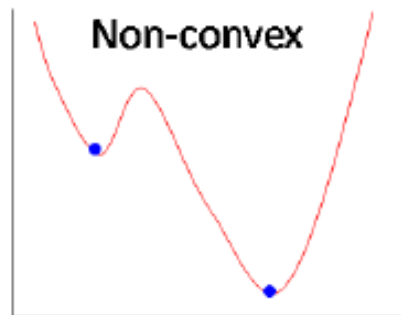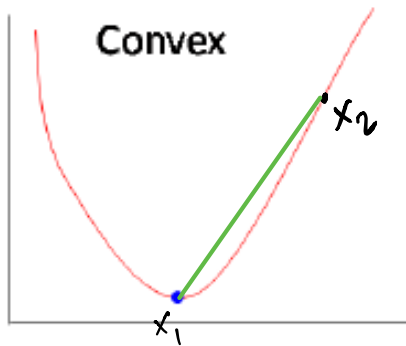- ▶ **Gradient descent:** choose initial $x_0 \in \mathbb{R}^d$ and repeat

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

- ▶ for $t = 1, ..., T$

# Convex vs Non-convex functions
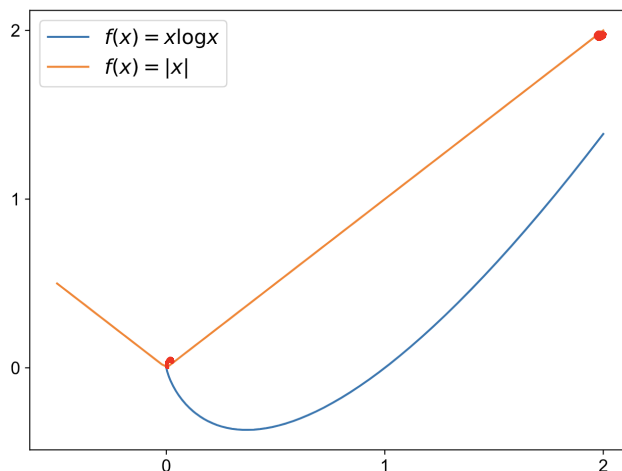
▶ a function $f$ is called **convex** if

$$\forall x_1, x_2 \in \mathcal{X},\ \forall t \in [0,1]: \quad f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$



Convex     Non-convex
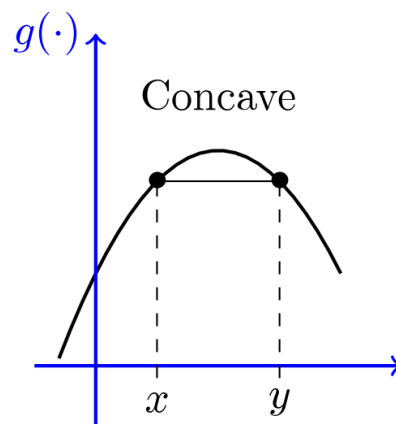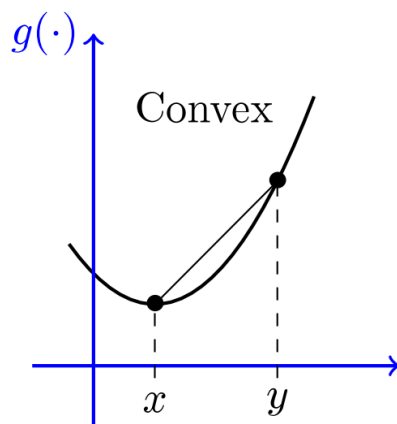
# Convex vs Non-convex functions

▶ a function $f$ is called **strictly convex** if

$$\forall x_1 \neq x_2 \in \mathcal{X}, \forall t \in [0,1]: \quad f(tx_1 + (1-t)x_2) < tf(x_1) + (1-t)f(x_2)$$

# Concave functions

▶ a function $f$ is called (strictly) **concave** if
  $-f$ is (strictly) convex

# Differentiable functions

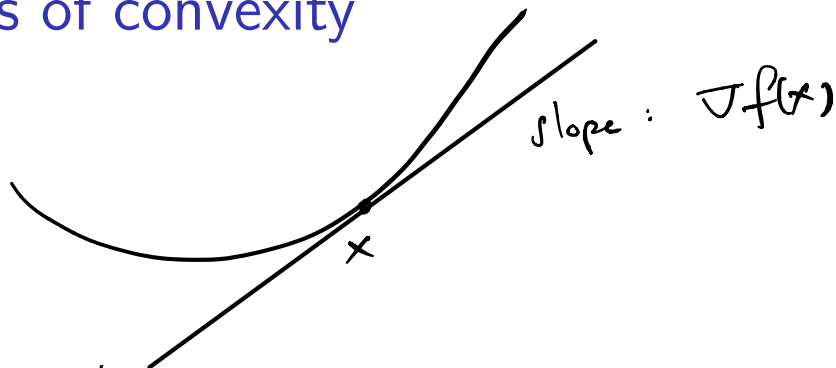▶ A one dimensional function $f : \mathbb{R} \to \mathbb{R}$ is differentiable if the derivative
$f'(x) := \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$ exists

▶ Suppose that all partial derivatives of $f : \mathbb{R}^d \to \mathbb{R}$ exists
The gradient $\nabla f(x)$ is the vector of partial derivatives
$[\nabla f(x)]_i = \frac{\partial}{\partial x_i} f(x)$

# Alternative definitions of convexity

slope : $\nabla f(x)$

x

▶ Assume that $f(x) : \mathbb{R}^d \to \mathbb{R}$ is differentiable. Then $f$ is convex, if and only if for every $x, y$ the inequality

$$f(y) \geq f(x) + \nabla f(x)^T (y - x)$$
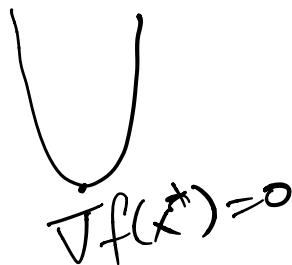
is satisfied

first order Taylor approx at x

# Twice differentiable functions

▶ Suppose that all second derivatives of $f : \mathbb{R}^d \to \mathbb{R}$
$\frac{\partial}{\partial x_i} \frac{\partial}{\partial x_j} f(x)$ exists

The Hessian $\nabla^2 f(x)$ is the matrix of partial derivatives
$[\nabla^2 f(x)]_{ij} = \frac{\partial}{\partial x_i} \frac{\partial}{\partial x_j} f(x)$

# Twice differentiable convex functions


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

- ▶ A twice differentiable function $f(x)$ is convex if and only if the Hessian $\nabla^2 f(x)$ is positive semi-definite for all $x \in \mathbb{R}^d$
- ▶ Suppose that $f$ is convex and differentiable, then $x^*$ is a global minimizer of $f$ if and only if $\nabla f(x^*) = 0$
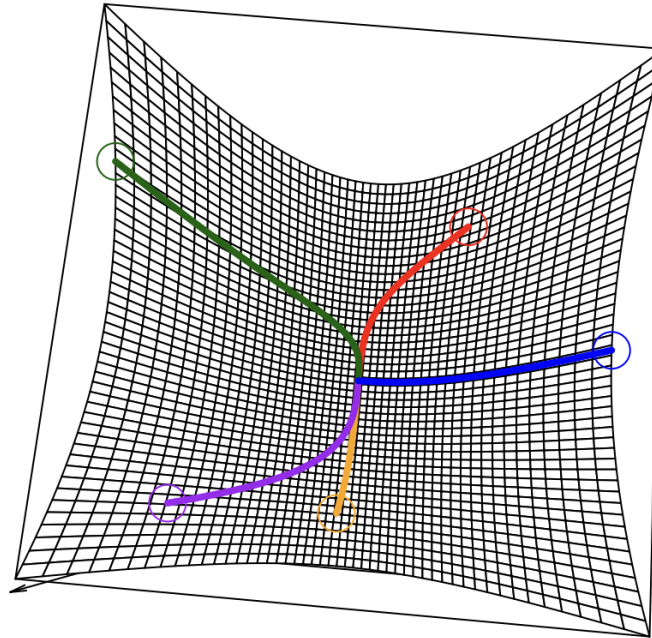
# Gradient descent for differentiable functions

▶ $-\nabla f(x)$ is the direction of largest instantaneous decrease
▶ Gradient Descent (GD):

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

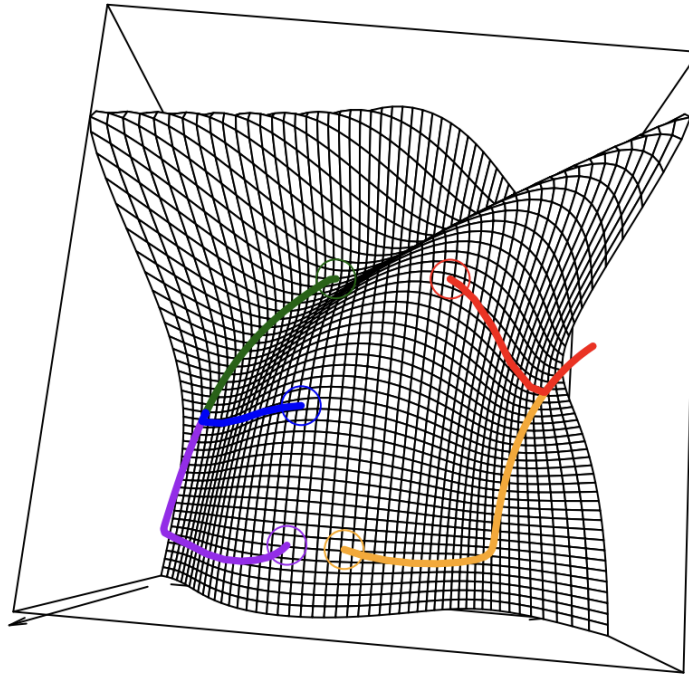▶ where $\mu_t$ is the step size at iteration $t$.
▶ if $\mu_t$ is sufficiently small and $\nabla f(x_t) \neq 0$, guaranteed to decrease the value of $f$
▶ If $f$ is convex, converges to **global minimum** under mild conditions
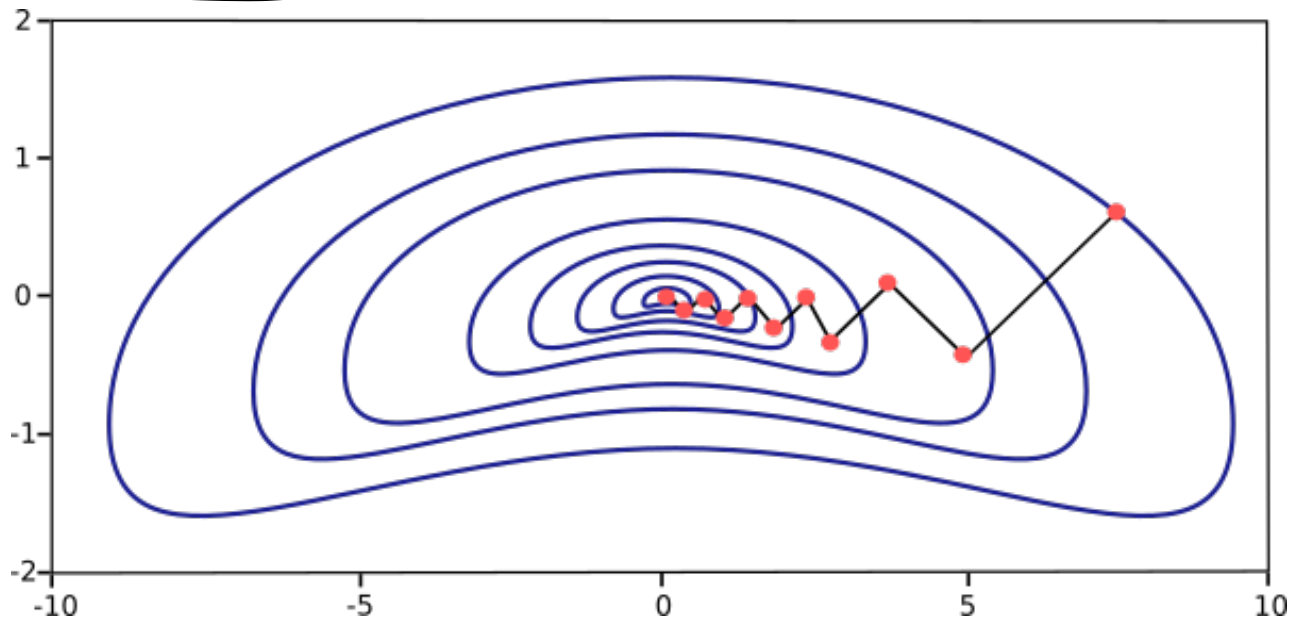
# Gradient descent for convex functions



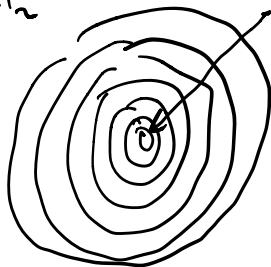slide credit: R. Tibshirani

# Gradient descent for non-convex functions

slide credit: R. Tibshirani

# Gradient descent iterations $J \sim \|Ax-b\|_2^{\sim}$

# Gradient descent on highly curved functions

▶ Rosenbrock function (non-convex)
  $$f(x_1, x_2) = (a - x_1)^2 + b(x_2 - x_1^2)^2$$
  where $a$ and $b$ are parameters, e.g., $a = 1, b = 100$
  has a global minimum at $(x_1, x_2) = (a, a^2)$

# Optimizing convex least squares cost

▶ Consider

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

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

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

▶ fixed step size $\mu_t = \mu$

# Optimizing convex least squares cost

$$\bar{A}^T b = \bar{A}^T A x^*$$

▶ Basic (in)equality method

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

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

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

$$\underbrace{x_{t+1} - x^*}_{\Delta_{t+1}} = \underbrace{x_t - x^*}_{\Delta_t} - \mu \, \bar{A}^T(Ax_t - b)$$

$$\underbrace{\bar{A}^T A x_t - \bar{A}^T A x^*}_{} = \bar{A}^T A \cdot \Delta_t$$

$$\boxed{\Delta_{t+1} = (I - \mu \cdot \bar{A}^T A) \cdot \Delta_t}$$

$$\Delta_1 = (I - \mu \bar{A}^T A) \Delta_0$$
$$\Delta_2 = (I - \mu \bar{A}^T A) \cdot \Delta_1 = (I - \mu \bar{A}^T A)^2 \cdot \Delta_0$$
$$\Delta_T = (I - \mu \cdot \bar{A}^T A)^T \cdot \Delta_0$$

# 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

$$\|\mathcal{O}\|_2 = \max_{\|x\|_2 \leq 1} \|\mathcal{O} \cdot x\|_2 \implies \|\Delta_M\|_2 \leq \|(I - \mu A^T A)^M\|_2 \|\Delta_0\|$$

$$\underbrace{\sigma_{max} \cdot \left((I - \mu A^T A)^M\right)}$$

$$= \max_k |\lambda_k \left((I - \mu A^T A)^M\right)|$$

- run gradient descent $M$ iterations, i.e., $t = 1, ..., M$
- $\Delta_M = (I - \mu A^T A)^M \Delta_0$

$$= \max_k \cdot |1 - \mu \lambda_k (A^T A)|^M$$

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

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

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

$$\max_i |1 - \mu \lambda_i| < 1$$

$$\|Ax - b\|_h^2 = x^T A^T A x - 2 b^T A x + b^T b$$

$$\nabla^2 = 2 \cdot A^T A$$

$$\max_{i=1,\ldots,d} \max\left(1 - \mu \lambda_i, \ \mu \lambda_i - 1\right) < 1$$

$$\text{convergence condition}$$

# Optimizing convex least squares cost

# Questions?