

# Optimization for Machine Learning

## Introduction and gradient descent

Clément Royer

CIMPA School “Control, Optimization and Model Reduction in Machine Learning”

February 24, 2025

**Dauphine**  
UNIVERSITÉ PARIS

| PSL 

**PR[AI]RIE**  
PaRis Artificial Intelligence Research InstitutE

## Who I am: Clément Royer

- *Maître de conférences* at Dauphine since 2019.
- Research topics: Optimization and applications.
- Email: `clement.royer@lamsade.dauphine.fr`
- Webpage: <https://www.lamsade.dauphine.fr/~croyer>

## Repository:

<https://tinyurl.com/3etmd46y>

## Learning goals

- Have an optimization toolbox for ML;
- Know the theoretical underpinnings;
- Practical experience.

- 1 Optimization problems in ML
- 2 Optimization theory
- 3 Gradient descent
- 4 Beyond gradient descent: Nonsmoothness
- 5 Beyond gradient descent: Regularization

- 1 Optimization problems in ML
- 2 Optimization theory
- 3 Gradient descent
- 4 Beyond gradient descent: Nonsmoothness
- 5 Beyond gradient descent: Regularization

# What's optimization?

- Operations research;
- Decision-making;
- Decision sciences;
- Mathematical programming;
- Mathematical optimization.

⇒ All of these can be considered as optimization.

# What's optimization?

- Operations research;
- Decision-making;
- Decision sciences;
- Mathematical programming;
- Mathematical optimization.

⇒ All of these can be considered as optimization.

## My definition

*The purpose of optimization is to make the best decision out of a set of alternatives.*

## Optimization $\not\subset$ Machine Learning

- Optimization is a mathematical tool;
- Used in many areas: Economics, Chemistry, Physics, Social sciences,...
- Appears in other branches of (applied) mathematics: Linear Algebra, PDEs, Statistics, etc.



## Optimization $\not\subset$ Machine Learning

- Optimization is a mathematical tool;
- Used in many areas: Economics, Chemistry, Physics, Social sciences,...
- Appears in other branches of (applied) mathematics: Linear Algebra, PDEs, Statistics, etc.

## Machine Learning $\not\subset$ Optimization

- Optimization targets a certain problem;
- ML is not just about this problem;
- Other features of ML (data cleaning, hardware,...) will not appear in the optimization.

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w})$$

# Formulation of an (unconstrained) optimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w})$$

- $\mathbf{w}$  represents the optimization variable(s);
- $d$  is the dimension of the problem (we will assume  $d \geq 1$ );
- $f(\cdot)$  is the **objective/cost/loss** function.

# Formulation of an (unconstrained) optimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w})$$

- $\mathbf{w}$  represents the optimization variable(s);
- $d$  is the dimension of the problem (we will assume  $d \geq 1$ );
- $f(\cdot)$  is the **objective/cost/loss** function.

*Maximizing  $f$  is equivalent to minimizing  $-f$ .*

# Example: SVM Classification

**Given:** A dataset  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ .

- $\mathbf{x}_i$  is a **feature** vector in  $\mathbb{R}^d$ ;
- $y_i$  is a **label**.

# Example: SVM Classification

**Given:** A dataset  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ .

- $\mathbf{x}_i$  is a **feature** vector in  $\mathbb{R}^d$ ;
- $y_i$  is a **label**.

## Motivation: text classification

Using  $d$  words for classification:

- $\mathbf{x}_i$  represents the words contained in a text document:

$$[\mathbf{x}_i]_j = \begin{cases} 1 & \text{if word } j \text{ is in document } i, \\ 0 & \text{otherwise.} \end{cases}$$

- $y_i$  is equal to  $+1$  if the document addresses a certain topic of interest, to  $-1$  otherwise.

## Example: SVM Classification (2)

### Learning process

- Given  $\{(\mathbf{x}_i, y_i)\}_i$ , discover a function  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $h(\mathbf{x}_i) \approx y_i \forall i = 1, \dots, n$ .
- Choose the predictor function  $h$  among a set  $\mathcal{H}$  parameterized by a vector  $\mathbf{w} \in \mathbb{R}^d$ :  $\mathcal{H} = \{h \mid h = h(\cdot; \mathbf{w}), \mathbf{w} \in \mathbb{R}^d\}$ ;

# Example: SVM Classification (2)

## Learning process

- Given  $\{(\mathbf{x}_i, y_i)\}_i$ , discover a function  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $h(\mathbf{x}_i) \approx y_i \forall i = 1, \dots, n$ .
- Choose the predictor function  $h$  among a set  $\mathcal{H}$  parameterized by a vector  $\mathbf{w} \in \mathbb{R}^d$ :  $\mathcal{H} = \{h \mid h = h(\cdot; \mathbf{w}), \mathbf{w} \in \mathbb{R}^{\hat{d}}\}$ ;

## Linear model for text classification

- We seek a hyperplane in  $\mathbb{R}^d$  separating the feature vectors associated with  $y_i = +1$  and those associated with  $y_i = -1$ ;
- This corresponds to a linear model  $h(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$ , and we want to choose  $\mathbf{w}$  such that:

$$\forall i = 1, \dots, n, \quad \begin{cases} \mathbf{x}_i^T \mathbf{w} \geq 1 & \text{if } y_i = +1 \\ \mathbf{x}_i^T \mathbf{w} \leq -1 & \text{if } y_i = -1. \end{cases}$$



## Example: SVM Classification (3)

### An objective to optimize over

- Our goal: penalize values of  $\mathbf{w}$  for which  $h(\mathbf{x}_i)$  does not predict  $y_i$  well enough.
- We use the **hinge loss function**

$$\forall (h, y) \in \mathbb{R}^2, \quad \ell(h, y) = \max \{1 - yh, 0\}.$$

### About the hinge loss

- $hy > 1 \Rightarrow \ell(h, y) = 0$ :  $h$  and  $y$  are of the same sign,  $|h| > 1$  so good prediction;
- $hy < -1 \Rightarrow \ell(h, y) > 2$ :  $h$  and  $y$  are of opposite sign and  $|h| > 1$  bad prediction);
- $|hy| \leq 1 \Rightarrow \ell(h, y) \in [0, 2]$ : small penalty (value of  $|h|$  makes the prediction less certain).

An optimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n \max \{1 - y_i(\mathbf{x}_i^T \mathbf{w}), 0\} \quad .$$

An optimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n \max \{1 - y_i(\mathbf{x}_i^T \mathbf{w}), 0\} \quad .$$

- Minimize the sum of the losses for all examples;

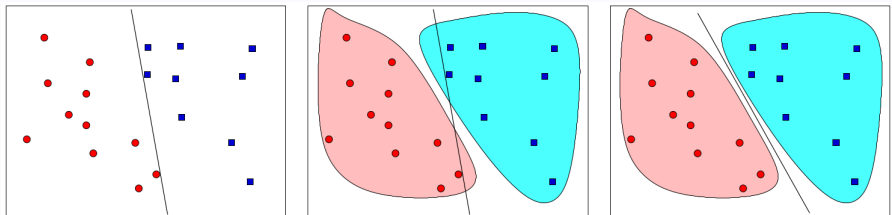
## An optimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n \max \{1 - y_i(\mathbf{x}_i^T \mathbf{w}), 0\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

for  $\lambda \geq 0$ .

- Minimize the sum of the losses for all examples;
- **Regularizing term** to promote small-norm solutions (more on that later).

# Example: SVM Classification (4)



Source: S. J. Wright & B. Recht, Optimization for Data Analysis, 2022.

- Red/Blue dots: data points labeled  $+1/-1$ ;
- Red/Blue clouds: distribution of the text documents;
- Two linear classifiers;
- Rightmost plot: maximal-margin solution.

# Typical optimization problem for ML

- **Data**, e.g.  $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$ .
- **Model class**  $\mathcal{H} = \{\mathbf{h}(\cdot; \mathbf{w}), \mathbf{w} \in \mathbb{R}^d\}$
- **Loss function**  $\ell$ .

## Empirical risk minimization

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(\mathbf{h}(\mathbf{x}_i, \mathbf{w}), \mathbf{y}_i)}_{f(\mathbf{w})} + \lambda \Omega(\mathbf{w})$$

- $f$ : Data-fitting term.
- $\Omega$ : Regularization term.

## Linear regression

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \frac{1}{2n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{w} - y_i)^2.$$

- Simplest data analysis task possible.
- $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ .
- Nontrivial to solve when  $n, d \gg 1$ .

## Linear regression

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \frac{1}{2n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{w} - y_i)^2.$$

- Simplest data analysis task possible.
- $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ .
- Nontrivial to solve when  $n, d \gg 1$ .

## Alternate losses for linear regression

- $\ell_1$  loss:  $\|\mathbf{X}\mathbf{w} - \mathbf{y}\|_1 = \sum_{i=1}^n |\mathbf{x}_i^T \mathbf{w} - y_i|$
- Chebyshev loss:  $\|\mathbf{X}\mathbf{w} - \mathbf{y}\|_\infty = \max_{1 \leq i \leq n} |\mathbf{x}_i^T \mathbf{w} - y_i|$ .
- And more!



## Binary classification (using CNNs)

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \text{CNN}(\mathbf{x}_i))) + \lambda \|\mathbf{w}\|_1.$$

- Cross-entropy/Logistic loss.
- $\mathbf{x}_i \in \mathbb{R}^{d_0 \times d_0 \times c_0}$  (image),  $y_i \in \{-1, 1\}$  (class).
- CNN :  $\mathbf{x}_i = \mathbf{z}^{(0)} \mapsto \mathbf{z}^{(1)} \mapsto \dots \mapsto \mathbf{z}^{(L)}$ , where

$$\mathbf{z}_{ijk}^{(l)} = \phi \left( \sum_{m,n,p} \mathbf{W}_{m,n,p,k}^{(l-1)} \mathbf{z}_{i-m,j-n,p}^{(l-1)} + \mathbf{b}_k^{(l-1)} \right).$$

$\phi(\mathbf{z}) = [\max(\mathbf{z}_i, 0)]_i$  (ReLU activation).

- $\mathbf{w}$  concatenates all  $(\mathbf{W}^l, \mathbf{b}^l)_{l=0 \dots (L-1)}$ .

**Generic form:** minimize  $\mathbf{w} \in \mathbb{R}^d$   $f(\mathbf{w}) + \lambda\Omega(\mathbf{w})$ .

## Common traits

- Defined based on data.
- Use continuous functions (linear, ReLU, log/exp).

## Distinctive aspects

- Model complexity/Number of parameters.
- Nonlinearity of operations.
- Regularization/Lack thereof.

- 1 Optimization problems in ML
- 2 Optimization theory**
- 3 Gradient descent
- 4 Beyond gradient descent: Nonsmoothness
- 5 Beyond gradient descent: Regularization

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w})$$

- $\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ : Set of solutions (can be empty).
- $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ : Optimal value (can be infinite).

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w})$$

- $\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ : Set of solutions (can be empty).
- $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ : Optimal value (can be infinite).

## Global and local minima

- $\mathbf{w}^*$  is a solution or a **global minimum** of  $f$  if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w} \in \mathbb{R}^d$ .
- $\mathbf{w}^*$  is a **local minimum** of  $f$  if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}, \|\mathbf{w} - \mathbf{w}^*\|_2 \leq \epsilon$  for some  $\epsilon > 0$ .

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w})$$

- $\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ : Set of solutions (can be empty).
- $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ : Optimal value (can be infinite).

## Global and local minima

- $\mathbf{w}^*$  is a solution or a **global minimum** of  $f$  if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w} \in \mathbb{R}^d$ .
  - $\mathbf{w}^*$  is a **local minimum** of  $f$  if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}, \|\mathbf{w} - \mathbf{w}^*\|_2 \leq \epsilon$  for some  $\epsilon > 0$ .
- 
- Finding global/local minima is hard in general!
  - Regularity of  $f$  is needed.

## Class of $\mathcal{C}^1$ functions

$f : \mathbb{R}^d \rightarrow \mathbb{R}$  is continuously differentiable/ $\mathcal{C}^1$  if

- For any  $\mathbf{w} \in \mathbb{R}^d$ , the **gradient**  $\nabla f(\mathbf{w})$  exists.
- $\nabla f : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is continuous.

$\Rightarrow f(\mathbf{v}) \approx f(\mathbf{w}) + \nabla f(\mathbf{w})^T(\mathbf{v} - \mathbf{w})$  for  $\mathbf{v}$  close to  $\mathbf{w}$ .

## Class of $\mathcal{C}^1$ functions

$f : \mathbb{R}^d \rightarrow \mathbb{R}$  is continuously differentiable/ $\mathcal{C}^1$  if

- For any  $\mathbf{w} \in \mathbb{R}^d$ , the **gradient**  $\nabla f(\mathbf{w})$  exists.
- $\nabla f : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is continuous.

$\Rightarrow f(\mathbf{v}) \approx f(\mathbf{w}) + \nabla f(\mathbf{w})^T(\mathbf{v} - \mathbf{w})$  for  $\mathbf{v}$  close to  $\mathbf{w}$ .

## Class of $\mathcal{C}_L^{1,1}$ functions ( $L > 0$ )

$f$  is  $\mathcal{C}_L^{1,1}$  if it is  $\mathcal{C}^1$  and  $\nabla f$  is  $L$ -Lipschitz continuous, i.e.

$$\forall (\mathbf{v}, \mathbf{w}) \in (\mathbb{R}^d)^2, \quad \|\nabla f(\mathbf{v}) - \nabla f(\mathbf{w})\| \leq L\|\mathbf{v} - \mathbf{w}\|.$$

*Ex) Linear regression, logistic regression, etc.*



# Smoothness and optimality conditions

**Problem:** minimize  $\mathbf{w} \in \mathbb{R}^d$   $f(\mathbf{w})$ ,  $f \in \mathcal{C}^1$ .

## First-order necessary condition

If  $\mathbf{w}^*$  is a local minimum of the problem, then

$$\|\nabla f(\mathbf{w}^*)\|_2 = 0.$$

- This condition is only necessary;
- A point such that  $\|\nabla f(\mathbf{w}^*)\| = 0$  can also be a local maximum or a saddle point.

# Smoothness and optimality conditions

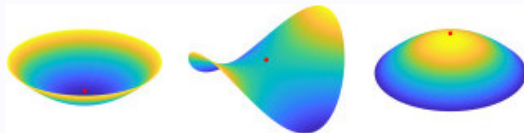
**Problem:** minimize  $\mathbf{w} \in \mathbb{R}^d$   $f(\mathbf{w})$ ,  $f \in \mathcal{C}^1$ .

## First-order necessary condition

If  $\mathbf{w}^*$  is a local minimum of the problem, then

$$\|\nabla f(\mathbf{w}^*)\|_2 = 0.$$

- This condition is only necessary;
- A point such that  $\|\nabla f(\mathbf{w}^*)\| = 0$  can also be a local maximum or a saddle point.



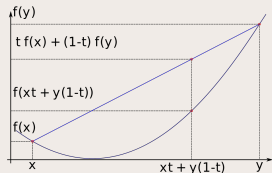
Picture from (Wright and Ma '22).

# Another notion of regularity: Convexity

Generic definition (+Wikicommons picture)

A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is **convex** if

$$\forall (\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2, \forall t \in [0, 1], \\ f(t\mathbf{u} + (1-t)\mathbf{v}) \leq t f(\mathbf{u}) + (1-t) f(\mathbf{v}).$$

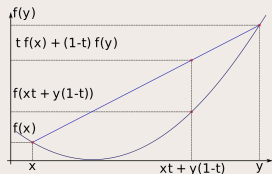


# Another notion of regularity: Convexity

## Generic definition (+Wikicommons picture)

A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is **convex** if

$$\forall (\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2, \forall t \in [0, 1], \\ f(t\mathbf{u} + (1-t)\mathbf{v}) \leq tf(\mathbf{u}) + (1-t)f(\mathbf{v}).$$



## Examples in ML

- Linear function  $\mathbf{w} \mapsto \mathbf{a}^T \mathbf{w} + b$
- Norms  $\|\mathbf{w}\|_2$ ,  $\|\mathbf{w}\|_1$ ,  $\|\mathbf{w}\|_2^2$ .
- Logistic loss.

## Convexity and gradient

A continuously differentiable function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex if and only if

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T (\mathbf{v} - \mathbf{u}).$$

## Convexity and gradient

A continuously differentiable function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex if and only if

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T (\mathbf{v} - \mathbf{u}).$$

**A key inequality in optimization.**

minimize  $f(\mathbf{w})$ ,  $f$  convex.  
 $\mathbf{w} \in \mathbb{R}^d$

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}), f \text{ convex.}$$

## Theorem

Every local minimum of  $f$  is a global minimum.



$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}), f \text{ convex.}$$

## Theorem

Every local minimum of  $f$  is a global minimum.

## Corollary

If  $f$  is  $\mathcal{C}^1$ ,

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} f(\mathbf{w}) = \{ \bar{\mathbf{w}} \mid \|\nabla f(\bar{\mathbf{w}})\|_2 = 0 \}.$$

*Any point with a zero gradient is a global minimum!*

## Definition

A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  in  $\mathcal{C}^1$  is  $\mu$ -strongly convex (or strongly convex of modulus  $\mu > 0$ ) if for all  $(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2$  and  $t \in [0, 1]$ ,

$$f(t\mathbf{u} + (1-t)\mathbf{v}) \leq tf(\mathbf{u}) + (1-t)f(\mathbf{v}) - \frac{\mu}{2}t(1-t)\|\mathbf{v} - \mathbf{u}\|_2^2.$$

## Definition

A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  in  $\mathcal{C}^1$  is  $\mu$ -strongly convex (or *strongly convex of modulus*  $\mu > 0$ ) if for all  $(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2$  and  $t \in [0, 1]$ ,

$$f(t\mathbf{u} + (1-t)\mathbf{v}) \leq tf(\mathbf{u}) + (1-t)f(\mathbf{v}) - \frac{\mu}{2}t(1-t)\|\mathbf{v} - \mathbf{u}\|_2^2.$$

## Theorem

Any strongly convex function in  $\mathcal{C}^1$  has a unique global minimizer.

## Gradient and strong convexity

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^1$ . Then,

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T(\mathbf{v} - \mathbf{u}) + \frac{\mu}{2}\|\mathbf{v} - \mathbf{u}\|_2^2.$$

- 1 Optimization problems in ML
- 2 Optimization theory
- 3 Gradient descent**
- 4 Beyond gradient descent: Nonsmoothness
- 5 Beyond gradient descent: Regularization

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}).$$

**Assumptions:**  $f$  smooth ( $\mathcal{C}^1$ ), bounded below.

# General optimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}).$$

**Assumptions:**  $f$  smooth ( $\mathcal{C}^1$ ), bounded below.

## Key properties

- Smoothness: We will exploit the gradient of  $f$ .
- In presence of convexity, get better guarantees.

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}), \quad f \in \mathcal{C}_L^{1,1}.$$

Consider any  $\mathbf{w} \in \mathbb{R}^d$ . Then, one of the two assertions below holds:

- 1 Either  $\mathbf{w}$  is a local minimum and  $\nabla f(\mathbf{w}) = 0$ ;
- 2 Or the function  $f$  decreases **locally** from  $\mathbf{w}$  in the direction of  $-\nabla f(\mathbf{w})$ .

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}), \quad f \in \mathcal{C}_L^{1,1}.$$

Consider any  $\mathbf{w} \in \mathbb{R}^d$ . Then, one of the two assertions below holds:

- 1 Either  $\mathbf{w}$  is a local minimum and  $\nabla f(\mathbf{w}) = 0$ ;
- 2 Or the function  $f$  decreases **locally** from  $\mathbf{w}$  in the direction of  $-\nabla f(\mathbf{w})$ .

Key argument (Taylor expansion)

$$f(\mathbf{v}) \approx f(\mathbf{w}) + \nabla f(\mathbf{w})^\top (\mathbf{v} - \mathbf{w}) \quad \text{for } \mathbf{v} \text{ close to } \mathbf{w}.$$



**Inputs:**  $\mathbf{w}_0 \in \mathbb{R}^d$ ,  $\alpha_0 > 0$ ,  $k = 0$ .

- 1 Evaluate  $\nabla f(\mathbf{w}_k)$ .
- 2 Set  $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)$ .
- 3 Increment  $k$  by 1 and go to Step 1.

**Inputs:**  $\mathbf{w}_0 \in \mathbb{R}^d$ ,  $\alpha_0 > 0$ ,  $k = 0$ .

- 1 Evaluate  $\nabla f(\mathbf{w}_k)$ .
- 2 Set  $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)$ .
- 3 Increment  $k$  by 1 and go to Step 1.

## Stopping criterion

- Convergence criterion (optional): Stop when  $\|\nabla f(\mathbf{w}_k)\|_2 < \varepsilon$ ;
- Budget criterion (optional): Stop when  $k = k_{\max}$ .

# Key: choosing the stepsize $\alpha_k$

## Constant stepsize

Set  $\alpha_k = \alpha > 0$  for all  $k$ .

- Must be chosen carefully (see lab session).
- Can be set according to properties of  $f$  (see theory).

## Decreasing stepsize

Choose  $\alpha_k$  such that  $\alpha_k \rightarrow 0$ .

- Guarantees that  $f$  will decrease eventually (for small stepsizes);
- But steps get smaller and smaller.

# Choosing the stepsize (2)

## What's done in optimization

- Line search: At every iteration,  $\alpha_k$  is obtained by *backtracking* on a subset of values (ex:  $1, \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots$ ).
- The chosen value must satisfy certain conditions (ex: decreasing the function value).

## What's done in optimization for ML

- Start with a fixed value until the method starts stalling (gradient gets small);
- Decrease the step size, then repeat.

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x}), \quad f \in \mathcal{C}_L^{1,1}.$$

## Gradient descent

- Iteration:  $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)$ , stop if  $\nabla f(\mathbf{w}_k) = 0$ .
- Typical choice in theory :  $\alpha_k = \frac{1}{L}$ .

## Theoretical analysis

- Convergence: Show that  $\|\nabla f(\mathbf{w}_k)\|_2 \rightarrow 0$ ;
- Convergence rate: Look at how fast  $\|\nabla f(\mathbf{w}_k)\|_2$  decreases.
- Worst-case complexity: Equivalent to convergence rate, measures the cost of satisfying  $\|\nabla f(\mathbf{w}_k)\|_2 \leq \epsilon$  for  $\epsilon > 0$ .

## Theorem

If  $f \in \mathcal{C}_L^{1,1}$  and  $\alpha_k = \frac{1}{L}$ ,

$$\min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\|_2 \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$$

after  $K \geq 1$  iterations.

## Theorem

If  $f \in \mathcal{C}_L^{1,1}$  and  $\alpha_k = \frac{1}{L}$ ,

$$\min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\|_2 \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$$

after  $K \geq 1$  iterations.

## A key inequality for the proof

$$\forall(\mathbf{v}, \mathbf{w}), \quad f(\mathbf{v}) \leq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top(\mathbf{v} - \mathbf{w}) + \frac{L}{2} \|\mathbf{v} - \mathbf{w}\|_2^2.$$

- Another key inequality in optimization.
- With  $\mathbf{v} = \mathbf{w}_{k+1}$  and  $\mathbf{w} = \mathbf{w}_k$ , gives decrease in  $\mathcal{O}(\|\nabla f(\mathbf{w}_k)\|_2^2)$ .

## Theorem

Let  $f \in \mathcal{C}_L^{1,1}$  be convex and  $\alpha_k = \frac{1}{L}$  in GD. Then, for  $K \geq 1$ ,

- 1 If  $f$  is convex,

$$f(\mathbf{w}_K) - f^* \leq \mathcal{O}\left(\frac{1}{K}\right).$$

- 2 If  $f$  is  $\mu$ -strongly convex,

$$f(\mathbf{w}_K) - f^* \leq \mathcal{O}\left(\left(1 - \frac{\mu}{L}\right)^K\right).$$



## Theorem

Let  $f \in \mathcal{C}_L^{1,1}$  be convex and  $\alpha_k = \frac{1}{L}$  in GD. Then, for  $K \geq 1$ ,

- 1 If  $f$  is convex,

$$f(\mathbf{w}_K) - f^* \leq \mathcal{O}\left(\frac{1}{K}\right).$$

- 2 If  $f$  is  $\mu$ -strongly convex,

$$f(\mathbf{w}_K) - f^* \leq \mathcal{O}\left(\left(1 - \frac{\mu}{L}\right)^K\right).$$

## Interpretation

Nonconvex	Convex	Strongly convex
$\mathcal{O}(1/\sqrt{K})$	$\mathcal{O}(1/K)$	$\mathcal{O}(t^K)$

*Stronger guarantees for convex problems at lower cost.*

# Conclusion: Gradient descent

## A versatile algorithm

- Applies as long as  $f$  has a gradient.
- Various implementations (stepsizes).
- Theoretical guarantees for convex/nonconvex problems.

## Going further

- What if the function does not have a gradient?
- What about the problem structure?

- 1 Optimization problems in ML
- 2 Optimization theory
- 3 Gradient descent
- 4 Beyond gradient descent: Nonsmoothness**
- 5 Beyond gradient descent: Regularization

## The linear SVM problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \max\{1 - y_i \mathbf{x}_i^T \mathbf{w}, 0\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

- The hinge loss is not continuously differentiable!
- But it is continuous and convex...

## Definition

A function is called **nonsmooth** if it is not differentiable everywhere.

NB: Nonsmooth  $\neq$  Discontinuous.

## Example of nonsmooth functions

- $w \mapsto |w|$  from  $\mathbb{R}$  to  $\mathbb{R}$ ;
- $w \mapsto \|w\|_1$  from  $\mathbb{R}^d$  to  $\mathbb{R}$ ;
- ReLU:  $w \mapsto \max\{w, 0\}$  from  $\mathbb{R}^d$  to  $\mathbb{R}$ .

## Definition

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a convex function. A vector  $\mathbf{g} \in \mathbb{R}^n$  is called a *subgradient* of  $f$  at  $\mathbf{w} \in \mathbb{R}^n$  if

$$\forall \mathbf{z} \in \mathbb{R}^n, \quad f(\mathbf{z}) \geq f(\mathbf{w}) + \mathbf{g}^T(\mathbf{z} - \mathbf{w}).$$

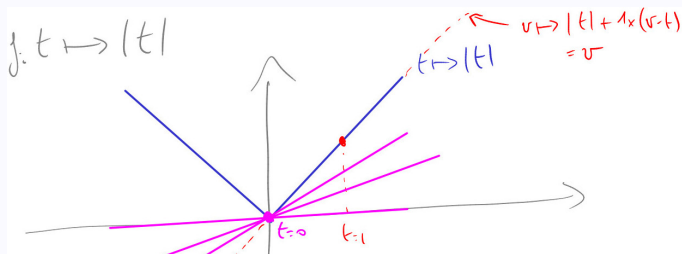
The set of all subgradients of  $f$  at  $\mathbf{w}$  is called the *subdifferential* of  $f$  at  $\mathbf{w}$ , and denoted by  $\partial f(\mathbf{w})$ .

- If  $f$  differentiable at  $\mathbf{w}$ ,  $\partial f(\mathbf{w}) = \{\nabla f(\mathbf{w})\}$ ;
- $0 \in \partial f(\mathbf{w}) \Leftrightarrow \mathbf{w}$  minimum of  $f$ !

**Example:** Let  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(w) = |w|$ .

$$\partial f(w) = \begin{cases} -1 & \text{if } w < 0 \\ 1 & \text{if } w > 0 \\ [-1, 1] & \text{if } w = 0. \end{cases}$$

# Subdifferential: Illustration



$$\partial(|\cdot|)(t) = \begin{cases} -1 & \text{if } t < 0 \\ 1 & \text{if } t > 0 \\ [-1, 1] & \text{if } t = 0. \end{cases}$$

## Iteration for nonsmooth convex $f$

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \mathbf{g}_k, \quad \mathbf{g}_k \in \partial f(\mathbf{w}_k).$$

- Depends on the subgradient: a subgradient can be a direction of increase!
- $\alpha_k$  typically constant or decreasing.



## Iteration for nonsmooth convex $f$

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \mathbf{g}_k, \quad \mathbf{g}_k \in \partial f(\mathbf{w}_k).$$

- Depends on the subgradient: a subgradient can be a direction of increase!
- $\alpha_k$  typically constant or decreasing.

## Guarantees

Let  $\bar{\mathbf{w}}_K = \frac{1}{\sum_{k=0}^{K-1}} \sum_{k=0}^{K-1} \alpha_k \mathbf{w}_k$ . Then,

$$f(\bar{\mathbf{w}}_K) - f^* \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right).$$

*Worst rate than gradient descent but a lot more general!*

- 1 Optimization problems in ML
- 2 Optimization theory
- 3 Gradient descent
- 4 Beyond gradient descent: Nonsmoothness
- 5 Beyond gradient descent: Regularization**

## The linear SVM problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \max\{1 - y_i \mathbf{x}_i^T \mathbf{w}, 0\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

- The problem is **regularized** (by a data-independent term);
- The purpose of regularization is to enforce specific properties/structure on a solution.

# General form of a regularized problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} \underbrace{f(\mathbf{w})}_{\text{loss function}} + \underbrace{\lambda \Omega(\mathbf{w})}_{\text{regularization term}} .$$

where  $\lambda > 0$  is called a regularization parameter.

## Example: Ridge regularization

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 .$$

Interpretations:

- Equivalent to enforcing a constraint on  $\|\mathbf{w}\|_2^2 = \sum_{i=1}^d w_i^2$ ;
- Penalizes  $\mathbf{w}$ s with large components;
- The variance of the solution w. r. t. the data is reduced;
- The objective function is strongly convex.

# Solving regularized problems

## Setup: Composite optimization

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) + \lambda \Omega(\mathbf{w}).$$

- $f \in C^{1,1}$ ;
- $\Omega$  convex but nonsmooth.

## Proximal approach

- Classical optimization paradigm: replace a problem by a sequence of easier (sub)problems;
- Exploit smoothness of  $f$ , use the structure of  $\Omega$  to solve the subproblems;
- Those should be solvable **efficiently**.

## Iteration of PGD

$$\mathbf{w}_{k+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2\alpha_k} \|\mathbf{w} - \mathbf{w}_k\|_2^2 + \lambda \Omega(\mathbf{w}) \right\}.$$

- If  $\Omega \equiv 0$ , the solution is  $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)$ : **This is the Gradient Descent iteration!**
- In general, the cost of an iteration is 1 gradient call + **1 proximal subproblem solve**.

## Properties

- Complexity bounds exist for nonconvex and mostly for convex  $f$ ;
- Stepsize choices can be designed based on those for GD.

## Sparsity-inducing regularizers

- Want solution  $\mathbf{w} \in \mathbb{R}^d$  with few nonzero components.
- For linear models, amounts to feature selection.

## A better approach: LASSO regularization

*LASSO=Least Absolute Shrinkage and Selection Operator*

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1, \quad \|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i|.$$

- $\|\cdot\|_1$  is convex, continuous, and a norm.
- Nonsmooth but subgradients can be computed.
- No close form even for linear regression  $\Rightarrow$  Proximal gradient!

## Context

- Solve minimize  $\mathbf{w} \in \mathbb{R}^d$   $f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$ .
- Common problem in image processing: Proximal gradient  $\equiv$  ISTA.
- Explicit form of the proximal subproblem solution.



## Context

- Solve minimize  $\mathbf{w} \in \mathbb{R}^d$   $f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$ .
- Common problem in image processing: Proximal gradient  $\equiv$  ISTA.
- Explicit form of the proximal subproblem solution.

## Iteration of ISTA: Iterative Soft-Thresholding Algorithm

Define  $\mathbf{w}_{k+1}$  componentwise: for any  $i \in \{1, \dots, d\}$ ,

$$[\mathbf{w}_{k+1}]_i = \begin{cases} [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i + \alpha_k \lambda & \text{if } [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i < -\alpha_k \lambda \\ [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i - \alpha_k \lambda & \text{if } [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i > \alpha_k \lambda \\ 0 & \text{if } [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i \in [-\alpha_k \lambda, \alpha_k \lambda]. \end{cases}$$

## Optimization problems in ML

- Common feature: Depend on data.
- Distinctive features: Convexity, smoothness, regularization.

## Gradient descent

- The basic block for optimization.
- Applies to convex and nonconvex functions.
- Some freedom in the implementation (see lab session).

## Beyond gradient descent

- Nonsmoothness  $\Rightarrow$  Subgradient methods!
- Regularization  $\Rightarrow$  Proximal methods!
- Data dependency?  $\Rightarrow$  See next lecture.

## Textbooks:

- A. Beck, *First-order methods in optimization*, MOS-SIAM Series on Optimization, 2017.  
⇒ *Chapter 10* is related to proximal methods, and contains many examples of explicit proximal step calculations.
- J. Wright and Y. Ma, *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications*, Cambridge University Press, 2022.  
⇒ Numerous applications, freely available online.
- S. J. Wright and B. Recht, *Optimization for Data Analysis*, Cambridge University Press, 2022.  
⇒ Textbook with full analysis for gradient descent.