## Outline:

Convex Optimization

Convergence of GD

## Optimization and machine learning

Data $\{(x_i, y_i)\}_{i=1\cdots n}$

Consider a model $\hat{y}_\theta(x_i)$

$$\min_\theta \sum_{i=1}^n \ell(\hat{y}_\theta(x_i), y_i)$$

## Optimization in general

$$\min_x f(x)$$
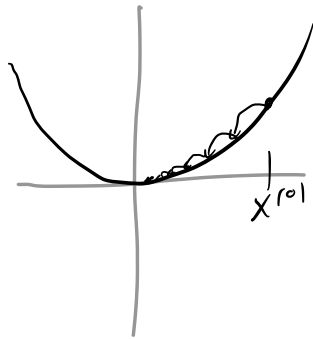
Gradient descent: Take successive steps "downhill"

$$x^{(i+1)} = x^{(i)} - \alpha \nabla f(x^{(i)})$$

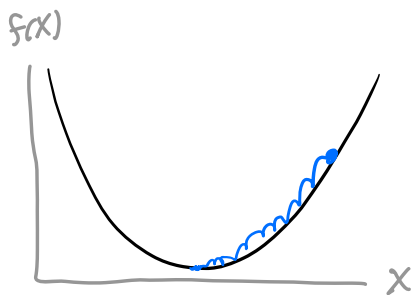iteration index — (step size, learning rate) — $-\nabla f$ points in direction of steepest descent

$$f: \mathbb{R} \to \mathbb{R}$$

Example: $f(x) = \frac{1}{2} L x^2$.



If $\alpha < \frac{2}{L}$, GD converges.

Picture:
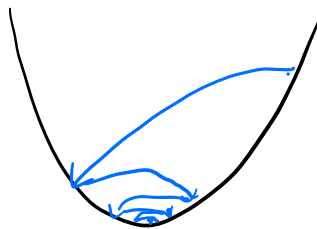


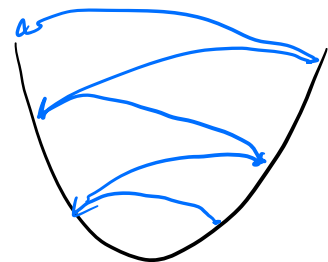Small learning rate

$$\alpha < \frac{1}{L}$$

medium learning rate

$$\frac{1}{L} < \alpha < \frac{2}{L}$$

high learning rate

$$\alpha > \frac{2}{L}$$

# Challenges of gradient descent in machine learning & minibatches

$$\min_\theta \underbrace{\frac{1}{n} \sum_{i=1}^{n} \ell(\hat{y}_\theta(x_i), y_i)}_{f(\theta)}$$

$$\theta^{k+1} = \theta^k - \alpha \nabla f(\theta) = \theta^k - \alpha \frac{1}{n} \sum_{i=1}^{n} \nabla_\theta \ell(\hat{y}_\theta(x_i), y_i)$$

To evaluate $\nabla f(\theta)$, one needs to loop through all data ( batch gradient descent )

- expensive
- not possible in some contexts

Idea: use minibatches
   Select a minibatch $B \subset \{1, 2, \cdots, n\}$

$$\theta^{k+1} = \theta^k - \alpha \underbrace{\frac{1}{|B|} \sum_{i \in B} \nabla_\theta \ell(\hat{y}_\theta(x_i), y_i)}$$
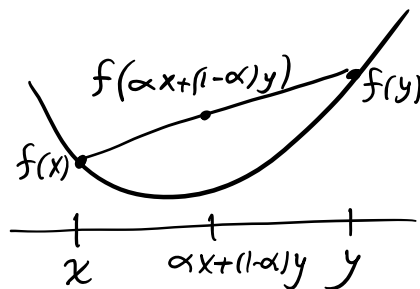
use as approximation of $\nabla_\theta f(\theta)$

# Convex Optimization

We say $f: \mathbb{R}^d \to \mathbb{R}$ is **Convex** if

$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha) f(y)$$

for all $0 \leq \alpha \leq 1$, $x, y$.

Convex Combo

"always curves up"

$$f(\alpha x + (1-\alpha)y) \quad f(y)$$
$$f(x)$$

$$x \qquad \alpha x + (1-\alpha)y \quad y$$

Examples:
$f: \mathbb{R} \to \mathbb{R}$
$f(x) = x^2$     is or is not convex

Fix a $c \in \mathbb{R}$.
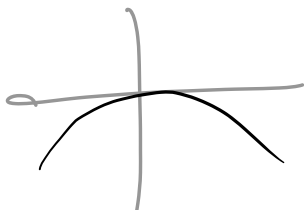$f: \mathbb{R} \to \mathbb{R}$
$f(x) = cx^2$

is or is not convex

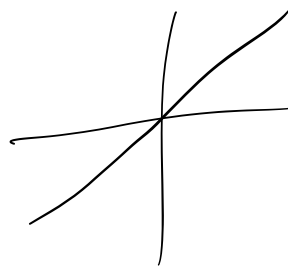If $c \geq 0$, yes
$c < 0$, no
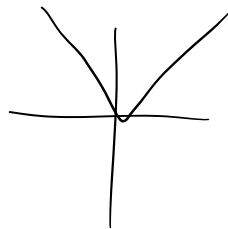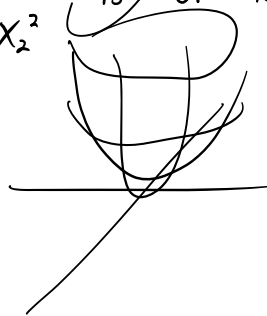
$f: \mathbb{R} \to \mathbb{R}$

$f(x) = x$

(is) or is not   convex

$f: \mathbb{R} \to \mathbb{R}$

$f(x) = |x|$

(is) or is not   convex

$f: \mathbb{R}^2 \to \mathbb{R}$

$f(x) = \|X\|^2 = X_1^2 + X_2^2$

(is) or is not   convex

$f: \mathbb{R}^2 \to \mathbb{R}$

$f(x) = X_1^2$

(is) or is not   convex

We will study the minimization of convex functions.

Does every convex function f have a minimal value?

$$\min_{x} f(x)$$

No ? $f(x) = x$

or $= e^x$

## All local minima of convex functions are global minima. [i]



local
min

global
min

not
convex

Suppose $X_*$ is a local min of f.
If $X \approx X_*$ the $f(x) \geqslant f(X_*)$.

Suppose $X_*$ is not a global min.

Suppose $f(\hat{x}) < f(x^*)$.

By convexity, $f(x)$ lies
below dotted line between
$X^*$ and $\hat{x}$. So $X^*$ not a
local min.
           Contradiction. ✗



$f(\hat{x})$

$f(\hat{x})$

$x^*$    $\hat{x}$

# Convexity and Second derivatives

## Functions of one variable

If $f: \mathbb{R} \to \mathbb{R}$ is twice differentiable everywhere, $f$ is convex if and only if $f''(x) \geq 0$ for all $x$.



## Functions of multiple variables

Let $f: \mathbb{R}^n \to \mathbb{R}$ be twice differentiable

$f$ is convex if $D^2 f = Hf$ is positive semidefinite everywhere

Hessian matrix

$$D^2f = Hf(x) = \begin{pmatrix} \dfrac{\partial^2 f}{\partial x_1^2} & \cdots & \dfrac{\partial^2 f}{\partial x_n \partial x_1} \\ \vdots & & \\ \dfrac{\partial^2 f}{\partial x_1 \partial x_n} & \cdots & \dfrac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

$H$ is <u>positive definite</u> if all eigenvalues are positive

$H$ is <u>positive semidefinite</u> if all eigenvalues are nonnegative

Eigenvalue Decomposition:

If $H \in \mathbb{R}^{n \times n}$ is symmetric ($H^t = H$), then

$H$ has an orthonormal basis of eigenvectors with real eigenvalues. So

$$H = U \Lambda U^t \qquad \text{where } U \text{ has orthonormal columns}$$

$n \times n$

diagonal $n \times n$

⟍ could have negative entries

We say $U_i$ is an eigenvector of $H$ with eigenvalue $\lambda_i$ if $H U_i = \lambda_i U_i$

$$H = \begin{pmatrix} | & | & & | \\ U_1 & U_2 & \cdots & U_n \\ | & | & & | \end{pmatrix} \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix} \begin{pmatrix} - & U_1^t & - \\ - & U_2^t & - \\ & \vdots & \\ - & U_n^t & - \end{pmatrix}$$

$$U \quad \cdot \quad \underline{\Lambda} \quad \cdot \quad U^t$$

Columns are
unit length eigenvectors
that are orthogonal to
each other

$$U_i \cdot U_j = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

We also have

$$H = \sum_{i=1}^{n} \lambda_i U_i U_i^t.$$

Why?
$$H = \begin{pmatrix} | & & | \\ U_1 & \cdots & U_n \\ | & & | \end{pmatrix} \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \begin{pmatrix} - & U_1^t & - \\ & \vdots & \\ - & U_n^t & - \end{pmatrix}$$

$$= \begin{pmatrix} | & & | \\ U_1 & \cdots & U_n \\ | & & | \end{pmatrix} \begin{pmatrix} - \lambda_1 U_1^t - \\ \vdots \\ - \lambda_n U_n^t - \end{pmatrix} = \sum_{i=1}^{n} U_i \left( \lambda_i U_i^t \right)$$

Theorem: H is positive semidefinite if and only if

$$z^t H z \geqslant 0 \quad \text{for all } z \in \mathbb{R}^n$$

Recall, because H is symmetric $(H = H^t)$, H has an orthonormal basis of eigenvectors with real eigenvalues. So

$$H = U \Lambda U^t \quad \text{where } U \text{ has orthonormal columns}$$

and

$\Lambda$ is diagonal

Proof of Theorem: • PSD $\Rightarrow$ $Z^t H Z \geqslant 0$ for all $Z$

As $H$ is PSD, $\Lambda$ has nonneg. diagonal entries. So $Z^t H Z = Z^t U \Lambda U^t Z$

$$= \sum_{i=1}^{n} \Lambda_{ii} \left( U^t Z \right)_i^2$$

$$\geqslant 0.$$

• $Z^t H Z \geqslant 0$ for all $Z \Rightarrow H$ is PSD

Suppose $H$ is not PSD. At least one eigenvalue is negative. Suppose $U_i$ is eigenvector w/ e-val $\lambda_i < 0$. Then let $Z = U_i$.

$$Z^t H Z = U_i^t H U_i = \lambda_i U_i^t U_i < 0.$$

**Many but not all ML optimization problems are convex.**

Convex Problems:  least squares regression, logistic regression,

Not convex: neural networks

locally convex

# How fast does gradient descent converge?

$$\min_{x} f(x), \qquad x^{(i+1)} = x^{(i)} - \alpha \nabla f(x^{(i)})$$

Suppose $x^{(i)} \to x^*$ as $i \to \infty$.

How long do you need to wait to get a certain accuracy $\varepsilon$?

Can gain understanding in some convex cases.

# Convergence of GD for quadratic functions

Let $f(x) = \frac{1}{2} x^t Q x - b^t x$

where $x \in \mathbb{R}^d$, $b \in \mathbb{R}^d$, $Q \in \mathbb{R}^{d \times d}$ is positive definite

Let $m = \lambda_{min}(Q)$, $M = \lambda_{max}(Q)$, $K = \frac{M}{m}$

$\llcorner$ condition number of $Q$

Consider GD w/ fixed step size $\alpha$

$$x^{k+1} = x^k - \alpha \nabla f(x^k)$$

Note: $x^* = Q^{-1} b$ is the unique global min of $f$

**Analytically show that this is the solution to the problem**

$$\nabla f(x) = Qx - b = 0$$

$$Qx = b \implies x = Q^{-1} b$$

Theorem: If $\alpha = \frac{2}{M+m}$, then GD

for $f(x) = \frac{1}{2} x^t Q x - b^t x$ satisfies

$$\| x^k - x^* \| \leq \left( \frac{1 - 1/k}{1 + 1/k} \right)^k \| x^0 - x^* \|$$

"first-order convergence"
Error decays exponentially

To get error $\varepsilon$, need $O\left( \log \left( \varepsilon^{-1} \right) \right)$ iterations


Proof: Note $\nabla f(x) = Qx - b$.
The global minimizer solves $Qx^* = b \Rightarrow x^* = Q^{-1} b$

$$x^{k+1} - x^* = x^k - \alpha \nabla f(x^k) - x^*$$
$$= x^k - \alpha (Qx^k - b) - x^*$$
$$= x^k - \alpha (Qx^k - Qx^*) - x^*$$
$$= (I - \alpha Q)(x^k - x^*)$$

So,
$$\| x^{k+1} - x^* \| \leq \| I - \alpha Q \| \; \| x^k - x^* \|$$

largest
Eigenvalue (in terms of abs value)
of $I - \alpha Q$

$= \max(\alpha M - 1, 1 - \alpha m)$

We choose $\alpha = \dfrac{2}{M+m}$.

So $\|I - \alpha Q\| = \dfrac{M-m}{M+m} = \dfrac{1 - 1/k}{1 + 1/k} < 1$

$\Rightarrow \|x^{k+1} - x^*\| \le \left(\dfrac{1 - 1/k}{1 + 1/k}\right) \|x^k - x^*\|$

$\Rightarrow \|x^k - x^*\| \le \left(\dfrac{1 - 1/k}{1 + 1/k}\right)^k \|x^0 - x^*\|$ ■

## Interpretation:

If $f$ doesn't curve up too much
and doesn't curve up too little,
then GD with fixed step size

can exhibit first order convergence
to the global minimizer

**Should we think of GD as converging "quickly"?**

If the function is quadratic, then GD (with the right step size) can converge very quickly.

If the function is not quadratic, then GD may converge slowly

Theorem: Let $f$ be convex and $\lambda_{max}(Hf(x)) \le M$ for all $x$. If $\alpha \le \frac{1}{M}$, then GD satisfies

$$f(x^{(i)}) - f(x^*) \le \frac{1}{2i\alpha} \|x^{(0)} - x^*\|^2$$

where $x^*$ is a minimizer of $f$.

— Error decays <u>slowly</u>

— To get error $\varepsilon$ from optimal value, need $O(\varepsilon^{-1})$ iterations

# Summary:

- Too large learning rate can lead to divergence
- In convex case, to get convergence $\alpha$ should be small relative to curvature of $f$
- Too small learning rate can lead to slow convergence

- For convex quadratic functions, convergence of GD can be first order (fast)

- For more general convex functions, convergence can be slow

- SGD w/ fixed step size is not expected to converge

- SGD with decaying step sizes may converge