

Convex optimization

1 / 26

Soft-margin SVMs

Soft-margin SVM optimization problem defined by training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$:

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

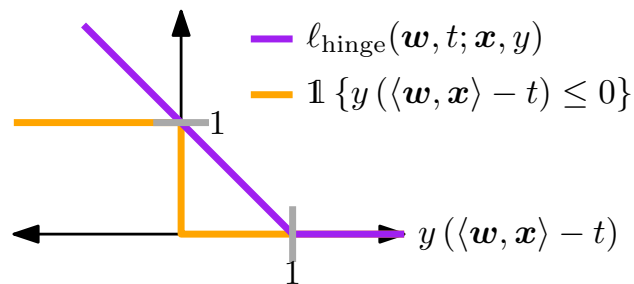
Compare with **Empirical Risk Minimization** (i.e., minimize training error rate):

$$\min_{\mathbf{w} \in \mathbb{R}^d, t \in \mathbb{R}} \quad \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle - t) \leq 0\}.$$

In both cases, i -th term in summation is function of $y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle - t)$.

2 / 26

Zero-one loss vs. hinge loss



Hinge loss: an upper-bound on **zero-one loss**.

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

Soft-margin SVM minimizes an upper-bound on the training error rate, plus a term that encourages large margins.

3 / 26

General form

Soft-margin SVM:

$$\min_{\mathbf{w} \in \mathbb{R}^d, t \in \mathbb{R}} \quad \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \ell_{\text{hinge}}(\mathbf{w}, t; \mathbf{x}_i, y_i).$$

Empirical risk minimization (i.e., minimize training error rate):

$$\min_{\mathbf{w} \in \mathbb{R}^d, t \in \mathbb{R}} \quad \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle - t) \leq 0\}.$$

Generic learning objective:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \quad R(\mathbf{w}) + \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{x}_i, y_i).$$

- **Regularization term:** encodes “learning bias” (e.g., prefer large margins).
- **Training loss term:** how poor is the “fit” to the training data.

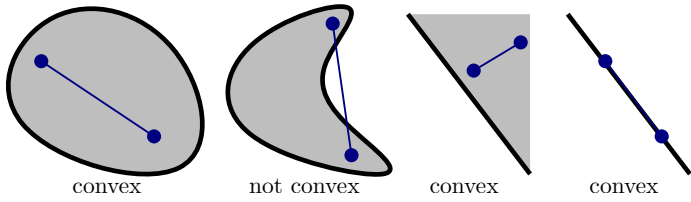
4 / 26

- ▶ Many different choices for regularization and loss.
 - ▶ $R(\mathbf{w}) = \lambda \|\mathbf{w}\|_2^2/2$: encourage large margins
 - ▶ $R(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$: encourage \mathbf{w} to be sparse
 - ▶ $R(\mathbf{w}) = \lambda \sum_{i=1}^{d-1} |w_{i+1} - w_i|$: useful for ordered features.
 - ▶ $\ell(\mathbf{w}; \mathbf{x}, y) = [\langle \mathbf{w}, \mathbf{x} \rangle - y]_+^2$
 - ▶ $\ell(\mathbf{w}, t; \mathbf{x}, y) = \ln(1 + \exp(-y\langle \mathbf{w}, \mathbf{x} \rangle))$ (logistic regression)
 - ▶ ...
 - ▶ Also used beyond classification (e.g., regression, parameter estimation, clustering)
- ▶ For classification problems, often want ℓ to be an upper-bound on zero-one loss—i.e., a **surrogate loss**.
- ▶ Trade-off parameter $\lambda > 0$: usually determine using cross validation.
- ▶ Computationally easier when overall objective function is **convex**: possible to efficiently find global minimizer in polynomial time.

Introduction to convexity

Convex sets

A set A is **convex** if, for every pair of points $\{\mathbf{x}, \mathbf{x}'\}$ in A , the line segment between \mathbf{x} and \mathbf{x}' is also contained in A .



Examples:

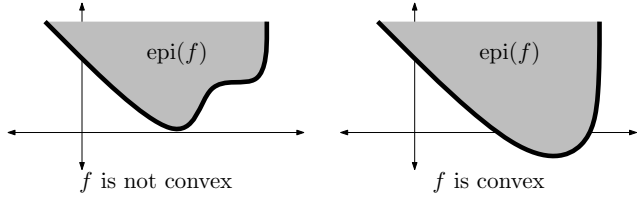
- ▶ All of \mathbb{R}^d .
- ▶ Empty set.
- ▶ Affine hyperplanes.
- ▶ Half-spaces: $\{\mathbf{a} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{x} \rangle - b \leq 0\}$.
- ▶ Intersections of convex sets.
- ▶ Convex hulls:
 $\text{conv}(S) := \{\sum_{i=1}^k \alpha_i \mathbf{x}_i : k \in \mathbb{N}, \mathbf{x}_i \in S, \alpha_i \geq 0, \sum_{i=1}^k \alpha_i = 1\}$.

Convex functions

For any function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, the **epigraph** of f , denoted $\text{epi}(f)$, is the set

$$\text{epi}(f) := \{(\mathbf{x}, b) \in \mathbb{R}^{d+1} : f(\mathbf{x}) \leq b\}.$$

A function f is **convex** if $\text{epi}(f)$ is a convex set in \mathbb{R}^{d+1} .



Examples:

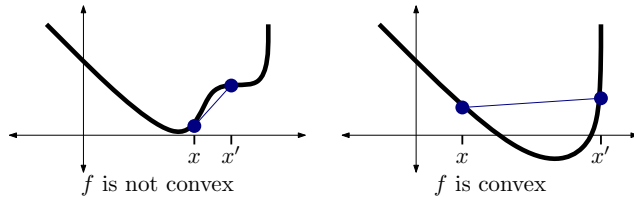
- ▶ $f(x) = c^x$ for any $c > 0$ (on \mathbb{R})
- ▶ $f(x) = |x|^c$ for any $c \geq 1$ (on \mathbb{R})
- ▶ $f(\mathbf{x}) = c$ for any constant $c \in \mathbb{R}$.
- ▶ $f(\mathbf{x}) = \langle \mathbf{a}, \mathbf{x} \rangle$ for any $\mathbf{a} \in \mathbb{R}^d$.
- ▶ $f(\mathbf{x}) = \|\mathbf{x}\|$ for any norm $\|\cdot\|$.
- ▶ $f(\mathbf{x}) = \langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle$ for symmetric positive semidefinite \mathbf{A} .

Jensen's inequality

Equivalent definition of convex functions

A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex iff for any $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ and $\alpha \in [0, 1]$,

$$f((1-\alpha)\mathbf{x} + \alpha\mathbf{x}') \leq (1-\alpha) \cdot f(\mathbf{x}) + \alpha \cdot f(\mathbf{x}').$$



Jensen's inequality

If $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex, then for any $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and $\alpha_1, \alpha_2, \dots, \alpha_n \in [0, 1]$ such that $\sum_{i=1}^n \alpha_i = 1$,

$$f\left(\sum_{i=1}^n \alpha_i \mathbf{x}_i\right) \leq \sum_{i=1}^n \alpha_i \cdot f(\mathbf{x}_i).$$

9 / 26

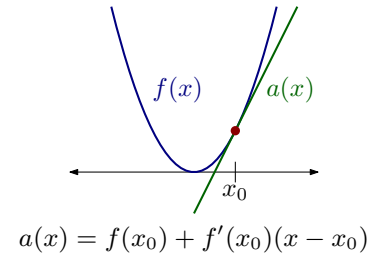
Convexity of differentiable functions

Differentiable functions

If $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is differentiable, then f is convex if and only if

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle$$

for all $\mathbf{x}, \mathbf{x}_0 \in \mathbb{R}^d$.



Twice-differentiable functions

If $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is twice-differentiable, then f is convex if and only if

$$\nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$$

for all $\mathbf{x} \in \mathbb{R}^d$ (i.e., the Hessian, or matrix of second-derivatives, is positive semidefinite for all \mathbf{x}).

10 / 26

Convex optimization problems

Optimization problems

A typical optimization problem (in standard form) is written as

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^d} \quad & f_0(\mathbf{x}) \\ \text{s.t.} \quad & f_i(\mathbf{x}) \leq 0 \quad \text{for all } i = 1, 2, \dots, n. \end{aligned}$$

- ▶ $f_0: \mathbb{R}^d \rightarrow \mathbb{R}$ is the **objective function**;
- ▶ $f_1, f_2, \dots, f_n: \mathbb{R}^d \rightarrow \mathbb{R}$ are the **constraint functions**;
- ▶ inequalities $f_i(\mathbf{x}) \leq 0$ are **constraints**;
- ▶ $A := \{\mathbf{x} \in \mathbb{R}^d : f_i(\mathbf{x}) \leq 0 \text{ for all } i = 1, 2, \dots, n\}$ is the **feasible region**.
- ▶ **The goal**: Find $\mathbf{x} \in A$ so that $f_0(\mathbf{x})$ is as small as possible.
- ▶ **(Optimal) value** of the optimization problem is the smallest such value of $f_0(\mathbf{x})$ achieved by a feasible point $\mathbf{x} \in A$.
- ▶ Point $\mathbf{x} \in A$ achieving the optimal value is a **(global) minimizer** of the problem.

12 / 26

Standard form of a **convex optimization problem**:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^d} \quad & f_0(\mathbf{x}) \\ \text{s.t.} \quad & f_i(\mathbf{x}) \leq 0 \quad \text{for all } i = 1, 2, \dots, n \end{aligned}$$

where $f_0, f_1, \dots, f_n: \mathbb{R}^d \rightarrow \mathbb{R}$ are *convex functions*.

Fact: the feasible set $A := \{\mathbf{x} \in \mathbb{R}^d : f_i(\mathbf{x}) \leq 0 \text{ for all } i = 1, 2, \dots, n\}$ is a convex set.

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, t \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad & \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle - t) \geq 1 - \xi_i \quad \text{for all } i = 1, 2, \dots, n, \\ & \xi_i \geq 0 \quad \text{for all } i = 1, 2, \dots, n. \end{aligned}$$

Bringing to standard form:

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, t \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad & \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \xi_i \quad (\text{sum of two convex functions, also convex}) \\ \text{s.t.} \quad & 1 - \xi_i - y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle - t) \leq 0 \quad \text{for all } i = 1, \dots, n \quad (\text{linear}) \\ & -\xi_i \leq 0 \quad \text{for all } i = 1, \dots, n \quad (\text{linear}) \end{aligned}$$

Consider an optimization problem (not necessarily convex):

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^d} \quad & f_0(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{x} \in A. \end{aligned}$$

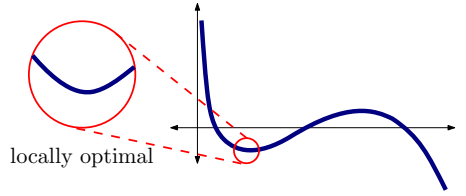
We say $\tilde{\mathbf{x}} \in A$ is a **local minimizer** if there is an open ball

$$U := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \tilde{\mathbf{x}}\|_2 < r\}$$

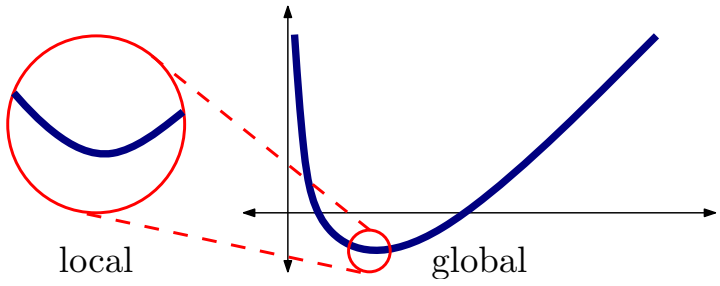
of positive radius $r > 0$ such that $\tilde{\mathbf{x}}$ is a global minimizer for

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^d} \quad & f_0(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{x} \in A \cap U. \end{aligned}$$

Nothing looks better than $\tilde{\mathbf{x}}$ in the immediate vicinity of $\tilde{\mathbf{x}}$.



If the optimization problem is **convex**, and $\tilde{\mathbf{x}} \in A$ is a **local minimizer**, then it is also a **global minimizer**.



Local optimization algorithms

Unconstrained convex optimization

Unconstrained convex optimization problem

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

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

Optimality condition for differentiable convex objectives

\mathbf{x}^* is a global minimizer if and only if $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

Unfortunately, can't always find closed-form solution to system of equations $\nabla f(\mathbf{x}) = \mathbf{0}$. \rightarrow Resort to iterative methods to find a solution.

18 / 26

Local optimization for convex objectives

Locally change $\mathbf{x} \rightarrow \mathbf{x} + \delta$.

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

How to pick δ ?

By convexity of f : $f(\mathbf{x} + \delta) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \delta \rangle$.

If $\langle \nabla f(\mathbf{x}), \delta \rangle \geq 0$, then

$$f(\mathbf{x} + \delta) \geq f(\mathbf{x}). \quad \text{Clearly a bad direction.}$$

Moral: to be useful, the change δ must satisfy

$$\langle \nabla f(\mathbf{x}), \delta \rangle < 0.$$

For example, $\delta := -\eta \nabla f(\mathbf{x})$ for some $\eta > 0$:

$$\langle \nabla f(\mathbf{x}), -\eta \nabla f(\mathbf{x}) \rangle = -\eta \|\nabla f(\mathbf{x})\|_2^2 < 0$$

as long as $\nabla f(\mathbf{x}) \neq \mathbf{0}$.

Gradient descent

Gradient descent for differentiable objectives

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

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

- ▶ Update:

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

Here, $\eta_1, \eta_2, \dots > 0$ are the **step sizes**. Common choices include:

1. Set $\eta_t := c$ for some constant $c > 0$.
2. Set $\eta_t := c/\sqrt{t}$ for some constant $c > 0$.
3. Set η_t using a line search procedure.

Backtracking line search

Goal: given $\mathbf{x} \in \mathbb{R}^d$ and $\boldsymbol{\lambda} = \nabla f(\mathbf{x}) \in \mathbb{R}^d$, find $\eta > 0$ so that $f(\mathbf{x} - \eta \boldsymbol{\lambda}) < f(\mathbf{x})$ by a reasonable amount.

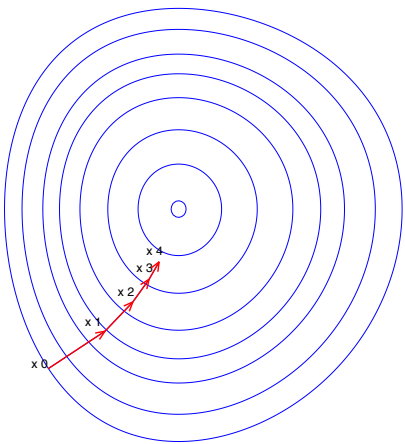
- ▶ Start with $\eta := 1$.
- ▶ While $f(\mathbf{x} - \eta \boldsymbol{\lambda}) > f(\mathbf{x}) - \frac{1}{2} \eta \|\boldsymbol{\lambda}\|_2^2$: Set $\eta := \frac{1}{2} \eta$.

Main idea: $f(\mathbf{x} - \eta \boldsymbol{\lambda}) \approx f(\mathbf{x}) - \eta \|\boldsymbol{\lambda}\|_2^2$ when η is small, so can optimistically hope to decrease value by about $\eta \|\boldsymbol{\lambda}\|_2^2$.

Settle for decreasing by $\frac{1}{2} \eta \|\boldsymbol{\lambda}\|_2^2$: upon termination of while-loop,

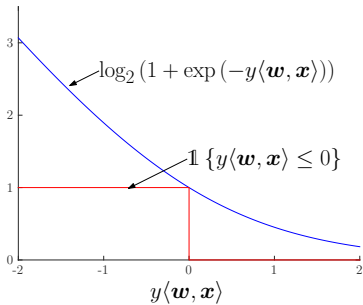
$$f(\mathbf{x} - \eta \boldsymbol{\lambda}) \leq f(\mathbf{x}) - \frac{1}{2} \eta \|\boldsymbol{\lambda}\|_2^2.$$

Many other line search methods are possible.



If f is convex (and satisfies some other smoothness and curvature conditions), then $f(\mathbf{x}^{(t)})$ converges to the optimal value at a geometric rate.

$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \quad \text{where} \quad f(\mathbf{w}) := \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \ln(1 + \exp(-y \langle \mathbf{w}, \mathbf{x} \rangle)).$



Easy to check that f is convex.

Question: How do we compute its gradient at a given point $\mathbf{w} \in \mathbb{R}^d$?

Gradient of f at \mathbf{w} :

$$\begin{aligned} \nabla f(\mathbf{w}) &= -\frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \frac{1}{1 + e^{y \langle \mathbf{w}, \mathbf{x} \rangle}} y \mathbf{x} \\ &= -\frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} (1 - P_{\mathbf{w}}(Y=y \mid \mathbf{X}=\mathbf{x})) y \mathbf{x}. \end{aligned}$$

Gradient descent algorithm 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.

$$\begin{aligned} \mathbf{w}^{(t+1)} &:= \mathbf{w}^{(t)} - \eta_t \nabla f(\mathbf{w}^{(t)}) \\ &= \mathbf{w}^{(t)} + \eta_t \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} (1 - P_{\mathbf{w}^{(t)}}(Y=y \mid \mathbf{X}=\mathbf{x})) y \mathbf{x}. \end{aligned}$$

- ▶ In many applications of (convex) optimization, care about solving problems to very high precision.

Example: stop when gradient is close enough to zero ($\|\nabla f(\mathbf{x})\|_2 \leq \epsilon$ for some small parameter $\epsilon > 0$).

- ▶ For machine learning applications: optimization problem based on training data often just a means-to-an-end.

We really just care about true error rate.

- ▶ Running gradient descent to convergence not strictly necessary: **may be beneficial to stop early (e.g., when hold-out error rate starts to increase significantly).**

1. Formulate learning with a general loss function and regularizer as an optimization problem.
2. Convex sets, convex functions, ways to check convexity of a function.
3. Standard form of convex optimization problems, concept of local and global minimizers.
4. Gradient descent algorithm (for unconstrained problems with differentiable objectives); high-level idea of backtracking line search.
5. Stopping condition for machine learning applications.