# Optimization for Machine Learning Part II: Introduction to SGD

Lionel Tondji

African Master's in Machine Intelligence

July 18, 2024

# Structure of Optimization Problems Arising in Training Supervised machine Learning Models



In this course, we are interested in the optimization problem

$$\min_{w \in \mathbb{R}^d} f(w) + R(w) \tag{1}$$



In this course, we are interested in the optimization problem

$$\min_{w \in \mathbb{R}^d} f(w) + R(w) \tag{1}$$

Typical structure of f:

Infinite sum

$$f(w) = E_{\zeta \sim \mathcal{D}}[f_{\zeta}(w)] \tag{2}$$



In this course, we are interested in the optimization problem

$$\min_{w \in \mathbb{R}^d} f(w) + R(w) \tag{1}$$

Typical structure of f:

Infinite sum

$$f(w) = E_{\zeta \sim \mathcal{D}}[f_{\zeta}(w)] \tag{2}$$

• Finite sum:

$$f(w) = \frac{1}{n} \sum_{i=1}^{n} f_i(w)$$
 (3)



In this course, we are interested in the optimization problem

$$\min_{w \in \mathbb{R}^d} f(w) + R(w) \tag{1}$$

Typical structure of f:

Infinite sum

$$f(w) = E_{\zeta \sim \mathcal{D}}[f_{\zeta}(w)] \tag{2}$$

Finite sum :

$$f(w) = \frac{1}{n} \sum_{i=1}^{n} f_i(w)$$
 (3)

Finite sum of Finite Sums :

$$f_i(w) = \frac{1}{m} \sum_{j=1}^{m} f_{ij}(w) \tag{4}$$

2/50

These problems are of keys importance in supervised learning theory and pratice.

Common feature: It is prohibitively expensive to compute the gradient of f, while an unbiased estimator of the gradient can be computed efficiently/cheaply.



In the stochastic optimization problem

$$f(w) = E_{\zeta \sim \mathcal{D}}[f_{\zeta}(w)]$$

w represents a machine learning model described by d
parameters/features (e.g logistic regression or a deep neural network).



In the stochastic optimization problem

$$f(w) = E_{\zeta \sim \mathcal{D}}[f_{\zeta}(w)]$$

- w represents a machine learning model described by d
  parameters/features (e.g logistic regression or a deep neural network).
- ullet D is an unknown distribution of labelled examples,



In the stochastic optimization problem

$$f(w) = E_{\zeta \sim \mathcal{D}}[f_{\zeta}(w)]$$

- w represents a machine learning model described by d
  parameters/features (e.g logistic regression or a deep neural network).
- ullet D is an unknown distribution of labelled examples,
- $f_{\zeta}(w)$  represents the loss of model w on a data point  $\zeta$ , and



In the stochastic optimization problem

$$f(w) = E_{\zeta \sim \mathcal{D}}[f_{\zeta}(w)]$$

- w represents a machine learning model described by d
  parameters/features (e.g logistic regression or a deep neural network).
- ullet D is an unknown distribution of labelled examples,
- $f_{\zeta}(w)$  represents the loss of model w on a data point  $\zeta$ , and
- f is the generalization error.

Problem (1) seeks to find the model w minimizing the generalization error

- ① In statistical learning theory one assumes that while  $\mathcal D$  is not known, samples  $\zeta \sim \mathcal D$  are available.
- ② In such case,  $\nabla f(w)$  is not computable, while  $\nabla f_{\zeta}(w)$ , which is an unbiased estimator of the gradient of f at w, is easily computable

#### Finite Sum Problems

In this course we will focus on functions f which arise as averages of very large number of (smooth) functions:

$$f(w) = \frac{1}{n} \sum_{i=1}^{n} f_i(w)$$



#### Finite Sum Problems

In this course we will focus on functions f which arise as averages of very large number of (smooth) functions:

$$f(w) = \frac{1}{n} \sum_{i=1}^{n} f_i(w)$$

- This problem often arises by approximation of the stochastic optimization loss function (2) via Monte Carlo Integration.
- Known as the empirical risk minimization (ERM) problem.
- ERM is currently the dominant paradigm for solving supervised learning problems.



#### Finite Sum Problems

In this course we will focus on functions f which arise as averages of very large number of (smooth) functions:

$$f(w) = \frac{1}{n} \sum_{i=1}^{n} f_i(w)$$

- This problem often arises by approximation of the stochastic optimization loss function (2) via Monte Carlo Integration.
- Known as the empirical risk minimization (ERM) problem.
- ERM is currently the dominant paradigm for solving supervised learning problems.
- If index i is chosen uniformly at random from  $[n] = \{1, 2, ..., n\}$ ,  $\nabla f_i(w)$  is an unbiased estimator of  $\nabla f(w)$ .
- Typically,  $\nabla f_i(w)$  is about n times less expensive to compute than  $\nabla f(w)$ .

# Distributed Training

In distributed of supervised models, one considers the finite sum problem (3), with n being the number of machines, and each  $f_i$ 

also having a finite sum structure, i.e,

$$f_i(w) = \frac{1}{m} \sum_{j=1}^{m} f_{ij}(w)$$
 (5)

where m corresponds to the number of training examples stored on machine i.



# Distributed Training

In distributed of supervised models, one considers the finite sum problem (3), with n being the number of machines, and each  $f_i$ 

• also having a finite sum structure, i.e,

$$f_i(w) = \frac{1}{m} \sum_{j=1}^{m} f_{ij}(w)$$
 (5)

where m corresponds to the number of training examples stored on machine i.

or an infinite-sum structure, i.e,

$$f_i(w) = E_{\zeta_i \sim \mathcal{D}_i}[f_{i\zeta_i}(w)] \tag{6}$$

where  $\mathcal{D}_i$  is the distribution of data stored on machine i.



#### **SGD**

Stochastic gradient descent (SGD) is a state-of-the-art algorithmic paradigm for solving optimization problem (1) in situations where f is either of structure (2) or (3).

In its generic form (proximal) SGD defines the new iterate by subtracting a multiple of a stochastic gradient from the current iterate, and subsequently applying the proximal operator of R:

$$w_{k+1} = prox_{\gamma R}(w_k - \gamma g^k) \tag{7}$$



#### **SGD**

Stochastic gradient descent (SGD) is a state-of-the-art algorithmic paradigm for solving optimization problem (1) in situations where f is either of structure (2) or (3).

In its generic form (proximal) SGD defines the new iterate by subtracting a multiple of a stochastic gradient from the current iterate, and subsequently applying the proximal operator of R:

$$w_{k+1} = prox_{\gamma R}(w_k - \gamma g^k) \tag{7}$$

•  $g^k$  is an unbiased estimator of the gradient (i.e a "stochastic gradient"):

$$E[g^k/w_k] = \nabla f(w_k) \tag{8}$$



#### **SGD**

Stochastic gradient descent (SGD) is a state-of-the-art algorithmic paradigm for solving optimization problem (1) in situations where f is either of structure (2) or (3).

In its generic form (proximal) SGD defines the new iterate by subtracting a multiple of a stochastic gradient from the current iterate, and subsequently applying the proximal operator of R:

$$w_{k+1} = prox_{\gamma R}(w_k - \gamma g^k) \tag{7}$$

•  $g^k$  is an unbiased estimator of the gradient (i.e a "stochastic gradient"):

$$E[g^k/w_k] = \nabla f(w_k) \tag{8}$$

 $prox_R(x) = arg \min_{u} \left\{ R(u) + \frac{1}{2} ||u - x||^2 \right\}$ 

### The Prox Operator

#### Some facts about the prox operator<sup>1</sup>:

- **1** single-valuedness:  $x \mapsto prox_R(x)$  is a function
- 2 non-expansiveness:

$$\|prox_R(x) - prox_R(y)\| \le \|x - y\|, \quad \forall x, y \in \mathbb{R}^d$$

Moreau decomposition:

$$prox_R(x) - prox_{R^*}(x) = x, \quad \forall x \in \mathbb{R}^d$$

Here  $R^*$  is the Fenchel conjugate<sup>2</sup> of R.

 $^{2}R^{*}(x) = sup_{y \in \mathbb{R}^{d}}\{\langle x, y \rangle - R(y)\}$ 

Lionel Tondji (AIMS-AMMI)

4 1 2 4 1 2 7 4 3

<sup>&</sup>lt;sup>1</sup>Assume  $R: \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$  is proper, closed and convex.

AIMS NECT EINSTEIN INITIATIVE

#### Stochastic Gradient

There are infinitely many ways of obtaining a random vector  $g^k$  satisfying (8)

- Prox: flexibility to construct stochastic gradients in various ways based on problem structure, and in order to target desirable properties such as:
  - convergence speed
  - iteration cost
  - overall complexity
  - parallelizability
  - suitability for given computing architecture
  - communication cost
  - generalization properties



#### Stochastic Gradient

There are infinitely many ways of obtaining a random vector  $g^k$  satisfying (8)

- Cons: A crazy ZOO of methods
  - Little hard to get into the fields, hard to keep up with new results
  - Considerable challenges in terms of convergence analysis. Indeed, if one aims to, as one should, obtain the sharpest bounds possible, dedicated analyses are needed to handle each of the particular variants of SGD.

# Batch SGD = Gradient Descent





#### Gradient Descent

We first describe the (proximal) gradient descent (GD) method for solving the regularized convex optimization problem

$$\min_{w \in \mathbb{R}^d} f(w) + R(w) \tag{9}$$

This is the most basic of all SGD methods, and a starting point for the development of more elaborate variants.

#### Algorithm GD

```
starting points x_0 \in \mathbb{R}^d, learning rate \gamma > 0 for k = 0, 1, 2, \cdots do Set g^k = \nabla f(w_k) w_{k+1} = prox_{\gamma R}(w_k - \gamma g^k) end for
```



The idea is f might be something complicated but the linear approximation is simple.

$$f(w) \simeq \text{linear function}$$
  
=  $f(w_0) + \langle \nabla f(w_0), w - w_0 \rangle$ 



The idea is f might be something complicated but the linear approximation is simple.

$$f(w) \simeq \text{linear function}$$
  
=  $f(w_0) + \langle \nabla f(w_0), w - w_0 \rangle$ 

GD is an iterative algorithm where at each step we minimize a linear approximation but plus an additional term that makes sure we do not roll all the way down hill.



The idea is f might be something complicated but the linear approximation is simple.

$$f(w) \simeq \text{linear function}$$
  
=  $f(w_0) + \langle \nabla f(w_0), w - w_0 \rangle$ 

GD is an iterative algorithm where at each step we minimize a linear approximation but plus an additional term that makes sure we do not roll all the way down hill.

#### Idea of GD:

- start w<sub>0</sub>
- $w_{k+1} = arg \min_{w} f(w_k) + \langle \nabla f(w_k), w w_k \rangle + \frac{1}{2\gamma} ||w w_k||_2^2$
- \* 1<sup>st</sup> term: linear function
- \*  $2^{nd}$  term : quadratic term that penalize w being very far from  $w_k$ .



$$0 + \nabla f(w_k) + \frac{1}{\gamma}(w - w_k) = 0$$
$$\Rightarrow w_{k+1} = w_k - \gamma \nabla f(w_k)$$

**Example**:  $f(x) = 3x^2 + 4x - 2$ 



$$0 + \nabla f(w_k) + \frac{1}{\gamma}(w - w_k) = 0$$
$$\Rightarrow w_{k+1} = w_k - \gamma \nabla f(w_k)$$

**Example**:  $f(x) = 3x^2 + 4x - 2$ 

• Can solve this directly:  $6x + 4 = 0 \Rightarrow x^* = -\frac{2}{3}$ 

$$0 + \nabla f(w_k) + \frac{1}{\gamma}(w - w_k) = 0$$
$$\Rightarrow w_{k+1} = w_k - \gamma \nabla f(w_k)$$

**Example**:  $f(x) = 3x^2 + 4x - 2$ 

- **1** Can solve this directly:  $6x + 4 = 0 \Rightarrow x^* = -\frac{2}{3}$
- ② Apply gradient descent. Initialize at  $x_1$ , take step size  $\gamma$ .
  - ► GD iteration:  $x^+ = x \gamma \nabla f(x) = x \gamma (6x + 4)$



$$0 + \nabla f(w_k) + \frac{1}{\gamma}(w - w_k) = 0$$
$$\Rightarrow w_{k+1} = w_k - \gamma \nabla f(w_k)$$

**Example**:  $f(x) = 3x^2 + 4x - 2$ 

- **1** Can solve this directly:  $6x + 4 = 0 \Rightarrow x^* = -\frac{2}{3}$
- ② Apply gradient descent. Initialize at  $x_1$ , take step size  $\gamma$ .

► GD iteration: 
$$x^+ = x - \gamma \nabla f(x) = x - \gamma (6x + 4)$$

$$\Rightarrow x_{k+1} = x_k - \gamma 6x_k - 4\gamma$$

$$= (1 - 6\gamma)x_k - 4\gamma$$

$$\vdots$$

$$= (1 - 6\gamma)^k x_1 - ((1 - 6\gamma)^{k-1} + (1 - 6\gamma)^{k-2} + \dots + 1)4\gamma$$

$$= (1 - 6\gamma)^k x_1 - \frac{1 - (1 - 6\gamma)^k}{1 - (1 - 6\gamma)}4\gamma$$
AIMS

Need  $|1 - 6\gamma| < 1 \Rightarrow (1 - 6\gamma)^k \to 0$  as  $k \to \infty$ .

$$x_{k+1} = (1 - 6\gamma)^k x_1 + \frac{(1 - 6\gamma)^k}{6\gamma} - \frac{1}{6\gamma} 4\gamma$$
  
=  $(1 - 6\gamma)^k \left[ x_1 + \frac{2}{3} \right] - \frac{2}{3} \to -\frac{2}{3}$  as  $k$  grows

 $x_k \to -\frac{2}{3}$  very quickly with linear rate.



Need  $|1 - 6\gamma| < 1 \Rightarrow (1 - 6\gamma)^k \to 0$  as  $k \to \infty$ .

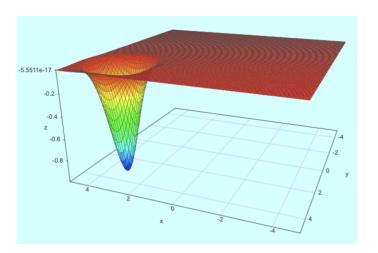
$$x_{k+1} = (1 - 6\gamma)^k x_1 + \frac{(1 - 6\gamma)^k}{6\gamma} - \frac{1}{6\gamma} 4\gamma$$
$$= (1 - 6\gamma)^k \left[ x_1 + \frac{2}{3} \right] - \frac{2}{3} \to -\frac{2}{3} \text{as } k \text{ grows}$$

 $x_k \to -\frac{2}{3}$  very quickly with linear rate.

- Good News: GD for this function, converges very quickly. Error goes down with  $(1-6\gamma)^k$  fast when  $\gamma<\frac{1}{6}$ .
- **2** Step size: small enough so that  $|1 6\gamma| < 1$
- **1 Improvement at every iteration**: Exercise: check that  $f(w_{k+1}) \le f(w_k)$



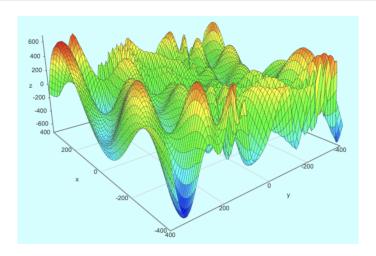
# Optimization is hard (in general)



$$f(x,y) = -\cos(x)\cos(y)\exp(-(x-\pi)^2 - (y-\pi)^2)$$
 in  $[-5,5]^2$  AIMS

Optimization for Machine Learning

# Optimization is hard (in general)



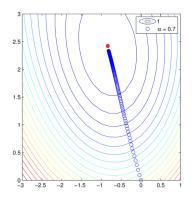
$$f(x,y) = -(y+47)\sin\sqrt{|\frac{x}{2}+(y+47)|} - x\sin\sqrt{|\frac{x}{2}-(y+47)|} \quad \text{in } [-400,400]$$

Lionel Tondji (AIMS-AMMI)

# Gradient Descent Example

A Logistic Regression problem using the fourclass labelled data from LIBSVM (n, d) = (862, 2). Logistic Regression :

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y^i \langle w, x^i \rangle}) + \lambda \|w\|_2^2$$

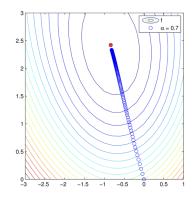




# Gradient Descent Example

A Logistic Regression problem using the fourclass labelled data from LIBSVM (n, d) = (862, 2). Logistic Regression :

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y^i \langle w, x^i \rangle}) + \lambda \|w\|_2^2$$



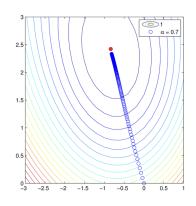
Can we prove that this always works?



# Gradient Descent Example

A Logistic Regression problem using the fourclass labelled data from LIBSVM (n, d) = (862, 2). Logistic Regression :

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y^i \langle \mathbf{w}, \mathbf{x}^i \rangle}) + \lambda \|\mathbf{w}\|_2^2$$



- Can we prove that this always works?
- No! There is no universal optimization method. The "no free lunch" of Optimization
  - Need assumptions: Convex and smooth training problems

### Main assumption

#### Nice property:

If 
$$\nabla f(w^*) = 0$$
 then  $f(w^*) \le f(w)$ ,  $\forall w \in \mathbb{R}^d$ 

 $\Rightarrow$  All stationary points are global minima.



### Main assumption

#### Nice property:

If 
$$\nabla f(w^*) = 0$$
 then  $f(w^*) \le f(w)$ ,  $\forall w \in \mathbb{R}^d$ 

⇒ All stationary points are global minima.

#### Lemma

Convexity  $\Rightarrow$  Nice property.

If 
$$f(w) \ge f(y) + \langle \nabla f(y), w - y \rangle$$
,  $\forall w, y \in \mathbb{R}^d$  then Nice property holds.

PROOF: Choose  $y = w^*$ .



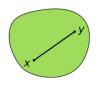
### Convex sets - Definition







### Convex sets - Definition





A set  $\mathbb{C} \subseteq \mathbb{R}^n$ , is convex if

$$\forall x, y \in \mathbb{C}, \ \forall \ \lambda \in [0, 1], \quad \lambda x + (1 - \lambda)y \in \mathbb{C}.$$
 (10)

Why it is important?





### Convex sets

#### Example

#### **Definition**

 $M \in \mathbb{R}^{n \times n}$  matrix is p.s.d if:

- Symmetric
- $x^{\top}Mx \geq 0, \ \forall \ x \in \mathbb{R}^n$

We denote by  $S_+^n$ , the set of symmetric p.s.d matrices.

The set of p.s.d matrices is a convex set. Let  $M_1, M_2 \in \mathcal{S}^n_+$ , we want to show that

$$M = \lambda M_1 + (1 - \lambda)M_2 \in \mathcal{S}_+^n$$
, we want to show that  $M = \lambda M_1 + (1 - \lambda)M_2 \in \mathcal{S}_+^n$ ,  $\lambda \in [0, 1]$ 



#### Convex sets

Copositive matrices: a hard convex set

A symmetric matrix M is **copositive** if

$$x^{\top} M x \ge 0, \ \forall \ x \in \mathbb{R}^n_+$$

We denote by  $C^n$ , the set of symmetric copositive matrices.

- **1** Q: Is the set of copositives matrices bigger or smaller than  $S_+^n$ ?
- ② In general it is intractable to determine whether a matrix M is copositive.

### Example

- Every matrix with only non negative entries is copositive Hence  $M = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$  is copositive but not p.s.d due to det(M) = -1.
- $oldsymbol{2}$  Every p.s.d is also copositive but the converse is false i.e.  $\mathcal{S}^n_+ \subseteq \mathcal{C}^n$

### Exercise

Characterize the triples 
$$(x, y, z) \in \mathbb{R}^3$$
 for which the matrix  $M = \begin{pmatrix} x & z \\ z & y \end{pmatrix}$  is

- copositive
- p.s.d





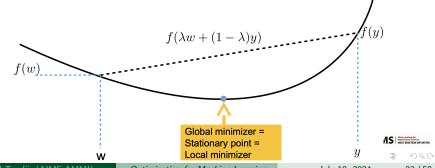
# Def 1: Convexity

#### Definition

We say that  $f: dom(f) \subseteq \mathbb{R}^d \to \mathbb{R}$  is **convex** if dom(f) is **convex** and

$$f(\lambda w + (1 - \lambda)y) \le \lambda f(w) + (1 - \lambda)f(y), \tag{11}$$

 $\forall w, y \in dom(f), \lambda \in [0, 1]$ 

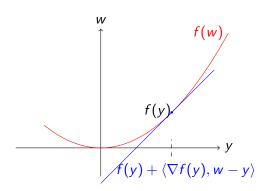


### Def 2: Convexity - First derivative

#### Definition

A differential function  $f: dom(f) \subseteq \mathbb{R}^d \to \mathbb{R}$  is **convex** iff

$$f(w) \ge f(y) + \langle \nabla f(y), w - y \rangle, \quad \forall \ w, y \in dom(f)$$
 (12)



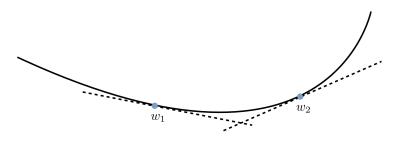


### Def 3: Convexity - Second derivative

#### Definition

A twice differentiable function  $f:dom(f)\subseteq\mathbb{R}^d\to\mathbb{R}$  is **convex** iff

$$\nabla^2 f(w) \succcurlyeq 0 \Leftrightarrow v^{\top} \nabla^2 f(w) v \ge 0, \quad \forall \ w, v \in dom(f)$$
 (13)



$$w_1 \le w_2 \Rightarrow f'(w_1) \le f'(w_2).$$
  
 $(\nabla^2 f(w) \text{ p.s.d i.e. all eigen values of } \nabla^2 f(w) \text{ are } \ge 0).$ 



### Convexity: Examples

- **1** Norms and squared norms :  $x \mapsto ||x||$ ,  $x \mapsto ||x||^2$
- ② Negative log and logistic :  $x \mapsto \log(x)$ ,  $x \mapsto \log(1 + e^{-y\langle a, x \rangle})$
- **3** Hinge Loss :  $x \mapsto max\{0, 1 yx\}$
- Negatives log determinant, exponentiation ... etc



# Def 2': Convex functions (non-smooth)

#### Definition

A function f is **convex** if  $\forall y$ ,  $\exists g$  such that

$$f(w) \ge f(y) + \langle g, w - y \rangle \tag{14}$$

i.e at every point of the function there exists some linear function which touch the function at that point but it is not necessarily unique as a gradient.

• If f is convex and differentiable, then  $g = \nabla f(y)$  satisfies this. It is unique.



# Subgradient and Subdifferential

#### Definition

For a convex function f, a vector g such that

$$f(w) \ge f(y) + \langle g, w - y \rangle, \quad \forall y$$

is called a subgradient.

The set of subgradients of f at a point w is called the subdifferential of f at w and denoted by  $\partial f(w)$ 



### More examples

- **1** Applying the three definitions for  $f(x) = x^{\top}Qx$ ,  $Q \geq 0$
- The max of convex function is convex.
- The min may not be.
- 4 Largest element of a vector :

$$f(x_1,\ldots,x_n)=$$
 maximum element (defined on  $\mathbb{R}^n_+$ )



### More examples

- **1** Applying the three definitions for  $f(x) = x^{\top}Qx$ ,  $Q \geq 0$
- 2 The max of convex function is convex.
- The min may not be.
- 4 Largest element of a vector :

$$f(x_1,\ldots,x_n)=\mathsf{maximum}$$
 element (defined on  $\mathbb{R}^n_+$ )

In fact

$$f(x_1, ..., x_n) = f(x)$$

$$= \max\{e_1^\top x, e_2^\top x, ..., e_n^\top x\}$$

$$= \max_{1 \le i \le n} f_i(x)$$

each  $f_i(x)$  is convex.



# More example: The largest eigenvalue of a symmetric matrix

The function *f* defined by:

$$f(Q) = \lambda_{\mathsf{max}}(Q),$$

is convex. We can show that

$$\lambda_{\sf max}({\it Q}) = {\sf max} \ {\sf of} \ {\sf convex} \ {\sf functions}$$

**Recall**: Q symmetric  $\Rightarrow x^{\top}Qx \leq \lambda_{\max}||x||_2^2$  and  $x^{\top}Qx = \lambda_{\max}||x||_2^2$  when

x= eigen vector corresponding to  $\lambda_{\max}$ 

because 
$$x^{\top}Qx = x^{\top}(\lambda_{\max}x) = \lambda_{\max}x^{\top}x = \lambda_{\max}||x||_2^2$$

$$\lambda_{\mathsf{max}} = \sup_{\|x\|_2 = 1} x^\top Q x$$

 $x^{\top}Qx$  as a function of Q is linear since  $\langle x^{\top}x, Q \rangle = x^{\top}Qx$ 

AIMS Metabolical Science Metabolica Metabolica Metabolica Metabolica Metabolica Metabo

# Monotonicity

If f is convex, the gradient (subgradient) of f is monotone i.e

• 
$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \ge 0$$

• 
$$\langle g_x - g_y, x - y \rangle \ge 0$$
 for  $g_x \in \partial f(x), g_y \in \partial f(y)$ 

Proof:



# Monotonicity

If f is convex, the gradient (subgradient) of f is monotone i.e

• 
$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \ge 0$$

• 
$$\langle g_x - g_y, x - y \rangle \ge 0$$
 for  $g_x \in \partial f(x), g_y \in \partial f(y)$ 

#### Proof:

$$f(y) \ge f(x) + \langle g_x, y - x \rangle$$
  

$$f(x) \ge f(y) + \langle g_y, x - y \rangle$$
  

$$f(x) + f(y) \ge f(x) + f(y) + \langle g_x - g_y, y - x \rangle$$

$$\Rightarrow \langle g_x - g_y, y - x \rangle \ge 0$$



# Equivalence: Convexity and monotonicity

- a) If f is convex (def 3), its gradient is monotone.
- b) If the gradient of f is monotone, then f is convex (def 2)

#### Proof:

Recall that  $\int_0^1 F'(t)dt = F(1) - F(0)$  and by assumption  $\nabla^2 f(x) \succcurlyeq 0$ 

$$\int_0^1 (x-y)^\top \nabla^2 f(tx+(1-t)y)dt = \int_0^1 \frac{d}{dt} (\nabla f(t(x-y)+y)dt)$$
$$= \nabla f(x) - \nabla f(y)$$

Taking the inner product with (x - y)

$$\int_0^1 (x-y)^\top \nabla^2 f(tx+(1-t)y)(x-y)dt = \langle \nabla f(x) - \nabla f(y), x-y \rangle \ge 0$$



### Equivalence: Convexity and monotonicity

$$\int_{0}^{1} \nabla f((y-x)t + x)^{\top} (y-x) dt = \int_{0}^{1} \frac{d}{dt} f((y-x)t + x) dt$$
$$= f(y) - f(x)$$

$$\Rightarrow f(y) = f(x) + \int_0^1 \nabla f((y-x)t + x)^\top (y-x) dt \ (\geq \langle \nabla f(x), y-x \rangle)$$

We can show this inequality holds, by showing that the integral is smallest at t=0. We will use monotone property.

$$\langle \nabla f((y-x)t+x) - \nabla f(x), (y-x)t+x-x \rangle \ge 0$$
  
 $\langle \nabla f((y-x)t+x) - \nabla f(x), (y-x) \rangle \ge 0$ 

Let 
$$h(t) = \nabla f((y-x)t + x)^{\top}(y-x)$$
, then  $h(t) - h(0) = \langle \nabla f((y-x)t + x) - \nabla f(x), (y-x) \rangle \ge 0$ 



ロト 4回ト 4 ミト 4 ミト ミーク90

### Equivalence: Convexity and monotonicity

$$h(t) = \nabla f((y-x)t + x)^{\top}(y-x),$$

then

$$h(t) - h(0) = \langle \nabla f((y-x)t + x) - \nabla f(x), (y-x) \rangle \ge 0$$

$$\Rightarrow \int h(t)dt \ge h(0) \cdot 1$$

$$\Rightarrow f(y) = f(x) + \int_0^1 h(t)dt$$

$$\ge f(x) + h(0)$$

$$= f(x) + \langle \nabla f(x), y - x \rangle$$
AIMS

### Example

**Sum of squares**: Let  $a_1, a_2, \ldots, a_n$ , find x to minimize :  $\frac{1}{n} \sum_{i=1}^{n} (a_i - x)^2$ .



### Example

**Sum of squares**: Let  $a_1, a_2, \ldots, a_n$ , find x to minimize :  $\frac{1}{n} \sum_{i=1}^{n} (a_i - x)^2$ . Take derivative, set to 0

$$-\frac{2}{n}\sum_{i=1}^{n}(a_i-x)=0 \Rightarrow \sum_{i=1}^{n}(a_i-x)=0$$
$$\Rightarrow \hat{x}=\frac{1}{n}\sum_{i=1}^{n}a_i$$



### Example

Ridge regression: Let  $(y_i, x_i)_{1 \le i \le n}, x_i \in \mathbb{R}^d$ , find  $\beta$  to minimize :

$$f(\beta) = \frac{1}{n} \sum_{i=1}^{n} (y_i - x_i^{\top} \beta)^2 + \lambda \sum_{i=1}^{n} \beta_i^2$$



### Example

Ridge regression: Let  $(y_i, x_i)_{1 \le i \le n}, x_i \in \mathbb{R}^d$ , find  $\beta$  to minimize :

$$f(\beta) = \frac{1}{n} \sum_{i=1}^{n} (y_i - x_i^{\top} \beta)^2 + \lambda \sum_{i=1}^{n} \beta_i^2$$

$$f(\beta) = \frac{1}{n} ||X\beta - y||_2^2 + \lambda ||\beta||_2^2$$
, take derivative, set to 0

$$\nabla f(\beta) = 0 \Rightarrow \frac{2}{n} X^{\top} (X\beta - y) + 2\lambda \beta = 0$$
$$\Rightarrow (\frac{1}{n} X^{\top} X + \lambda I) \beta = \frac{1}{n} X^{\top} y$$
$$\Rightarrow \beta = (X^{\top} X + \lambda I)^{-1} X^{\top} y$$



#### Non-differentiable functions

Let f be a convex function and we want to solve  $\min_x f(x)$ .

#### Definition

 $\hat{x}$  is **an** optimal solution if  $0 \in \partial f(\hat{x})$ 

$$g_x \in \partial f(x) \Leftrightarrow f(y) \ge f(x) + \langle g_x, y - x \rangle, \quad \forall y$$

If  $0 \in \partial f(\hat{x})$ , the definition reads:

$$f(y) \ge f(\hat{x}) + 0, \quad \forall y$$

i.e.  $\hat{x}$  is **an** optimal solution.



# Example: Sum of absolute values

### Example

Let  $a_1, a_2, \ldots, a_n$ , find x to minimize :  $\frac{1}{n} \sum_{i=1}^n |a_i - x|$ .  $\hat{x}$  optimal if  $0 \in \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \hat{x}} (|a_i - x|)$ 

Note that

$$\frac{\partial}{\partial z}|z| = \begin{cases} -1, & \text{if } z < 0, \\ 1, & \text{if } z > 0, \\ [-1, 1], & \text{if } z = 0 \end{cases}$$

**Exercise**: Apply this definition of subdifferential to the above problem to find the solution.



### Example: Sum of absolute values

### Example

Let  $a_1, a_2, \ldots, a_n$ , find x to minimize :  $\frac{1}{n} \sum_{i=1}^n |a_i - x|$ .  $\hat{x}$  optimal if  $0 \in \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \hat{x}} (|a_i - x|)$ 

Note that

$$\frac{\partial}{\partial z}|z| = \begin{cases} -1, & \text{if } z < 0, \\ 1, & \text{if } z > 0, \\ [-1, 1], & \text{if } z = 0 \end{cases}$$

**Exercise**: Apply this definition of subdifferential to the above problem to find the solution.

**Answer**:  $\hat{x} = \text{median}\{a_1, a_2, \dots, a_n\}$ 



# Example: Lasso regression

### Example

**Data**: Let  $\{(y_i, x_i)_{1 \le i \le n}, y_i \in \mathbb{R}, x_i \in \mathbb{R}^d\}$ .

Find  $\beta$  to minimize :

$$f(\beta) = \frac{1}{n} ||X\beta - y||_2^2 + \lambda ||\beta||_1, \quad ||\beta||_1 = \sum_{i=1}^n |\beta_i|.$$

$$\beta \text{ optimal if } 0 \in \frac{\partial}{\partial x_i} (\frac{1}{2} ||X\beta - y||_2^2 + \lambda ||\beta||_1).$$

$$eta$$
 optimal if  $0 \in \frac{\partial}{\partial eta}(\frac{1}{n}\|Xeta - y\|_2^2 + \lambda\|eta\|_1)$  
$$0 \in \frac{2}{n}X^{\top}(Xeta - y) + \lambda\frac{\partial}{\partial eta}(\|eta\|_1) = 0$$

# Example: Lasso regression

### Example

Data: Let  $\{(y_i, x_i)_{1 \le i \le n}, y_i \in \mathbb{R}, x_i \in \mathbb{R}^d\}$ .

Find  $\beta$  to minimize :

$$\begin{split} f(\beta) &= \frac{1}{n} \|X\beta - y\|_2^2 + \lambda \|\beta\|_1, \quad \|\beta\|_1 = \sum_{i=1}^n |\beta_i|. \\ \beta \text{ optimal if } 0 &\in \frac{\partial}{\partial \beta} (\frac{1}{n} \|X\beta - y\|_2^2 + \lambda \|\beta\|_1) \\ 0 &\in \frac{2}{n} X^\top (X\beta - y) + \lambda \frac{\partial}{\partial \beta} (\|\beta\|_1) = 0 \end{split}$$

$$\hat{\beta}$$
 is optimal if  $\exists z$  with  $z_i = \begin{cases} sign(\hat{\beta}_i), & \text{if } \hat{\beta}_i \neq 0, \\ [-1,1], & \text{if } \hat{\beta}_i = 0 \end{cases}$  and  $0 = \frac{2}{\pi} X^{\top} (X\beta - y) + \lambda z.$ 

# Smoothness: "Self tuning" property

#### Definition

We say that  $f: \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$  is **smooth** with constant L if

$$\|\nabla f(x) - \nabla f(y)\|_2 \le \|x - y\|_2, \quad \forall \ x, y \in \mathbb{R}^d$$
 (15)

If a twice differentiable  $f: \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$  is *L*-smooth then

"Self tuning property":  $\|\nabla f(x)\| \to 0$ , as  $x \to x^*$ .



Fact: f is L-smooth iff  $g(x) = \frac{L}{2} ||x||_2^2 - f(x)$  is convex.

 $\mathsf{recall}: \mathsf{monotone} \Leftrightarrow \mathsf{convex}, \mathsf{so} \mathsf{\ it\ is\ enough\ to\ show\ that\ } \nabla g \mathsf{\ is\ monotone}.$ 

#### Proof.

 $\Rightarrow$ ) we want to proof that  $\langle \nabla g(x) - \nabla g(y), x - y \rangle \ge 0$ .



Fact: f is L-smooth iff  $g(x) = \frac{L}{2}||x||_2^2 - f(x)$  is convex. recall : monotone  $\Leftrightarrow$  convex, so it is enough to show that  $\nabla g$  is monotone.

#### Proof.

 $\Rightarrow$ ) we want to proof that  $\langle \nabla g(x) - \nabla g(y), x - y \rangle \geq 0$ .

$$\langle \nabla g(x) - \nabla g(y), x - y \rangle = L \|x - y\|_{2}^{2} - \langle \nabla f(x) - \nabla f(y), x - y \rangle$$

$$\geq L \|x - y\|_{2}^{2} - \|\nabla f(x) - \nabla f(y)\|_{2} \cdot \|x - y\|_{2}, \text{ (csi)}$$

$$\geq L \|x - y\|_{2}^{2} - L \|x - y\|_{2}^{2} = 0. \text{ (smoothness)}$$

It follows that  $g(x) = \frac{L}{2}||x||_2^2 - f(x)$  is convex.



Fact: f is L-smooth iff  $g(x) = \frac{L}{2} ||x||_2^2 - f(x)$  is convex.

#### Proof.

 $\Leftarrow$ ) since g is convex we have using def 2:

$$g(y) \ge g(x) + \langle \nabla g(x), y - x \rangle$$



**Fact**: f is L-smooth iff  $g(x) = \frac{L}{2}||x||_2^2 - f(x)$  is convex.

#### Proof.

 $\Leftarrow$ ) since g is convex we have using def 2:

$$g(y) \ge g(x) + \langle \nabla g(x), y - x \rangle$$

$$\Rightarrow \frac{L}{2} \|y\|_{2}^{2} - f(y) \ge \frac{L}{2} \|x\|_{2}^{2} - f(x) + \langle Lx - \nabla f(x), y - x \rangle$$
$$f(y) \le f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|_{2}^{2}$$

If f is twice differentiable, then g is also twice differentiable and g convex  $\Leftrightarrow \nabla^2 g(x) \succcurlyeq 0 \Leftrightarrow L \cdot I_d - \nabla^2 f(x) \succcurlyeq 0 \Leftrightarrow \nabla^2 f(x) \preccurlyeq L \cdot I_d$  AIMS

# Smoothness: Examples

### Example

- Convex quadratics  $x \mapsto x^{\top}Ax + b^{\top}x + c$  is smooth for any  $L \ge 2\lambda_{\max}(A)$ .
- 2 Logistic:  $x \mapsto log(1 + e^{-y\langle a, x \rangle})$
- **3** Trigonometric :  $x \mapsto sin(x), cos(x)$ .

#### Proof.

Exercise!

**Exercise**: Using the fact that  $\sigma_{\max}^2(X) \|d\|_2^2 \ge \|X^\top d\|_2^2$ , show that

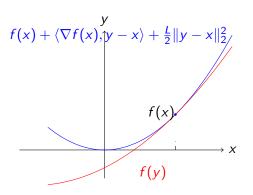
$$f(w) = \frac{1}{2} ||X^{\top}w - b||_2^2$$

is  $\sigma_{\max}^2(X)$ -smooth.



### Smoothness: convex example

$$f(y) = \|y\|_2^2$$

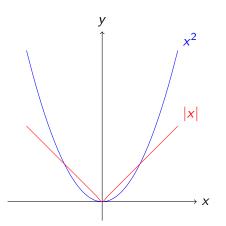


$$f(y) \le f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} ||y - x||_2^2, \quad \forall \ x, y \in \mathbb{R}^d$$



### Smoothness: convex counter-example

$$f(y) = \|y\|_1$$



False: 
$$f(y) \le f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} ||y - x||_2^2, \quad \forall x, y \in \mathbb{R}^d$$

### Strong Convexity

#### Definition

A function f is **strongly convex** with parameter  $\mu$  if  $g(x) = f(x) - \frac{\mu}{2} ||x||_2^2$  is convex.

A function f is convex if it is strongly convex with parameter  $\mu = 0$ .

$$g \text{ convex} \Rightarrow g(y) \ge g(x) + \langle \nabla g(x), y - x \rangle$$

$$\Rightarrow f(y) \ge f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} ||y - x||_2^2$$
(16)

$$f$$
 is  $\mu$ -strongly convex  $\Rightarrow g(x) = f(x) - \frac{\mu}{2} ||x||_2^2$  is convex

If f is twice differentiable, then g is also twice differentiable and g convex  $\Leftrightarrow \nabla^2 g(x) \succcurlyeq 0 \Leftrightarrow \nabla^2 f(x) - \mu \cdot I_d \succcurlyeq 0 \Leftrightarrow \mu \cdot I_d \preccurlyeq \nabla^2 f(x)$ 



### Recap

#### Definition

f is  $\mu$ -strongly convex if

**1** 
$$f(y) \ge f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} ||y - x||_2^2$$

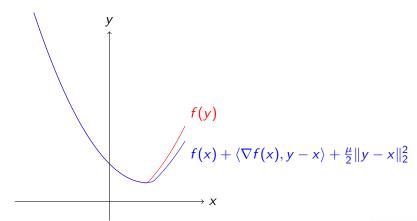
If f is  $\mu$ -strongly convex and L-smooth we have:

$$\mu \cdot I_d \preccurlyeq \nabla^2 f(x) \preccurlyeq L \cdot I_d \tag{17}$$



### Strong convexity example

Hinge Loss + L2 :  $f(y) = \max(0, 1 - y) + \frac{1}{2} ||y||_2^2$  is a strongly convex function.





# Strong convexity example

**1** Using the fact that  $\sigma_{\min}^2(X) \|d\|_2^2 \leq \|X^\top d\|_2^2$ , show that

$$f(w) = \frac{1}{2} \|X^{\top}w - b\|_2^2$$

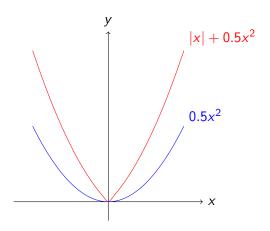
is  $\sigma_{\min}^2(X)$ -strongly convex.

- ②  $f(x) = x^{\top} Qx$  is  $\mu$ -strongly convex with  $\mu = 2\lambda_{\min}(Q)$
- $(x) = \frac{1}{2} ||x||_2^2 + ||x||_1$



# Example of strongly convex function and non-smooth

 $f(y) = \frac{1}{2} ||x||_2^2 + ||x||_1$  is a strongly convex function but non smooth.





# Acknowledgements

# Thanks for your attention!



