

Day 16 - Gradient Descent and Stochastic Gradient Descent

Outline:

- Gradient Descent (GD)
- Stochastic Gradient Descent (SGD)
- Convex Optimization
- Convergence of GD

Optimization and machine learning

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

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

$$\min_x f(x)$$

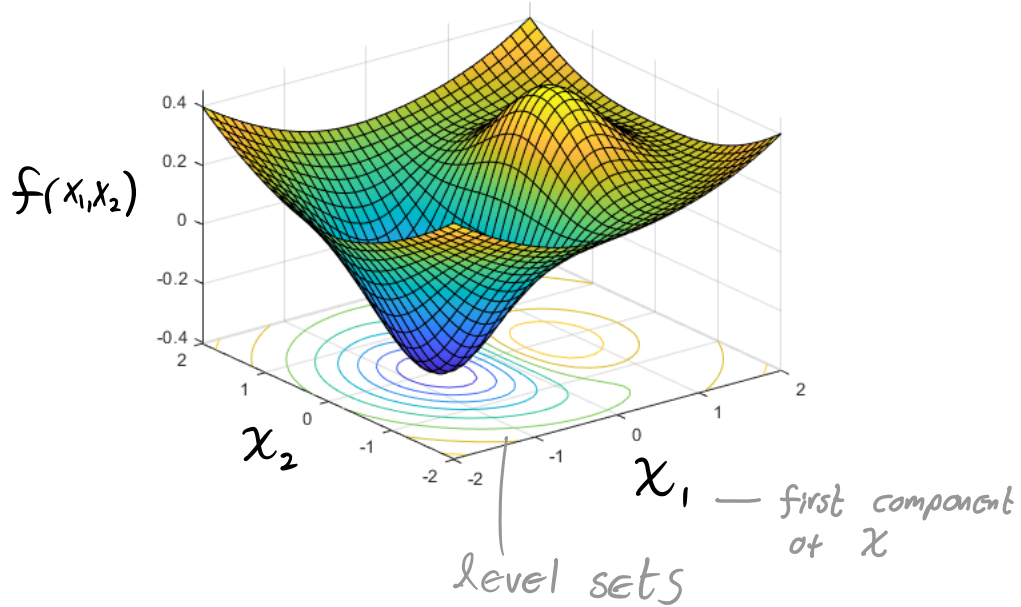
Gradient descent: Take successive steps "downhill"

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

iteration
index

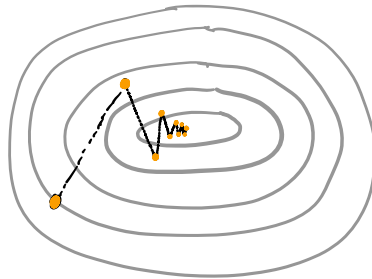
step size,
learning rate

$-\nabla f$ points in direction
of steepest descent



Depiction of gradient descent

Top down view:



Recall: gradient is orthogonal to level sets

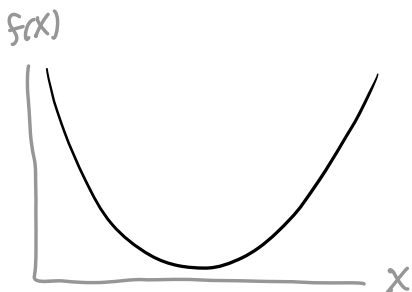
Example: Suppose $f: \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = x$.

What is sequence of points given by GD if starting from x^0 ?

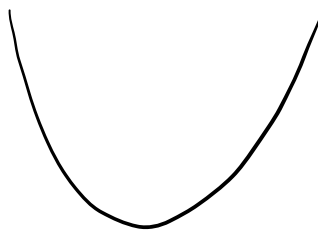
Example: $f: \mathbb{R} \rightarrow \mathbb{R}$
 $f(x) = \frac{1}{2} L x^2$. If GD is initialized at $x^{(0)}$,
what is value of $x^{(n)}$?

When does $x^{(n)}$ converge to minimizer
of f as $n \rightarrow \infty$?

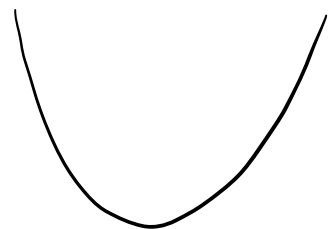
Picture:



Small learning
rate



medium learning
rate



high learning
rate

Example: $f: \mathbb{R} \rightarrow \mathbb{R}$
 $f(x) = |x|$. If GD is initialized at $x^{(0)}$,
Describe what GD will do.

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 $B \subset \{1, 2, \dots, n\}$

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

use as approximation
of $\nabla_{\theta} f(\theta)$

If you try to generate a minibatch by selecting a random subset of B data points uniformly, what practical challenges arise?

What considerations would affect the minibatch size you should use?

If the minibatch is chosen randomly,
on average, the gradient of a minibatch
is the full gradient

⇒ Stochastic gradient descent

Stochastic Gradient Descent

Want to solve $\min_x f(x)$

Instead of having access to $\nabla f(x)$,
suppose only have $G(x)$ w/ $E[G(x)] = \nabla f(x)$.

Write SGD as

$$x^{k+1} = x^k - \alpha_k G(x^k)$$

- on average, move in direction of steepest descent
- may move further from minimizer

Simple model: additive noise

$$G(x) = \nabla f(x) + w, \quad w \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$$

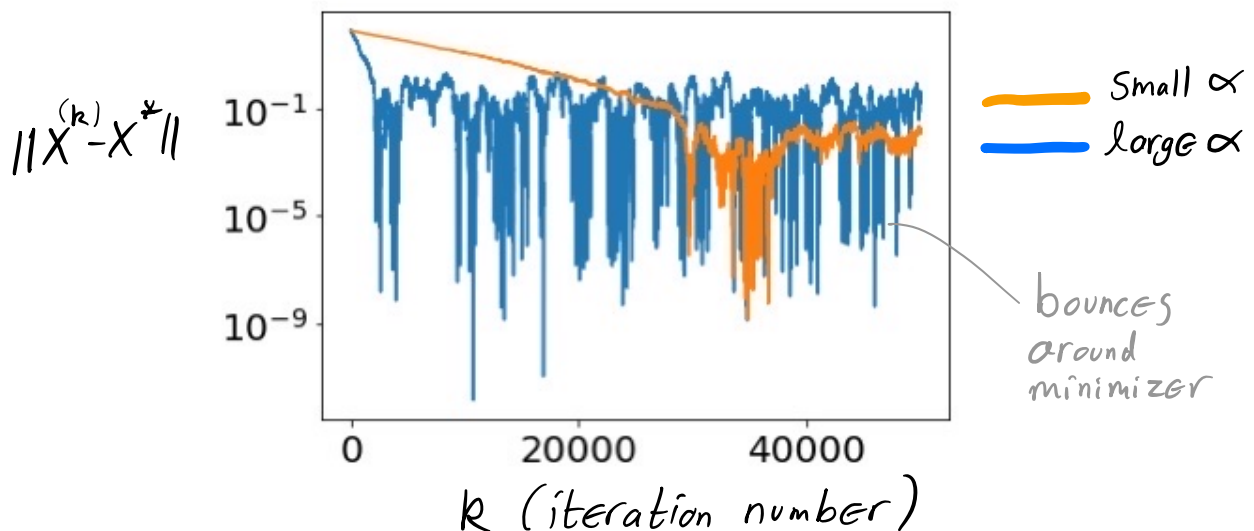
Use in ML: minibatches

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

$$G(\theta) = \frac{1}{|B|} \sum_{i \in B} \nabla_\theta \ell(\hat{y}_\theta(x_i), y_i) \quad \text{for random subset } B$$

Qualitatively,

with fixed step size α , $x^{(k)}$ will move close to x^* but will bounce around due to stochasticity



large $\alpha \Rightarrow$ fast initial convergence
large error

Small $\alpha \Rightarrow$ slow initial convergence
smaller error

Can formalize these observations w/ theory

How to choose step sizes/learning rates?

{ Run at a large value for α while
{ Shrink learning rate
{ Repeat

{ Have schedule of α_k decaying in k

In these cases can hope for convergence

Challenges w/ GD and SGD in Deep Learning

Non convexity and non smoothness

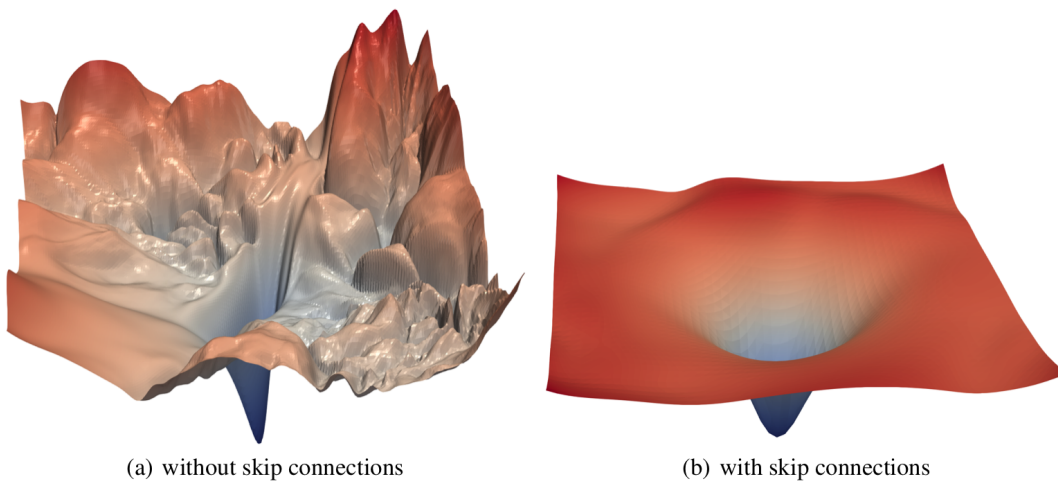


Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.

(Li et al. 2018)

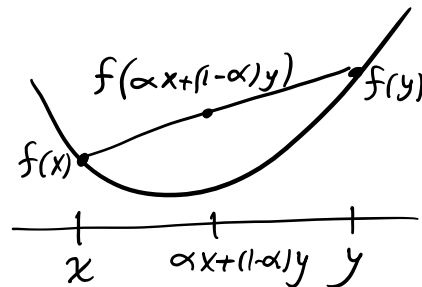
may be stuck in a local minimum,
so may want to temporarily increase
learning rate to get unstuck.

Convex Optimization

We say $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is **convex** if

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

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



"always curves up"

Examples: $f: \mathbb{R} \rightarrow \mathbb{R}$
 $f(x) = x^2$

is or is not convex

Fix a $c \in \mathbb{R}$.

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

$f(x) = cx^2$

is or is not convex

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

$$f(x) = x$$

is or is not CONVEX

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

$$f(x) = |x|$$

is or is not CONVEX

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

$$f(x) = \|x\|^2 = x_1^2 + x_2^2$$

is or is not CONVEX

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

$$f(x) = x_1^2$$

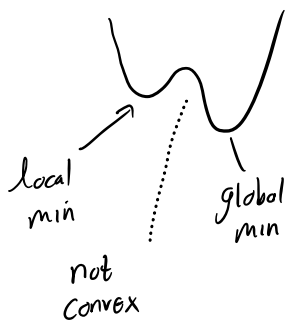
is or is not CONVEX

We will study the minimization of convex functions.

Does every convex function f have a minimal value?

$$\min_x f(x)$$

All local minima of convex functions are global minima.

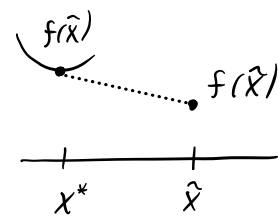


Suppose x^* is a local min of f .

If $x \approx x^*$ then $f(x) \geq f(x^*)$.

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

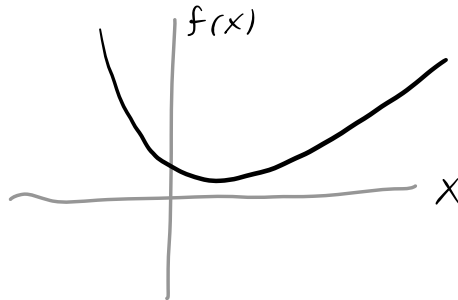
By convexity, $f(x)$ lies below dotted line between x^* and \tilde{x} . So x^* not a local min.



Convexity and Second derivatives

Functions of one variable

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



Functions of multiple variables

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$

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

Hessian matrix

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

H is positive definite if all eigenvalues are positive

H is positive semidefinite if all eigenvalues are nonnegative

Theorem: H is positive semidefinite if and only if

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

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 Λ is diagonal

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

$$H = \begin{pmatrix} | & | & \dots & | \\ v_1 & v_2 & \dots & v_n \\ | & | & \dots & | \end{pmatrix} \begin{pmatrix} \lambda_1 & & & 0 \\ & \ddots & & \\ 0 & & \ddots & \\ & & & \lambda_n \end{pmatrix} \begin{pmatrix} - v_1^t - \\ - v_2^t - \\ \vdots \\ - v_n^t - \end{pmatrix}$$

$U \cdot \Lambda \cdot U^t$

/

Columns are
unit length eigenvectors
that are orthogonal to
each other

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

Proof of Theorem: • PSD $\Rightarrow z^t H z \geq 0$ for all z

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

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

$$\geq 0$$

• $z^t H z \geq 0$ for all $z \Rightarrow H$ is PSD

Suppose H is not PSD. At least one eigenvalue is negative.

Suppose U_i is eigenvector w/ $\lambda_i < 0$. 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.

How fast does gradient descent converge?

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

Suppose $X^{(i)} \rightarrow X^*$ as $i \rightarrow \infty$.

How long do you need to wait to get a certain accuracy ϵ ?

Can gain understanding in some CONVEX cases.

Convergence of GD for quadratic functions

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

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

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

condition number of Q

Consider GD w/ fixed step size α

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

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

Analytically show that this is the solution to the problem

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

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

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

"first-order convergence"

Error decays exponentially

To get error ϵ , need $O(\log(\epsilon^{-1}))$ iterations

Proof: Note $\nabla f(x) = Qx - b$.

The global minimizer solves $Qx^* = b \Rightarrow x^* = Q^{-1}b$

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

So,

$$\|x^{k+1} - x^*\| \leq \underbrace{\|I - \alpha Q\|}_{\max(\alpha M - 1, 1 - \alpha m)} \|x^k - x^*\|$$

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

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

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

$$\Rightarrow \|X^k - X^*\| \leq \left(\frac{1 - 1/k}{1 + 1/k} \right)^k \|X^0 - X^*\| \quad \blacksquare$$

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"?

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

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

where x^* is a minimizer of f .

- Error decays slowly
- To get error ϵ from optimal value,
need $O(\epsilon^{-1})$ iterations

Summary :

- Too large learning rate can lead to divergence
- In convex case, to get convergence α 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

