Local, Unconstrained Function Optimization

COMPSCI 371D — Machine Learning

< 🗇 > < 🖻 > < 🖻

Outline

- 1 Motivation and Scope
- 2 First Order Methods
- Gradient, Hessian, and Convexity
- Gradient Descent
- 5 Descent Rate Selection Methods
- 6 Termination
- Is Gradient Descent a Good Strategy?
- 8 Stochastic Gradient Descent

Motivation and Scope

- Most estimation problems are solved by optimization
- Machine learning:
 - Parametric predictor: $h(\mathbf{x} ; \mathbf{v}) : \mathbb{R}^d \times \mathbb{R}^m \to Y$
 - Training set $T = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ and *loss* = $\ell(y_n, y)$
 - Risk: $L_T(\mathbf{v}) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, h(\mathbf{x}_n; \mathbf{v})) : \mathbb{R}^m \to \mathbb{R}$
 - Training: $\hat{\mathbf{v}} \in \arg\min_{\mathbf{v} \in \mathbb{R}^m} L_T(\mathbf{v})$
- "Solving" the system of equations e(z) = 0 can be viewed as

 $\hat{\boldsymbol{z}} = \in \arg\min_{\boldsymbol{z}} \|\boldsymbol{e}(\boldsymbol{z})\|$

Only Local Minimization

 $\hat{\mathbf{z}} = \arg\min_{\mathbf{z}\in\mathbf{?}} f(\mathbf{z})$

- All we know about *f* is a "black box" (think Python function)
- For many problems, f has many local minima
- Start somewhere (z₀), and take steps "down"
 f(z_{k+1}) < f(z_k)
- When we get stuck at a local minimum, we declare success
- · We would like global minima, but all we get is local ones
- For some problems, f has a unique minimum...
- ... or at least a single connected set of minima

Gradient

$$abla f(\mathbf{z}) = rac{\partial f}{\partial \mathbf{z}} = \left[egin{array}{c} rac{\partial f}{\partial z_1} \ dots \ rac{\partial f}{\partial z_m} \end{array}
ight]$$

 $\mathbf{z} \in \mathbb{R}^m$ with *m* possibly very large

- If ∇f(z) exists everywhere, the condition ∇f(z) = 0
 is necessary and sufficient for a stationary point (max, min, or saddle)
- Warning: only *necessary* for a minimum!
- Reduces to first derivative when $f : \mathbb{R} \to \mathbb{R}$

First Order Taylor Expansion

 $f(\mathbf{z}) \approx g_1(\mathbf{z}) = f(\mathbf{z}_0) + [\nabla f(\mathbf{z}_0)]^T(\mathbf{z} - \mathbf{z}_0)$

approximates $f(\mathbf{z})$ near \mathbf{z}_0 with a (hyper)plane through \mathbf{z}_0



 $\nabla f(\mathbf{z}_0)$ points to direction of steepest *increase* of *f* at \mathbf{z}_0

- If we want to find z₁ where f(z₁) < f(z₀), going along
 −∇f(z₀) seems promising
- This is the general idea of gradient descent

Hessian



• Symmetric matrix because of Schwarz's theorem:

$$\frac{\partial^2 f}{\partial z_i \partial z_j} = \frac{\partial^2 f}{\partial z_j \partial z_i}$$

- Eigenvalues are real because of symmetry
- Reduces to $\frac{d^2f}{dz^2}$ for $f : \mathbb{R} \to \mathbb{R}$

Convexity



- Strongly convex *everywhere*: For all \mathbf{z}, \mathbf{z}' in the (open) domain of f and for all $u \in (0, 1)$ $f(u\mathbf{z} + (1 - u)\mathbf{z}') < uf(\mathbf{z}) + (1 - u)f(\mathbf{z}')$
- Weak convexity: Replace "<" with "≤"

< A >

Convexity and Hessian

- Things become operational for twice-differentiable functions
- The function f(z) is strongly convex everywhere iff H(z) > 0 for all z
- " \succ " means *positive definite*: $\mathbf{v}^T H(\mathbf{z})\mathbf{v} > 0$ for all $\mathbf{v} \in \mathbb{R}^m$
- Above is *definition* of $H(\mathbf{z}) \succ 0$
- To check computationally: All eigenvalues are positive
- $H(\mathbf{z}) \succ 0$ reduces to $\frac{d^2f}{dz^2} > 0$ for $f : \mathbb{R} \to \mathbb{R}$

4 周 5 4 3 5 4 3 5 5

Local Convexity

- Definition: f is (strongly or weakly) convex at z₀ if it is (strongly or weakly) convex everywhere in some open neighborhood of z₀
- For *f* twice differentiable with continuous Hessian everywhere
 - *H*(**z**₀) ≻ 0 is *sufficient* (not necessary) for strong convexity at **z**₀
 - *H*(**z**₀) ≥ 0 is *necessary* (not sufficient) for weak convexity at **z**₀
- Examples:
 - $f(z) = z^2$ is strongly convex at $z_0 = 0$ and $H_f(0) = 2$
 - $f(z) = z^4$ is strongly convex at $z_0 = 0$ and $H_f(0) = 0$
 - $f(z) = z^3$ has a saddle at $z_0 = 0$ and $H_f(0) = 0$ (every neighborhood of $z_0 = 0$ has points (any z < 0) where $H_f(z) = 6z < 0$ so that $H_f(z) \prec 0$)

Some Uses of Convexity

- If ∇f(z₀) = 0 and f is (strongly or weakly) convex at z₀ then z₀ is a (strong or weak) minimum (as opposed to a maximum or a saddle)
- If *f* is globally convex then the value of the minimum is unique and the points where the minimum is achieved form a convex set
- Faster optimization methods (Newton) can be used when $f : \mathbb{R}^m \to \mathbb{R}$ is convex and *m* is not too large

A Template for Local Minimization

• Unconstrained minimization template:

• For some methods (Newton) the step

$$\mathbf{s}_k = \mathbf{z}_{k+1} - \mathbf{z}_k = \alpha_k \mathbf{p}_k$$

is the result of a single computation

• The *step size* is $\|\alpha_k \mathbf{p}_k\|$

Design Decisions

(**z**₀ given) k = 0while **z**_k is not a minimum compute step direction **p**_k compute descent rate $\alpha_k > 0$ **z**_{k+1} = **z**_k + α_k **p**_k k = k + 1

end

- In what direction to proceed (**p**_k)
- How long a step to take in that direction (α_k)
- When to stop ("while z_k is not a minimum")
- Different decisions lead to different methods

Gradient Descent

- In what direction to proceed: $\mathbf{p}_k = -\nabla f(\mathbf{z}_k)$
- "Gradient descent"
- Problem reduces to one dimension:
 h(α) = f(z_k + αp_k) with α_k > 0
- $\alpha = \mathbf{0} \Leftrightarrow \mathbf{z} = \mathbf{z}_k$
- Find $\alpha = \alpha_k > 0$ such that $f(\mathbf{z}_k + \alpha_k \mathbf{p}_k) < f(\mathbf{z}_k)$
- How to find α_k ?

-

Descent Rate

- Simplest idea: $\alpha_k = \alpha$ (fixed)
 - Step size $\| \alpha \nabla f(\mathbf{z}_k) \|$ decreases because $\nabla(\mathbf{z}_k) \to \mathbf{0}$
 - Small α leads to slow progress
 - Large α can miss minima



- Scheduling α :
 - Start with α relatively large (say $\alpha = 10^{-3}$)
 - Decrease α over time
 - Determine decrease rate by trial and error (good asymptotic guarantees with $\alpha_k \propto 1/(k+1)$).

Momentum

• Sometimes **z**_k meanders around on shallow plateaus



$$f(\mathbf{z}_k)$$
 versus k

- α is too small, direction is still promising
- Add momentum

$$\begin{aligned} \mathbf{v}_0 &= \mathbf{0} \\ \mathbf{v}_{k+1} &= \mu_k \mathbf{v}_k - \alpha_k \nabla f(\mathbf{z}_k) \\ \mathbf{z}_{k+1} &= \mathbf{z}_k + \mathbf{v}_{k+1} \end{aligned} \qquad (\mathbf{0} \leq \mu_k < \mathbf{1}) \end{aligned}$$

Line Search

• Find a local minimum in the search direction **p**_k

 $h(\alpha) = f(\mathbf{z}_k + \alpha \mathbf{p}_k)$, a one-dimensional problem

- Bracketing triple:
- a < b < c, $h(a) \ge h(b)$, $h(b) \le h(c)$
- Contains a (local) minimum!
- Split the bigger of [a, b] and [b, c] in half with a point u
- Find a new, narrower bracketing triple involving *u* and two out of *a*, *b*, *c*
- Stop when the bracket is narrow enough (say, 10⁻⁶)
- Pinned down a minimum to within 10⁻⁶

-

Phase 1: Find a Bracketing Triple



Phase 2: Shrink the Bracketing Triple



if
$$b - a > c - b$$

 $u = (a + b)/2$
if $h(u) > h(b)$
 $(a, b, c) = (u, b, c)$
otherwise
 $(a, b, c) = (a, u, b)$
end
otherwise
 $u = (b + c)/2$
if $h(u) > h(b)$
 $(a, b, c) = (a, b, u)$
otherwise
 $(a, b, c) = (b, u, c)$
end
end

Ξ.

< ロ > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > <

Termination

- Are we still making "significant progress"?
- Check $f(\mathbf{z}_{k-1}) f(\mathbf{z}_k)$? (We want this to be strictly positive)
- Check $\|\mathbf{z}_{k-1} \mathbf{z}_k\|$? (We want this to be large enough)
- Second is more stringent close the the minimum because ∇f(z) ≈ 0

• Stop when
$$\|\mathbf{z}_{k-1} - \mathbf{z}_k\| < \delta$$

Is Gradient Descent a Good Strategy?

- "We are going in the direction of fastest descent"
- "We choose an optimal descent rate by line search"
- "Must be good, no?" Not so fast!
- An example for which we know the answer:

$$f(\mathbf{z}) = \mathbf{c} + \mathbf{a}^T \mathbf{z} + \frac{1}{2} \mathbf{z}^T Q \mathbf{z}$$

 $Q \succeq 0$ (convex paraboloid)

All smooth functions look like this close enough to z*



isocontours

くぼう くきり くきり

Skating to a Minimum



- Many 90-degree turns slow down convergence
- There are methods that take fewer iterations, but each iteration takes more time and space
- We will stick to gradient descent
- See appendices in the notes for more efficient methods for problems in low-dimensional spaces

Stochastic Gradient Descent

• A special case of gradient descent, SGD works for *averages* of many terms (*N* very large):

$$f(\mathbf{z}) = \frac{1}{N} \sum_{n=1}^{N} \phi_n(\mathbf{z})$$

- Computing $\nabla f(\mathbf{z}_k)$ is too expensive
- Partition B = {1,..., N} into J random mini-batches B_j each of about equal size

$$f(\mathbf{z}) \approx f_j(\mathbf{z}) = rac{1}{|B_j|} \sum_{n \in B_j} \phi_n(\mathbf{z}) \quad \Rightarrow \quad \nabla f(\mathbf{z}) \approx \nabla f_j(\mathbf{z}) \;.$$

Mini-batch gradients are correct on average

SGD and Mini-Batch Size

- SGD iteration: $\mathbf{z}_{k+1} = \mathbf{z}_k \alpha_k \nabla f_j(\mathbf{z}_k)$
- Mini-batch gradients are correct on average
- One cycle through all the mini-batches is an epoch
- Repeatedly cycle through all the data (Scramble data before each epoch)
- *Asymptotic* convergence can be proven with suitable descent-rate schedule
- Small batches \Rightarrow low storage but high gradient variance
- Make batches as big as will fit in memory for minimal variance
- In deep learning, memory is GPU memory

-