

Optimization methods

Unconstrained convex optimization

Unconstrained convex optimization problem

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

(f is the convex objective function; feasible region is \mathbb{R}^d).

Gradient descent for differentiable objectives

- ▶ Start with some initial $\mathbf{w}^{(1)} \in \mathbb{R}^d$.
- ▶ For $t = 1, 2, \dots$ until some stopping condition is satisfied.
 - ▶ Compute gradient of f at $\mathbf{w}^{(t)}$:

$$\boldsymbol{\lambda}^{(t)} := \nabla f(\mathbf{w}^{(t)}).$$

- ▶ Update:

$$\mathbf{w}^{(t+1)} := \mathbf{w}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$

($\eta_t > 0$ are step sizes.)

1 / 25

2 / 25

Two issues

(Projected) (sub)gradient descent

1. Objective function not necessarily differentiable everywhere.
2. Feasible region not necessarily \mathbb{R}^d .

4 / 25

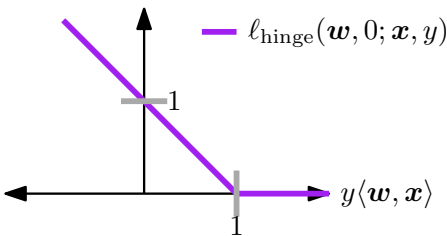
Non-differentiability

Non-differentiable convex objectives

Some convex functions f are not differentiable everywhere;
gradient descent not even well-specified for these problems.

Example: hinge loss

$$f(\mathbf{w}) = \ell_{\text{hinge}}(\mathbf{w}, 0; \mathbf{x}, y) = \left[1 - y\langle \mathbf{w}, \mathbf{x} \rangle \right]_+.$$



Not differentiable at $\mathbf{w} \in \mathbb{R}^d$ where $y\langle \mathbf{w}, \mathbf{x} \rangle = 1$.

5 / 25

Dealing with non-differentiability

Local optimization

- ▶ Locally change $\mathbf{x} \rightarrow \mathbf{x} + \delta$.
- ▶ Hopefully improve objective value $f(\mathbf{x}) \rightarrow f(\mathbf{x} + \delta)$.

Subgradients

We say $\boldsymbol{\lambda} \in \mathbb{R}^d$ is a **subgradient** of a function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ at $\mathbf{w}_0 \in \mathbb{R}^d$ if

$$f(\mathbf{w}) \geq f(\mathbf{w}_0) + \langle \boldsymbol{\lambda}, \mathbf{w} - \mathbf{w}_0 \rangle \quad \text{for all } \mathbf{w} \in \mathbb{R}^d,$$

i.e., $\boldsymbol{\lambda}$ specifies an **affine lower bound** on the function.

Although not every function f is differentiable everywhere,
every **convex function** f has **subgradients everywhere**.

(Some technical conditions apply.)

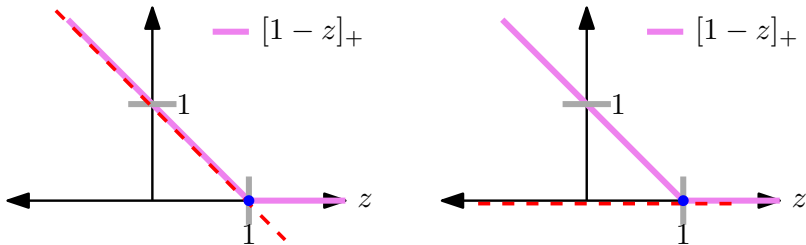
Call the entire set the **subdifferential of f at \mathbf{w}_0** ; denote it by

$$\partial f(\mathbf{w}_0) := \{ \boldsymbol{\lambda} \in \mathbb{R}^d : \boldsymbol{\lambda} \text{ is a subgradient of } f \text{ at } \mathbf{w}_0 \}.$$

6 / 25

Example: max of two affine functions

Consider one-dimensional function $f(z) := [1 - z]_+ = \max\{0, 1 - z\}$.



Two subgradients of f at $z = 1$: -1 and 0 .

$$\begin{aligned} f(z) &\geq f(1) + (-1) \cdot (z - 1) = 1 - z; \\ f(z) &\geq f(1) + (0) \cdot (z - 1) = 0. \end{aligned}$$

In fact, every $\lambda \in [-1, 0]$ satisfies

$$f(z) \geq f(1) + \lambda \cdot (z - 1).$$

We have $\partial f(1) = [-1, 0]$.

7 / 25

Subgradient descent

Subgradient descent for general convex objectives

- ▶ Start with some initial $\mathbf{w}^{(1)} \in \mathbb{R}^d$.
- ▶ For $t = 1, 2, \dots$ until some stopping condition is satisfied.
 - ▶ Compute *any* subgradient $\boldsymbol{\lambda}^{(t)} \in \partial f(\mathbf{w}^{(t)})$.
 - ▶ Update:

$$\mathbf{w}^{(t+1)} := \mathbf{w}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$

($\eta_t > 0$ are step sizes.)

8 / 25

Example: soft-margin (homogeneous) SVM

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}), \quad f(\mathbf{w}) := \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+$$

Question: How do we compute a subgradient \mathbf{g} of f at a given point $\mathbf{w}_0 \in \mathbb{R}^d$?

$$\mathbf{g} = \lambda \mathbf{w}_0 + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \begin{cases} \mathbf{0} & \text{if } 1 - y\langle \mathbf{w}_0, \mathbf{x} \rangle < 0; \\ -y\mathbf{x} & \text{if } 1 - y\langle \mathbf{w}_0, \mathbf{x} \rangle > 0; \\ \text{any conv. comb. of above} & \text{if } 1 - y\langle \mathbf{w}_0, \mathbf{x} \rangle = 0. \end{cases}$$

With “continuous” data, exact equalities like $1 - y\langle \mathbf{w}_0, \mathbf{x} \rangle = 0$ never come up.

9 / 25

Example: soft-margin (homogeneous) SVM

Subgradient descent algorithm for soft-margin SVM:

- ▶ Start with some initial $\mathbf{w}^{(1)} \in \mathbb{R}^d$.
- ▶ For $t = 1, 2, \dots$ until some stopping condition is satisfied.

$$\begin{aligned} \mathbf{w}^{(t+1)} &:= \mathbf{w}^{(t)} - \eta_t \left(\lambda \mathbf{w}^{(t)} + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \begin{cases} \mathbf{0} & \text{if } 1 - y\langle \mathbf{w}^{(t)}, \mathbf{x} \rangle < 0; \\ -y\mathbf{x} & \text{if } 1 - y\langle \mathbf{w}^{(t)}, \mathbf{x} \rangle \geq 0 \end{cases} \right) \\ &= (1 - \lambda \eta_t) \mathbf{w}^{(t)} + \eta_t \frac{1}{|S|} \sum_{\substack{(\mathbf{x}, y) \in S: \\ y\langle \mathbf{w}^{(t)}, \mathbf{x} \rangle \leq 1}} y\mathbf{x}. \end{aligned}$$

Note effect of regularization term $\frac{\lambda}{2} \|\mathbf{w}\|_2^2$ (whenever $\eta_t < 1/\lambda$):
Shrink $\mathbf{w}^{(t)}$ by a factor $1 - \lambda \eta_t$ before updating with subgradient of loss term.
Helps prevent $\mathbf{w}^{(t)}$ from becoming too long.

10 / 25

Constrained convex optimization

$$\begin{array}{ll} \min_{\mathbf{w} \in \mathbb{R}^d} & f(\mathbf{w}) \\ \text{s.t.} & \mathbf{w} \in A \end{array}$$

for convex function f and convex feasible set A (possibly a strict subset of \mathbb{R}^d).

Projected subgradient descent

- ▶ Start with some initial $\mathbf{w}^{(1)} \in A$.
- ▶ For $t = 1, 2, \dots$ until some stopping condition is satisfied.
 - ▶ Compute any subgradient $\boldsymbol{\lambda}^{(t)} \in \partial f(\mathbf{w}^{(t)})$.
 - ▶ Update:
$$\mathbf{w}^{(t+0.5)} := \mathbf{w}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$
 - ▶ Project to feasible region A :

$$\mathbf{w}^{(t+1)} := \text{Proj}_A(\mathbf{w}^{(t+0.5)}).$$

Projection onto convex set A : $\text{Proj}_A(\mathbf{w}) := \arg \min_{\mathbf{w}' \in A} \|\mathbf{w} - \mathbf{w}'\|_2^2$
(another convex optimization problem, but hopefully simpler objective!).

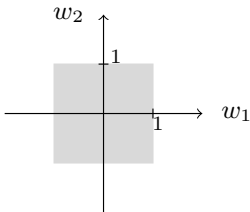
11 / 25

Example: box constraints

Suppose you want to prevent any feature from having a high weight:

- ▶ Constraints:
$$|w_i| \leq 1 \quad \text{for all } i = 1, 2, \dots, d.$$
- ▶ Feasible region:

$$A := [-1, +1]^d$$



- ▶ $\text{Proj}_A(\mathbf{w})$:
 - ▶ For each $i = 1, 2, \dots, d$:
 - ▶ If $w_i \in [-1, +1]$, leave w_i alone.
 - ▶ If $w_i > +1$, set $w_i := 1$.
 - ▶ If $w_i < -1$, set $w_i := -1$.

12 / 25

Convergence of (projected) subgradient descent

In general, subgradient descent (with possibly non-differentiable convex objectives) **can converge relatively slowly**.

(Can help to return the *average* of the iterates $\frac{1}{T} \sum_{t=1}^T \mathbf{w}^{(t)}$.)

Optimality condition for general convex optimization problems also considerably more involved (e.g., more than just “ $\nabla f_0(\mathbf{w}) = \mathbf{0}$ ”).

But for machine learning purposes, **we can use alternative criteria for stopping** (e.g., **hold-out error rate**).

Stochastic gradient method

In machine learning, we typically have objectives of the following form:

$$F(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \phi_i(\mathbf{w})$$

for convex functions $\phi_i: \mathbb{R}^d \rightarrow \mathbb{R}$ for $i = 1, 2, \dots, n$.

Example: ϕ_i is loss on i -th training example.

Computational cost of computing $\nabla F(\mathbf{w})$: linear in n . (Typically $O(nd)$).

Key idea

Let I be a uniform random variable in $\{1, 2, \dots, n\}$.

Then for any \mathbf{w} ,

$$F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \phi_i(\mathbf{w}) = \mathbb{E}[\phi_I(\mathbf{w})]$$

and by linearity,

$$\nabla F(\mathbf{w}) = \nabla \left\{ \frac{1}{n} \sum_{i=1}^n \phi_i(\mathbf{w}) \right\} = \frac{1}{n} \sum_{i=1}^n \nabla \phi_i(\mathbf{w}) = \mathbb{E}[\nabla \phi_I(\mathbf{w})].$$

Upshot: $\nabla \phi_I(\mathbf{w})$ is an *unbiased* estimate of gradient of F at \mathbf{w} .

Computational cost of computing $\nabla \phi_I(\mathbf{w})$: **independent of n** . (E.g., $O(d)$).

Stochastic gradient method

Consider convex functions $\phi_i: \mathbb{R}^d \rightarrow \mathbb{R}, i = 1, 2, \dots, n$;

$$\min_{\mathbf{w} \in A} F(\mathbf{w}), \quad F(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \phi_i(\mathbf{w}).$$

Stochastic (projected) (sub)gradient “descent” method

- ▶ Start with some initial $\mathbf{w}^{(1)} \in A$.
- ▶ For $t = 1, 2, \dots$ until some stopping condition is satisfied.
 - ▶ Pick I_t uniformly at random from $\{1, 2, \dots, n\}$.
 - ▶ Compute *any* subgradient $\boldsymbol{\lambda}^{(t)} \in \partial \phi_{I_t}(\mathbf{w}^{(t)})$.
 - ▶ Update:

$$\mathbf{w}^{(t+0.5)} := \mathbf{w}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$

- ▶ Project to feasible region A :

$$\mathbf{w}^{(t+1)} := \text{Proj}_A(\mathbf{w}^{(t+0.5)}).$$

17 / 25

Example: logistic regression

Maximum likelihood estimation for logistic regression:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)).$$

Stochastic gradient method for logistic regression:

- ▶ Start with some initial $\mathbf{w}^{(1)} \in \mathbb{R}^d$.
- ▶ For $t = 1, 2, \dots$ until some stopping condition is satisfied.
 - ▶ Pick I_t uniformly at random from $\{1, 2, \dots, n\}$;

$$\mathbf{w}^{(t+1)} := \mathbf{w}^{(t)} + \eta_t (1 - P_{\mathbf{w}^{(t)}}(Y=y_{I_t} \mid \mathbf{X}=\mathbf{x}_{I_t})) y_{I_t} \mathbf{x}_{I_t}.$$

18 / 25

Example: soft-margin (homogeneous) SVMs

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \left(\frac{\lambda}{2} \|\mathbf{w}\|_2^2 + [1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle]_+ \right)$$

Stochastic gradient method for soft-margin SVMs:

- ▶ Start with some initial $\mathbf{w}^{(1)} \in \mathbb{R}^d$.
- ▶ For $t = 1, 2, \dots$ until some stopping condition is satisfied.
 - ▶ Pick I_t uniformly at random from $\{1, 2, \dots, n\}$;

$$\begin{aligned} \mathbf{w}^{(t+1)} &:= \mathbf{w}^{(t)} - \eta_t \left(\lambda \mathbf{w}^{(t)} + \begin{cases} \mathbf{0} & \text{if } 1 - y_{I_t} \langle \mathbf{w}^{(t)}, \mathbf{x}_{I_t} \rangle < 0; \\ -y_{I_t} \mathbf{x}_{I_t} & \text{if } 1 - y_{I_t} \langle \mathbf{w}^{(t)}, \mathbf{x}_{I_t} \rangle \geq 0 \end{cases} \right) \\ &= (1 - \lambda \eta_t) \mathbf{w}^{(t)} + \begin{cases} \mathbf{0} & \text{if } y_{I_t} \langle \mathbf{w}^{(t)}, \mathbf{x}_{I_t} \rangle > 1; \\ y_{I_t} \mathbf{x}_{I_t} & \text{if } y_{I_t} \langle \mathbf{w}^{(t)}, \mathbf{x}_{I_t} \rangle \leq 1. \end{cases} \end{aligned}$$

19 / 25

Why should this work?

Consider differentiable convex functions $\phi_i: \mathbb{R}^d \rightarrow \mathbb{R}, i = 1, 2, \dots, n$:

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}), \quad F(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \phi_i(\mathbf{w}).$$

Then

$$\boldsymbol{\lambda}^{(t)} = \nabla \phi_{I_t}(\mathbf{w}^{(t)}) = \nabla F(\mathbf{w}^{(t)}) + \underbrace{\left(\nabla \phi_{I_t}(\mathbf{w}^{(t)}) - \mathbb{E}[\nabla \phi_{I_t}(\mathbf{w}^{(t)})] \right)}_{\text{mean-zero “noise”}}.$$

“Noise” is “averaged away” when step sizes η_t are small (e.g., $\eta_t \sim 1/\sqrt{t}$).

20 / 25

Convergence theory

Consider convex functions $\phi_i: \mathbb{R}^d \rightarrow \mathbb{R}$, $i = 1, 2, \dots, n$;

$$\min_{\mathbf{w} \in A} F(\mathbf{w}), \quad F(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \phi_i(\mathbf{w}).$$

Run stochastic gradient method for T iterations, get iterates $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{T+1} \in A$.

Under some suitable conditions (e.g., on step sizes η_t),

$$\mathbb{E}[F(\mathbf{w}_{T+1})] \leq \min_{\mathbf{w} \in A} F(\mathbf{w}) + O\left(\frac{\log T}{\sqrt{T}}\right).$$

Let $\bar{\mathbf{w}} := \frac{1}{T+1} \sum_{t=1}^{T+1} \mathbf{w}_t$ (average of the iterates). Under same conditions,

$$\mathbb{E}[F(\bar{\mathbf{w}})] \leq \min_{\mathbf{w} \in A} F(\mathbf{w}) + O\left(\frac{1}{\sqrt{T}}\right).$$

Often really helps in practice!

21 / 25

Efficient implementation trick

Special case:

- ▶ $\phi_i(\mathbf{w}) = \lambda \|\mathbf{w}\|_2^2 / 2 + \ell(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$ for some convex function $\ell: \mathbb{R} \rightarrow \mathbb{R}$.
- ▶ \mathbf{x}_i has few non-zero entries (i.e., \mathbf{x}_i is “sparse”).

Gradient of ϕ_i at \mathbf{w} :

$$\nabla \phi_i(\mathbf{w}) = \lambda \mathbf{w} + \ell'(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) y_i \mathbf{x}_i.$$

Update ($I_t = i$):

$$\mathbf{w}_{t+1} := (1 - \lambda \eta_t) \mathbf{w}_t - \eta_t \ell'(y_i \langle \mathbf{w}_t, \mathbf{x}_i \rangle) y_i \mathbf{x}_i.$$

- ▶ Represent weight vector \mathbf{w}_t as $c_t \mathbf{v}_t$ for scalar $c_t \in \mathbb{R}$ and vector $\mathbf{v}_t \in \mathbb{R}^d$.
- ▶ Sparsity-respecting update:

$$\begin{aligned} g_t &:= \ell'(y_t c_t \langle \mathbf{v}_t, \mathbf{x}_i \rangle) \\ c_{t+1} &:= (1 - \lambda \eta_t) c_t \\ \mathbf{v}_{t+1} &:= \mathbf{v}_t - \frac{\eta_t g_t}{c_{t+1}} y_i \mathbf{x}_i \end{aligned}$$

22 / 25

Other practical tips

- ▶ Pick I_1, I_2, \dots, I_t u.a.r. **without replacement** from $\{1, 2, \dots, n\}$.
I.e., randomly shuffle order of training data, then process in that order.
If possible, re-shuffle after each pass through training data.
- ▶ For debugging purposes, periodically check overall objective value.
(Should generally be improving, though not necessarily at every step.)
- ▶ Use hold-out error rate as a stopping criterion.
(Stop when hold-out error rate does not improve after a while.)
- ▶ Use hold-out set to tune step sizes.

23 / 25

Other solvers

Many other algorithms for solving convex optimization problems

- ▶ *Newton-Raphson*: use Hessian to pick better descent directions.
- ▶ *Quasi-Newton methods* (e.g., conjugate gradient, “BFGS”, “L-BFGS”):
use efficient approximations of Hessians.
- ▶ Techniques for dealing with constraints:
 - ▶ *Barrier methods*: add penalties for constraint violations, slowly relax.
 - ▶ *Primal-dual methods*: start with dual-feasible point, iteratively improve until corresponding primal point is feasible.
 - ▶ ...
- ▶ ...

But remember: end goal in machine learning is *not* to minimize training error rate or training surrogate loss.

24 / 25

Key takeaways

1. Concept of subgradients, and generalizing gradient descent to non-differentiable objectives.
2. Subgradients for maximum of affine functions.
3. Extending gradient descent to constrained optimization problems.
4. Concept of stochastic gradients; stochastic gradient descent method; high-level idea of convergence theory.
5. Efficient implementation for sparse gradients.