

## LEC 1

## 1.1 Prerequisites

Allowed coding languages: MATLAB, R, PYTHON, JULIA.

## 1.2 Optimization Problems

## 1.2.1 Unconstrained Optimization

$$\min_{\mathbf{x}} f(\mathbf{x})$$

where  $\mathbf{x}$  is the unknown *variable* (also called the *decision variable*) and

$$f : \mathbb{R}^n \rightarrow \mathbb{R}$$

is called the *objective function*.

## 1.2.2 Constrained Optimization

$$\min_{\mathbf{x} \in C} f(\mathbf{x})$$

where  $C \subseteq \mathbb{R}^n$  is called the *feasible region*.

Commonly,  $C$  is specified by inequality and equality constraints:

$$C = \left\{ \mathbf{x} \in \mathbb{R}^n : g_1(\mathbf{x}) \geq 0, \dots, g_k(\mathbf{x}) \geq 0; h_1(\mathbf{x}) = 0, \dots, h_j(\mathbf{x}) = 0 \right\}.$$

## 1.3 Supervised Learning

## 1.3.1 Definitions

**Training data.** We are given labeled examples, typically as pairs

$$(\mathbf{a}_1, \mathbf{y}_1), \dots, (\mathbf{a}_m, \mathbf{y}_m),$$

where

$$\begin{aligned} \mathbf{a}_i &\in \mathbb{R}^d, \quad i = 1, \dots, m \quad (\text{“feature vectors”}), \\ \mathbf{y}_i &\in \mathbb{R}^e, \quad i = 1, \dots, m \quad (\text{“labels”}). \end{aligned}$$

We seek a function  $\Phi$  such that

$$\Phi(\mathbf{a}_i) \approx \mathbf{y}_i, \quad i = 1, \dots, m.$$

**Why find such a function  $\Phi$ ?**

- Reason: To apply it to future *unseen* feature vectors  $\mathbf{a} \in \mathbb{R}^d$ .

### 1.3.2 Usual Optimization-Based Approach

Assume that  $\Phi$  comes from a parametric family:

$$\Phi(\mathbf{a}) \equiv f(\mathbf{a}, \mathbf{x}),$$

where  $f$  is a fixed known function, and  $\mathbf{x}$  is a vector of (say)  $p$  parameters, initially unknown.

**Loss function.** Define a *loss function*  $l$  to measure the discrepancy between the prediction  $f(\mathbf{a}_i, \mathbf{x})$  and the label  $\mathbf{y}_i$ .

#### Example

**Least squares.** A common choice is the squared  $\ell_2$  loss:

$$l(\mathbf{a}, \mathbf{y}, \mathbf{x}) = \frac{1}{2} \|f(\mathbf{a}, \mathbf{x}) - \mathbf{y}\|^2.$$

Let the training set be

$$D = \{(\mathbf{a}_1, \mathbf{y}_1), \dots, (\mathbf{a}_m, \mathbf{y}_m)\}.$$

Define the *empirical loss*:

$$L_D(\mathbf{x}) = \sum_{i=1}^m l(\mathbf{a}_i, \mathbf{y}_i, \mathbf{x}).$$

Then the learning problem becomes the optimization problem

$$\min_{\mathbf{x}} L_D(\mathbf{x}).$$

#### How to find an optimal $\mathbf{x}$ ?

This is the role of optimization *algorithms*, which will be introduced later in the course.

#### Generalization and test set.

- Question: How do we know that  $f(\mathbf{a}, \mathbf{x})$  is a good predictor for  $\mathbf{y}$  when  $\mathbf{a} \notin \{\mathbf{a}_1, \dots, \mathbf{a}_m\}$ ?

Hold aside additional labeled data items

$$(\mathbf{a}_{m+1}, \mathbf{y}_{m+1}), \dots, (\mathbf{a}_{m'}, \mathbf{y}_{m'}),$$

called the *test set*. We evaluate the learned parameter vector  $\mathbf{x}$  by computing the *test loss* on this set:

$$\sum_{i=m+1}^{m'} l(\mathbf{a}_i, \mathbf{y}_i, \mathbf{x}).$$

## 1.4 Unsupervised Learning

### 1.4.1 Clustering

**Classic version: Clustering.**

Given data points

$$\mathbf{a}_1, \dots, \mathbf{a}_m \in \mathbb{R}^d,$$

we want to partition them into  $k$  clusters. This can be expressed by selecting labels

$$y_1, \dots, y_m \quad \text{such that} \quad y_i \in \{1, \dots, k\}, \quad \forall i = 1, \dots, m,$$

and enforcing that

$$\|\mathbf{a}_i - \mathbf{a}_j\| \quad \text{is} \quad \begin{cases} \text{small,} & y_i = y_j, \\ \text{large,} & y_i \neq y_j. \end{cases}$$

### 1.4.2 K-means Objective Function

Introduce additional variables (cluster centers)

$$\mathbf{c}_1, \dots, \mathbf{c}_k \in \mathbb{R}^d.$$

#### Example

**K-means objective.** The K-means clustering problem can be written as

$$\min_{\{y_1, \dots, y_m\}} \min_{\{\mathbf{c}_1, \dots, \mathbf{c}_k\}} \sum_{i=1}^m \|\mathbf{a}_i - \mathbf{c}_{y_i}\|^2.$$

- This is a *partially discrete* optimization problem (e.g., K-means): the labels  $y_i$  are discrete, while the centers  $\mathbf{c}_j$  are continuous.
- In contrast, an example of a purely *continuous* optimization problem is linear least squares.

## 1.5 Linear Least Squares

Consider data points

$$(\mathbf{a}_1, y_1), \dots, (\mathbf{a}_m, y_m), \quad \mathbf{a}_i \in \mathbb{R}^n, \quad y_i \in \mathbb{R}.$$

**Linear model hypothesis.** Assume there exists  $\mathbf{x} \in \mathbb{R}^n$  such that

$$y_i \approx \mathbf{a}_i^\top \mathbf{x}, \quad \forall i = 1, 2, \dots, m.$$

**Inner product and transpose.**

For  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ ,

$$\mathbf{a}^\top \mathbf{b} = a_1 b_1 + \dots + a_n b_n$$

is the inner product. The superscript  $^\top$  denotes transpose.

#### Example

**Least-squares fit as supervised learning.**

In the earlier supervised-learning notation, we have

$$\Phi(\mathbf{a}) \equiv f(\mathbf{a}, \mathbf{x}) \equiv \mathbf{a}^\top \mathbf{x}.$$

Use the loss

$$l(\mathbf{a}, \mathbf{y}, \mathbf{x}) = \frac{1}{2}(\mathbf{a}^\top \mathbf{x} - \mathbf{y})^2,$$

so that

$$L_D(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{x} - y_i)^2.$$

Now define the data matrix

$$\mathbf{A} = \begin{pmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_m^\top \end{pmatrix} \in \mathbb{R}^{m \times n}.$$

Usually, vectors are treated as column vectors, which means  $\mathbf{x}^\top$  is the corresponding row vector with the same entries.

#### Fact

**Least squares in matrix form.**

$$\frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 = \frac{1}{2} \sum_{i=1}^m (\mathbf{A}\mathbf{x} - \mathbf{y})_i^2 = \frac{1}{2} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{x} - y_i)^2 = L_D(\mathbf{x}),$$

where  $\mathbf{A}$  is an  $m \times n$  matrix,  $\mathbf{x}$  is  $n$ -dimensional, and  $\mathbf{y}$  is  $m$ -dimensional.

## 1.6 Linear Algebra Review

Vectors are (by default) *column* vectors in  $\mathbb{R}^n$ .

**Special vectors.**

$\mathbf{0}$                       all-zero vector,  
 $\mathbf{e}$                         all-one vector,  
 $\mathbf{e}_1, \dots, \mathbf{e}_n$         standard basis vectors.

### 1.6.1 Linear Dependence

Say  $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathbb{R}^n$  are **dependent** if there exist scalars  $\alpha_1, \dots, \alpha_k \in \mathbb{R}$ , not all 0, such that

$$\alpha_1 \mathbf{v}_1 + \dots + \alpha_k \mathbf{v}_k = \mathbf{0}.$$

The left-hand side is called a *linear combination* of  $\mathbf{v}_1, \dots, \mathbf{v}_k$ .

If the *only* solution is  $\alpha_1 = \dots = \alpha_k = 0$ , then  $\mathbf{v}_1, \dots, \mathbf{v}_k$  are **independent**.

### 1.6.2 Subspaces and Span

A set  $S \subseteq \mathbb{R}^n$  is a **subspace** if:

$$\mathbf{0} \in S,$$

and for all  $\mathbf{x}, \mathbf{y} \in S$  and all  $\alpha, \beta \in \mathbb{R}$ ,

$$\alpha \mathbf{x} + \beta \mathbf{y} \in S.$$

Two extreme examples of subspaces are  $\{\mathbf{0}\}$  and  $\mathbb{R}^n$ .

Given  $k$  vectors  $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathbb{R}^n$ , the set of all linear combinations is their *span*:

$$\text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}.$$

#### Fact

- A span is always a subspace.
- If  $S \subseteq \mathbb{R}^n$  is a subspace, then there exist  $\mathbf{v}_1, \dots, \mathbf{v}_k \in S$  such that

$$S = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\},$$

and  $\mathbf{v}_1, \dots, \mathbf{v}_k$  are independent. The set  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  is called a *basis* of  $S$ , where  $k = \dim(S)$ .

Note that

$$\text{span}\{\} = \{\mathbf{0}\}.$$

- For a subspace  $S$ , all bases of  $S$  have the same cardinality. This cardinality is called the *dimension* of  $S$  and is denoted  $\dim(S)$ .
- If  $S, T \subseteq \mathbb{R}^n$  are subspaces with  $S \subseteq T$ , then

$$\dim(S) \leq \dim(T).$$

If  $\dim(S) = \dim(T)$ , then  $S = T$ .

Furthermore, if  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  is a basis of  $S$  with  $k = \dim(S)$ , then there exist vectors

$$\mathbf{v}_{k+1}, \dots, \mathbf{v}_{\dim(T)} \in T$$

such that  $\{\mathbf{v}_1, \dots, \mathbf{v}_{\dim(T)}\}$  is a basis of  $T$ . This is the *Basis Extension Theorem*.

### 1.6.3 Special Matrices

**Special matrices.**

- $\mathbf{0}$ : all-zero matrix.
- $\mathbf{I}$ : identity matrix.

## LEC 2

## 2.1 Matrix Concepts

## 2.1.1 Matrix–Vector Product

Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{x} \in \mathbb{R}^n$ . Define  $\mathbf{y} \in \mathbb{R}^m$  by

$$y_i = \sum_{j=1}^n A(i, j) x_j, \quad i = 1, \dots, m.$$

Then  $\mathbf{y}$  is called the *matrix–vector product* of  $\mathbf{A}$  and  $\mathbf{x}$ , and we write

$$\mathbf{y} = \mathbf{A}\mathbf{x}.$$

**Fact**

For all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  and all scalars  $\alpha \in \mathbb{R}$ :

- $\mathbf{A}(\mathbf{x} + \mathbf{y}) = \mathbf{A}\mathbf{x} + \mathbf{A}\mathbf{y}$ ,
- $\mathbf{A}(\alpha\mathbf{x}) = \alpha(\mathbf{A}\mathbf{x})$ .

A mapping  $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is called a **linear transformation** if and only if it satisfies

$$T(\mathbf{x} + \mathbf{y}) = T(\mathbf{x}) + T(\mathbf{y}), \quad T(\alpha\mathbf{x}) = \alpha T(\mathbf{x})$$

for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  and all  $\alpha \in \mathbb{R}$ .

**Fact**

For any linear transformation  $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , there exists a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  such that

$$T(\mathbf{x}) \equiv \mathbf{A}\mathbf{x}, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

**Interpretations of  $\mathbf{y} = \mathbf{A}\mathbf{x}$ :**

$$\mathbf{y} = \mathbf{A}\mathbf{x} \iff \mathbf{y} = \mathbf{A}(:, 1)x_1 + \dots + \mathbf{A}(:, n)x_n \iff \mathbf{y} = \begin{pmatrix} \mathbf{A}(1, :) \mathbf{x} \\ \vdots \\ \mathbf{A}(m, :) \mathbf{x} \end{pmatrix},$$

where  $\mathbf{A}(:, j)$  denotes the  $j$ -th column of  $\mathbf{A}$ , and  $\mathbf{A}(i, :)$  denotes the  $i$ -th row.

**Example**

**Identity matrix.** For the identity matrix  $\mathbf{I} \in \mathbb{R}^{n \times n}$ ,

$$\mathbf{I}\mathbf{x} = \mathbf{x}, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

### 2.1.2 Transpose

For  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , its **transpose**  $\mathbf{A}^\top \in \mathbb{R}^{n \times m}$  is defined by

$$\mathbf{A}^\top(j, i) = A(i, j), \quad \forall i = 1, \dots, m, \quad \forall j = 1, \dots, n.$$

For a vector  $\mathbf{x} \in \mathbb{R}^n$ , the transpose  $\mathbf{x}^\top$  is a row vector with the same entries as  $\mathbf{x}$ .

### 2.1.3 Inner Product

For  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , the **inner product** is

$$\mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i.$$

### 2.1.4 Range and Nullspace

For any  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , define:

- (i)  $\text{Range}(\mathbf{A}) = \{\mathbf{Ax} : \mathbf{x} \in \mathbb{R}^n\} \subseteq \mathbb{R}^m$ ,
- (ii)  $\text{Null}(\mathbf{A}) = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{Ay} = \mathbf{0}\} \subseteq \mathbb{R}^n$ .

#### Fact

$\text{Range}(\mathbf{A})$  and  $\text{Null}(\mathbf{A})$  are subspaces.

Intuitively:

- $\text{Range}(\mathbf{A})$  is the *output space* (all vectors that can appear as  $\mathbf{Ax}$ ),
- $\text{Null}(\mathbf{A})$  is the space of vectors that *disappear* under the transformation (they are mapped to  $\mathbf{0}$ ).

The **rank** of  $\mathbf{A}$  is

$$\text{rank}(\mathbf{A}) = \dim(\text{Range}(\mathbf{A})).$$

#### Fact

$$\text{rank}(\mathbf{A}) \leq \min(m, n).$$

We can express the range as

$$\text{Range}(\mathbf{A}) = \text{span}\{\mathbf{A}(:, 1), \dots, \mathbf{A}(:, n)\}.$$

#### Fact

$$\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^\top).$$

**Example**

If  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is **diagonal** (i.e.,  $A(i, j) = 0$  whenever  $i \neq j$ ), then

$$\text{rank}(\mathbf{A}) = \text{number of nonzero diagonal entries of } \mathbf{A}.$$

**2.1.5 Matrix Multiplication**

Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{B} \in \mathbb{R}^{n \times p}$ . Their product  $\mathbf{C} = \mathbf{AB} \in \mathbb{R}^{m \times p}$  is defined by

$$C(i, j) = \sum_{k=1}^n A(i, k) B(k, j), \quad \forall i = 1, \dots, m, \quad \forall j = 1, \dots, p.$$

**Fact**

Matrix multiplication corresponds to composition of linear transformations:

$$\forall \mathbf{x} \in \mathbb{R}^p, \quad \mathbf{A}(\mathbf{B}\mathbf{x}) = (\mathbf{AB})\mathbf{x}.$$

**Fact**

Matrix multiplication is associative:

$$\mathbf{A}(\mathbf{BC}) = (\mathbf{AB})\mathbf{C}.$$

In general, matrix multiplication is *not* commutative:

$$\mathbf{AB} \neq \mathbf{BA}.$$

**2.1.6 Square Matrices and Invertibility****Fact**

Given  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , the following are equivalent:

- (i) There exists a matrix  $\mathbf{A}^{-1} \in \mathbb{R}^{n \times n}$  such that

$$\mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$$

(then  $\mathbf{A}^{-1}$  is called the *inverse* of  $\mathbf{A}$ ).

- (ii)  $\text{rank}(\mathbf{A}) = n \iff \text{Range}(\mathbf{A}) = \mathbb{R}^n$ .

- (iii)  $\text{Null}(\mathbf{A}) = \{\mathbf{0}\}$ .

- (iv) The columns of  $\mathbf{A}$  are linearly independent.

If any (and hence all) of these hold, we say  $\mathbf{A}$  is **invertible** or **nonsingular**.



### 2.1.7 Properties of the Inverse

#### Fact

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$ .

(i) If  $\mathbf{A}$  is invertible, then  $\mathbf{A}^{-1}$  is uniquely determined and

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}.$$

(ii) If  $\mathbf{A}$  is invertible, then

$$(\mathbf{A}^{-1})^{-1} = \mathbf{A}.$$

(iii) If  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$  are invertible, then  $\mathbf{AB}$  is invertible and

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}.$$

Similarly for transposes:

$$(\mathbf{AB})^{\top} = \mathbf{B}^{\top}\mathbf{A}^{\top}.$$

(iv) If  $\mathbf{A}$  is invertible, then  $\mathbf{A}^{\top}$  is invertible and

$$(\mathbf{A}^{-1})^{\top} = (\mathbf{A}^{\top})^{-1}.$$

Notation:

$$\mathbf{A}^{-\top} := (\mathbf{A}^{-1})^{\top}.$$

## 2.2 Quadratic Functions

### 2.2.1 Definition of Quadratic Functions

A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is called **quadratic** if it is a sum of terms, each of which is

- constant (e.g. 7),
- linear (e.g.  $b_i x_i$ ),
- or quadratic (e.g.  $7x_i^2$  or  $8x_i x_j$ ).

#### Fact

If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is quadratic, then it can be written in the **standard form**

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^{\top} \mathbf{H} \mathbf{x} + \mathbf{g}^{\top} \mathbf{x} + d,$$

where

- $\mathbf{H} \in \mathbb{S}^n$  is a **symmetric** matrix:

$$\mathbf{H} = \mathbf{H}^{\top}, \quad H(i, j) = H(j, i), \quad i, j = 1, \dots, n,$$

- $\mathbf{g} \in \mathbb{R}^n$  is a vector;
- $d \in \mathbb{R}$  is a scalar.

## Example

Let

$$f(\mathbf{x}) = 2x_1^2 - 7x_1x_2 - 3x_1 - 9.$$

Then we can write

$$f(\mathbf{x}) = \frac{1}{2} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}^\top \begin{pmatrix} 4 & -7 \\ -7 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} -3 \\ 0 \end{pmatrix}^\top \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - 9.$$

Note that  $\frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x}$  is a scalar.

## 2.2.2 Positive (Semi)definite Matrices

Let  $\mathbf{A} \in \mathbb{S}^n$  (symmetric).

- $\mathbf{A}$  is **positive semidefinite** (psd) if

$$\forall \mathbf{x} \in \mathbb{R}^n, \quad \mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0.$$

We write  $\mathbf{A} \succeq \mathbf{0}$ .

- $\mathbf{A}$  is **positive definite** (pd) if

$$\forall \mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}, \quad \mathbf{x}^\top \mathbf{A} \mathbf{x} > 0.$$

We write  $\mathbf{A} \succ \mathbf{0}$ .

## Example

$$\mathbf{A} \text{ pd} \Rightarrow \mathbf{A} \text{ psd},$$

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \text{ is pd},$$

$$\begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix} \text{ is psd, not pd},$$

$$\begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} \text{ is pd},$$

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \text{ is psd, not pd},$$

$$\begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix} \text{ is not psd},$$

$$\begin{pmatrix} 1 & -1 \\ -1 & 2 \end{pmatrix} \text{ is pd}.$$

## 2.3 Theorem on Quadratic Minimization

### Theorem

Consider the quadratic function

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} + \mathbf{g}^\top \mathbf{x} + d,$$

where  $\mathbf{H} \in \mathbb{S}^n$ . Then  $f$  is **bounded below** and has a **minimizer** if and only if

- $\mathbf{H}$  is positive semidefinite, and
- $\mathbf{g} \in \text{Range}(\mathbf{H})$ .

In this case, any solution to

$$\mathbf{H} \mathbf{x} = -\mathbf{g}$$

is a minimizer of  $f$ .

If  $\mathbf{H}$  is positive definite, then all the above conditions hold automatically:

$$\mathbf{H} \text{ pd} \Rightarrow \mathbf{H} \text{ psd},$$

$$\mathbf{H} \text{ pd} \Rightarrow \mathbf{H} \text{ nonsingular} \Rightarrow \text{Range}(\mathbf{H}) = \mathbb{R}^n \Rightarrow \mathbf{g} \in \text{Range}(\mathbf{H}).$$

### Why a positive definite matrix is nonsingular?

If  $\mathbf{H}$  were singular, then there would exist  $\mathbf{x} \in \text{Null}(\mathbf{H}) \setminus \{\mathbf{0}\}$  with  $\mathbf{H} \mathbf{x} = \mathbf{0}$ . Then

$$\mathbf{x}^\top \mathbf{H} \mathbf{x} = \mathbf{x}^\top \mathbf{0} = 0,$$

contradicting the condition for positive definiteness  $\mathbf{x}^\top \mathbf{H} \mathbf{x} > 0$  for all  $\mathbf{x} \neq \mathbf{0}$ .

### Fact

If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuous and  $f$  has a minimizer, then  $f$  is bounded below. (The converse is not always true.)

### 2.3.1 Proof of Theorem (Forward Direction)

#### Proof

Assume  $\mathbf{H}$  is positive semidefinite and  $\mathbf{g} \in \text{Range}(\mathbf{H})$ . Then  $-\mathbf{g} \in \text{Range}(\mathbf{H})$ , so there exists  $\mathbf{w} \in \mathbb{R}^n$  such that

$$\mathbf{H} \mathbf{w} = -\mathbf{g}.$$

Let  $\mathbf{x} \in \mathbb{R}^n$  be arbitrary. Then

$$\begin{aligned} f(\mathbf{x}) &= \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} + \mathbf{g}^\top \mathbf{x} + d \\ &= \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} - \mathbf{w}^\top \mathbf{H} \mathbf{x} + d \quad (\text{since } \mathbf{g} = -\mathbf{H} \mathbf{w}) \\ &= \frac{1}{2} (\mathbf{x} - \mathbf{w})^\top \mathbf{H} (\mathbf{x} - \mathbf{w}) - \frac{1}{2} \mathbf{w}^\top \mathbf{H} \mathbf{w} + d \quad (\text{completing the square}). \end{aligned}$$

Note that  $\mathbf{w}^\top \mathbf{H} \mathbf{x} = \mathbf{x}^\top \mathbf{H} \mathbf{w}$  since both are scalars.  
Because  $\mathbf{H}$  is positive semidefinite,

$$\frac{1}{2}(\mathbf{x} - \mathbf{w})^\top \mathbf{H}(\mathbf{x} - \mathbf{w}) \geq 0 \quad \text{for all } \mathbf{x},$$

and the minimum value 0 is attained at  $\mathbf{x} = \mathbf{w}$ .

The remaining term  $-\frac{1}{2} \mathbf{w}^\top \mathbf{H} \mathbf{w} + d$  is constant, so  $f$  is bounded below and attains its minimum at  $\mathbf{x} = \mathbf{w}$ .

### 2.3.2 Tools for the Converse Direction (Unfinished)

Let  $S \subseteq \mathbb{R}^n$  be a subspace. Its **orthogonal complement** is

$$S^\perp = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x}^\top \mathbf{y} = 0, \forall \mathbf{y} \in S\}.$$

If  $\mathbf{x}^\top \mathbf{y} = 0$ , we say that  $\mathbf{x}$  and  $\mathbf{y}$  are **orthogonal**.

#### Fact

For a subspace  $S \subseteq \mathbb{R}^n$ :

- $S^\perp$  is a subspace.
- $(S^\perp)^\perp = S$ .
- $\dim(S) + \dim(S^\perp) = n$ .
- For every  $\mathbf{w} \in \mathbb{R}^n$ , there exist unique  $\mathbf{x} \in S$  and  $\mathbf{y} \in S^\perp$  such that

$$\mathbf{w} = \mathbf{x} + \mathbf{y}.$$

#### Fact

**Fundamental Theorem of Linear Algebra.** For all  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,

$$\text{Range}(\mathbf{A})^\perp = \text{Null}(\mathbf{A}^\top).$$

**Contrapositive statement (for the theorem).**

If  $\mathbf{H}$  is not positive semidefinite, *or*  $\mathbf{H}$  is positive semidefinite and  $\mathbf{g} \notin \text{Range}(\mathbf{H})$ , then the quadratic function

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} + \mathbf{g}^\top \mathbf{x} + d$$

is **unbounded below**. (Proof of this converse direction is finished in next lecture.)

## LEC 3

## 3.1 Quadratic Functions

## 3.1.1 The Minimizer

Given  $\mathbf{H} \in \mathbb{S}^n$ ,  $\mathbf{g} \in \mathbb{R}^n$ ,  $d \in \mathbb{R}$ , consider

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} + \mathbf{g}^\top \mathbf{x} + d.$$

**Theorem**

The function  $f$  is bounded below and has a minimizer if and only if:

- a)  $\mathbf{H}$  is positive semidefinite, and
- b)  $\mathbf{g} \in \text{Range}(\mathbf{H})$ .

In this case, any solution of

$$\mathbf{H} \mathbf{x} = -\mathbf{g}$$

is a minimizer.

**Fact**

**Contrapositive Theorem.** If  $\mathbf{H}$  is not positive semidefinite, or  $\mathbf{H}$  is positive semidefinite and  $\mathbf{g} \notin \text{Range}(\mathbf{H})$ , then  $f$  is unbounded below.

## 3.1.2 Proof of the Contrapositive

**Proof**

**Case 1: (a) fails:  $\mathbf{H}$  is not positive semidefinite.**

Then there exists  $\mathbf{z} \in \mathbb{R}^n$  such that

$$\mathbf{z}^\top \mathbf{H} \mathbf{z} < 0.$$

Define the univariate function

$$q(t) = f(t\mathbf{z}) = \frac{1}{2} t^2 \mathbf{z}^\top \mathbf{H} \mathbf{z} + t \mathbf{g}^\top \mathbf{z} + d.$$

This is a quadratic function in  $t$  with *negative* leading coefficient  $\frac{1}{2} \mathbf{z}^\top \mathbf{H} \mathbf{z} < 0$ , so  $q(t)$  attains arbitrarily negative values as  $t \rightarrow \pm\infty$ . Therefore  $f$  has no lower bound and has no minimizer.

**Case 2: (a) holds but (b) fails:  $\mathbf{H}$  is positive semidefinite and  $\mathbf{g} \notin \text{Range}(\mathbf{H})$ .**

By the Fundamental Theorem of Linear Algebra for symmetric  $\mathbf{H}$ ,

$$\mathbb{R}^n = \text{Range}(\mathbf{H}) \oplus \text{Null}(\mathbf{H}),$$

with  $\text{Range}(\mathbf{H}) \perp \text{Null}(\mathbf{H})$ . Thus there exist

$$\mathbf{g}_1 \in \text{Range}(\mathbf{H}), \quad \mathbf{g}_2 \in \text{Null}(\mathbf{H})$$

such that

$$\mathbf{g} = \mathbf{g}_1 + \mathbf{g}_2, \quad \mathbf{g}_1^\top \mathbf{g}_2 = 0.$$

Define

$$q(t) := f(t\mathbf{g}_2).$$

Since  $\mathbf{H}\mathbf{g}_2 = \mathbf{0}$  (because  $\mathbf{g}_2 \in \text{Null}(\mathbf{H})$ ), we have

$$\begin{aligned} q(t) &= \frac{1}{2}t^2\mathbf{g}_2^\top \mathbf{H}\mathbf{g}_2 + t\mathbf{g}^\top \mathbf{g}_2 + d \\ &= t\mathbf{g}^\top \mathbf{g}_2 + d. \end{aligned}$$

Next,

$$\mathbf{g}^\top \mathbf{g}_2 = (\mathbf{g}_1 + \mathbf{g}_2)^\top \mathbf{g}_2 = \mathbf{g}_2^\top \mathbf{g}_2 > 0,$$

because  $\mathbf{g} \notin \text{Range}(\mathbf{H}) \Rightarrow \mathbf{g} \neq \mathbf{g}_1 \Rightarrow \mathbf{g}_2 \neq \mathbf{0}$ .

Hence  $q(t)$  is a linear function of  $t$  with positive slope. Therefore

$$q(t) \rightarrow -\infty \quad \text{as} \quad t \rightarrow -\infty,$$

so  $f$  is unbounded below.

## 3.2 Linear Least Squares

### 3.2.1 Norms

For  $\mathbf{x} \in \mathbb{R}^n$ , the Euclidean norm (or 2-norm) is

$$\|\mathbf{x}\| = \sqrt{x(1)^2 + \cdots + x(n)^2}.$$

A **norm**  $\|\cdot\|$  on  $\mathbb{R}^n$  satisfies, for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  and all  $\lambda \in \mathbb{R}$ :

(i) Positivity:

$$\|\mathbf{x}\| > 0 \text{ for } \mathbf{x} \neq \mathbf{0}, \quad \text{and} \quad \|\mathbf{x}\| = 0 \Rightarrow \mathbf{x} = \mathbf{0}.$$

(ii) Positive 1-homogeneity:

$$\|\lambda\mathbf{x}\| = |\lambda| \cdot \|\mathbf{x}\|.$$

(iii) Triangle inequality:

$$\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|.$$

### 3.2.2 Linear Least Squares (LLS)

Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{y} \in \mathbb{R}^m$ , we want  $\mathbf{x} \in \mathbb{R}^n$  to solve

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{y}\|.$$

This is equivalent to

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{y}\|^2 \quad \text{or} \quad \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|^2.$$

Compute:

$$\begin{aligned}\frac{1}{2}\|\mathbf{Ax} - \mathbf{y}\|^2 &= \frac{1}{2}(\mathbf{Ax} - \mathbf{y})^\top (\mathbf{Ax} - \mathbf{y}) \\ &= \frac{1}{2}(\mathbf{x}^\top \mathbf{A}^\top - \mathbf{y}^\top)(\mathbf{Ax} - \mathbf{y}) \\ &= \frac{1}{2}\mathbf{x}^\top \mathbf{A}^\top \mathbf{Ax} - \frac{1}{2}\mathbf{y}^\top \mathbf{Ax} - \frac{1}{2}\mathbf{x}^\top \mathbf{A}^\top \mathbf{y} + \frac{1}{2}\mathbf{y}^\top \mathbf{y} \\ &= \frac{1}{2}\mathbf{x}^\top \mathbf{A}^\top \mathbf{Ax} - \mathbf{y}^\top \mathbf{Ax} + \frac{1}{2}\mathbf{y}^\top \mathbf{y},\end{aligned}$$

since  $\mathbf{x}^\top \mathbf{A}^\top \mathbf{y} = \mathbf{y}^\top \mathbf{Ax}$  (both scalars).

This is a quadratic function in  $\mathbf{x}$  with

$$\mathbf{H} := \mathbf{A}^\top \mathbf{A}, \quad \mathbf{g} := -\mathbf{A}^\top \mathbf{y}, \quad d := \frac{1}{2}\mathbf{y}^\top \mathbf{y}.$$

#### Theorem

For all  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , the matrix  $\mathbf{A}^\top \mathbf{A}$  is positive semidefinite.

#### Proof

$$\mathbf{x}^\top \mathbf{A}^\top \mathbf{Ax} = (\mathbf{Ax})^\top \mathbf{Ax} = \|\mathbf{Ax}\|^2 \geq 0, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

#### Theorem

For all  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and all  $\mathbf{y} \in \mathbb{R}^m$ ,

$$\mathbf{A}^\top \mathbf{y} \in \text{Range}(\mathbf{A}^\top \mathbf{A}).$$

#### Proof

(Idea.) Otherwise, by the quadratic theorem for

$$\frac{1}{2}\|\mathbf{Ax} - \mathbf{y}\|^2 = \frac{1}{2}\mathbf{x}^\top \mathbf{A}^\top \mathbf{Ax} - \mathbf{y}^\top \mathbf{Ax} + \frac{1}{2}\mathbf{y}^\top \mathbf{y},$$

the objective would be unbounded below, contradicting the fact that squared norms are always nonnegative. Therefore a minimizer exists and must satisfy

$$\mathbf{A}^\top \mathbf{Ax} = \mathbf{A}^\top \mathbf{y}.$$

This system is called the **system of normal equations**. The solution is unique if  $\mathbf{A}^\top \mathbf{A}$  is positive definite.

#### Theorem

For  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , the following are equivalent:

$$\mathbf{A}^\top \mathbf{A} \text{ is positive definite} \iff \mathbf{A}^\top \mathbf{A} \text{ is nonsingular} \iff \text{rank}(\mathbf{A}) = n.$$

**Proof**

(Sketch.) If  $\text{rank}(\mathbf{A}) = n$ , the columns of  $\mathbf{A}$  are independent, so

$$\mathbf{A}\mathbf{x} \neq \mathbf{0} \quad \forall \mathbf{x} \neq \mathbf{0}.$$

Hence

$$\mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} = \|\mathbf{A}\mathbf{x}\|^2 > 0 \quad \forall \mathbf{x} \neq \mathbf{0},$$

so  $\mathbf{A}^\top \mathbf{A}$  is positive definite and thus nonsingular. The other implications follow from standard linear algebra facts about positive definite and nonsingular matrices.

### 3.3 Orthogonal Matrices

Recall: if  $\mathbf{x}^\top \mathbf{y} = 0$ , we say  $\mathbf{x}$  and  $\mathbf{y}$  are **orthogonal**.

If  $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$  is an orthonormal set with

$$\mathbf{x}_i^\top \mathbf{x}_j = \delta_{ij},$$

and  $\mathbf{U}$  is the matrix having these vectors as columns, then

$$\mathbf{U}^\top \mathbf{U} = \mathbf{I}, \quad \text{rank}(\mathbf{U}) = k \leq n.$$

If  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  is square and satisfies

$$\mathbf{Q}^\top \mathbf{Q} = \mathbf{I},$$

then  $\mathbf{Q}$  is an **orthogonal matrix**. In this case,

$$\mathbf{Q}^{-1} = \mathbf{Q}^\top,$$

and the rows of  $\mathbf{Q}$  are also orthonormal.

**Fact**

- If  $\mathbf{Q}$  is orthogonal, then  $\mathbf{Q}^\top$  is also orthogonal.
- If  $\mathbf{Q}_1, \mathbf{Q}_2$  are orthogonal, then  $\mathbf{Q}_1 \mathbf{Q}_2$  is orthogonal as well. (The set

$$O(n) = \{\mathbf{Q} \in \mathbb{R}^{n \times n} : \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}\}$$

is a *Lie group*.)

- For any orthogonal  $\mathbf{Q}$  and any  $\mathbf{x} \in \mathbb{R}^n$ ,

$$\|\mathbf{Q}\mathbf{x}\| = \|\mathbf{x}\|.$$

Orthogonal matrices represent rotations or reflections.



### 3.4 Eigendecomposition

#### 3.4.1 Spectral Theorem

##### Theorem

For any  $\mathbf{A} \in \mathbb{S}^n$ , there exists an orthogonal matrix  $\mathbf{Q}$  and a diagonal matrix  $\mathbf{D}$  such that

$$\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^\top.$$

The diagonal entries of  $\mathbf{D}$  are the eigenvalues of  $\mathbf{A}$  (unique up to ordering), and the columns of  $\mathbf{Q}$  are the corresponding eigenvectors.

Equivalently,

$$\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^\top \iff \mathbf{A}\mathbf{Q} = \mathbf{Q}\mathbf{D}.$$

#### 3.4.2 Characterizations of Positive Semidefiniteness

##### Theorem

For  $\mathbf{A} \in \mathbb{S}^n$ , the following are equivalent:

- (i)  $\mathbf{A}$  is positive semidefinite.
- (ii) All eigenvalues of  $\mathbf{A}$  are  $\geq 0$ .
- (iii) There exists a matrix  $\mathbf{G}$  such that  $\mathbf{A} = \mathbf{G}\mathbf{G}^\top$ .

#### 3.4.3 Characterizations of Positive Definiteness

##### Theorem

For  $\mathbf{A} \in \mathbb{S}^n$ , the following are equivalent:

- (i)  $\mathbf{A}$  is positive definite.
- (ii) All eigenvalues of  $\mathbf{A}$  are  $> 0$ .
- (iii) There exists a matrix  $\mathbf{G}$  such that  $\mathbf{A} = \mathbf{G}\mathbf{G}^\top$  and  $\text{rank}(\mathbf{G}) = n$ .

##### Notation.

- $\mathbf{A} \succeq \mathbf{B}$  means  $\mathbf{A} - \mathbf{B}$  is positive semidefinite.
- $\mathbf{A} \succ \mathbf{B}$  means  $\mathbf{A} - \mathbf{B}$  is positive definite.
- $\mathbb{S}_+^n$  denotes the set of positive semidefinite matrices.
- $\mathbb{S}_{++}^n$  denotes the set of positive definite matrices.

### 3.5 Topology Review

Point-set topology of  $\mathbb{R}^n$ .

The **open ball** of radius  $r > 0$  around  $\mathbf{x} \in \mathbb{R}^n$  is

$$\mathbb{B}(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| < r\}.$$

The **closed ball** of radius  $r > 0$  around  $\mathbf{x} \in \mathbb{R}^n$  is

$$\overline{\mathbb{B}}(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| \leq r\}.$$

#### Fact

- A set  $S \subseteq \mathbb{R}^n$  is **open** if for every  $\mathbf{x} \in S$  there exists  $r > 0$  such that

$$\mathbb{B}(\mathbf{x}, r) \subseteq S.$$

- A set  $S \subseteq \mathbb{R}^n$  is **closed** if every convergent sequence in  $S$  has its limit in  $S$ .

A set that contains only part of its boundary is, in general, neither open nor closed.

## LEC 4

## 4.1 Topology Review

A set  $S \subseteq \mathbb{R}^n$  is **compact** if it is closed and bounded (there is a more general notion of compactness in pure math).

A set  $S \subseteq \mathbb{R}^n$  is **bounded** if  $\exists r > 0$  such that

$$S \subseteq \mathbb{B}(\mathbf{0}, r).$$

**Theorem**

If  $S \subseteq \mathbb{R}^n$  is compact, nonempty, and  $f : S \rightarrow \mathbb{R}$  is continuous, then there exist  $\mathbf{x}_1, \mathbf{x}_2 \in S$  such that

$$f(\mathbf{x}_1) \leq f(\mathbf{x}) \leq f(\mathbf{x}_2), \quad \forall \mathbf{x} \in S.$$

(That is,  $f$  attains both its minimum and maximum on  $S$ .)

**Theorem**

If  $\mathbf{x}_1, \mathbf{x}_2, \dots$  is an arbitrary infinite sequence in a compact set  $S$ , then there exists a subsequence

$$\mathbf{x}_{k_1}, \mathbf{x}_{k_2}, \dots \quad \text{with} \quad k_1 < k_2 < \dots$$

that converges to some  $\mathbf{x} \in S$ .

## 4.2 Basics of Continuous Optimization

## 4.2.1 Minimizers

Let  $\Omega \subseteq \mathbb{R}^n$ , and let  $f : \Omega \rightarrow \mathbb{R}$  (we will also write  $\text{dom}(f) = \Omega$ ).

We say  $\mathbf{x}^* \in \text{dom}(f)$  is a **local minimizer** of  $f$  if there exists  $r > 0$  such that

$$f(\mathbf{x}^*) \leq f(\mathbf{x}), \quad \forall \mathbf{x} \in \text{dom}(f) \cap \mathbb{B}(\mathbf{x}^*, r).$$

We say  $\mathbf{x}^*$  is a **global minimizer** if

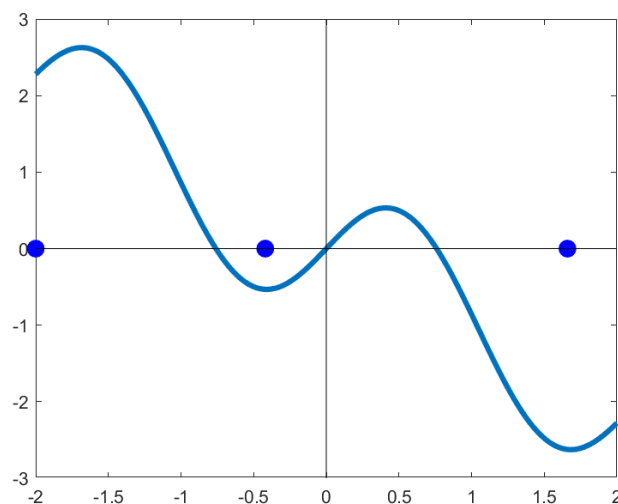
$$f(\mathbf{x}^*) \leq f(\mathbf{x}), \quad \forall \mathbf{x} \in \text{dom}(f).$$

**Example**

Let  $n = 1$  and  $\text{dom}(f) = [a, b]$ . In the picture below, the circled blue points are local minimizers; the third local minimizer is also the global minimizer.

We say  $\mathbf{x}^* \in \Omega$  is a **strict local minimizer** if there exists  $r > 0$  such that

$$f(\mathbf{x}^*) < f(\mathbf{x}), \quad \forall \mathbf{x} \in (\text{dom}(f) \cap \mathbb{B}(\mathbf{x}^*, r)) \setminus \{\mathbf{x}^*\}.$$



Remarks:

- $\Omega (= \text{dom}(f))$  could be all of  $\mathbb{R}^n$ .
- $\Omega$  could be specified by constraints.
- $\Omega$  could be the set where  $f$  is naturally defined (for example, if  $f(x) = x - \ln x$ , then  $\text{dom}(f) = \{x : x > 0\}$ ).

#### 4.2.2 Derivatives

Let  $f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subseteq \mathbb{R}^n$ , and let  $\mathbf{x} \in \Omega$ .

We say  $f$  is **differentiable** at  $\mathbf{x}$  if:

- $\mathbf{x} \in \text{int}(\Omega)$ , i.e., there exists  $r > 0$  such that  $\mathbb{B}(\mathbf{x}, r) \subseteq \Omega$ , and
- there exists  $\mathbf{g} \in \mathbb{R}^n$  (the *gradient* or *derivative* of  $f$  at  $\mathbf{x}$ ) such that

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}, \mathbf{h} \neq \mathbf{0}} \frac{f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) - \mathbf{g}^\top \mathbf{h}}{\|\mathbf{h}\|} = 0.$$

Equivalently, differentiability means that there exists a function  $\Phi_{\mathbf{x}}(\mathbf{h})$  such that

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \mathbf{g}^\top \mathbf{h} + \Phi_{\mathbf{x}}(\mathbf{h}),$$

with

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}, \mathbf{h} \neq \mathbf{0}} \frac{\Phi_{\mathbf{x}}(\mathbf{h})}{\|\mathbf{h}\|} = 0.$$

## Fact

If  $f$  is differentiable at  $\mathbf{x}$ , then the derivative is uniquely determined; we denote it by  $\nabla f(\mathbf{x})$ .

Given  $\mathbf{x} \in \text{int}(\Omega)$ , the **partial derivatives** are defined by

$$\frac{\partial f}{\partial x_i}(\mathbf{x}) = \lim_{h \rightarrow 0, h \neq 0} \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h}, \quad i = 1, \dots, n.$$

We say  $f : \Omega \rightarrow \mathbb{R}$  is  $\mathcal{C}^1$  or **continuously differentiable** at  $\mathbf{x}$  if all partial derivatives

$$\frac{\partial f}{\partial x_i}(\mathbf{y})$$

exist and are continuous

$$\forall i = 1, \dots, n, \quad \forall \mathbf{y} \in \mathbb{B}(\mathbf{x}, r) \text{ for some } r > 0.$$

A similar correction applies to the definition of  $\mathcal{C}^2$ .

## Theorem

If  $f$  is  $\mathcal{C}^1$  at  $\mathbf{x}$ , then

$$\nabla f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\mathbf{x}) \end{pmatrix}.$$

## Example

(Quadratic.)

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} + \mathbf{g}^\top \mathbf{x} + d,$$

where  $\mathbf{H} \in \mathbb{S}^n$ ,  $\mathbf{g} \in \mathbb{R}^n$ ,  $d \in \mathbb{R}$ .

Equivalently,

$$f(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n H(i, j) x(i) x(j) + \sum_{i=1}^n g(i) x(i) + d,$$

so

$$\frac{\partial f}{\partial x_i}(\mathbf{x}) = \sum_{k=1}^n H(i, k) x(k) + g(i),$$

and

$$\nabla f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\mathbf{x}) \end{pmatrix} = \begin{pmatrix} H(1, :) \mathbf{x} + g(1) \\ \vdots \\ H(n, :) \mathbf{x} + g(n) \end{pmatrix} = \mathbf{H} \mathbf{x} + \mathbf{g}.$$

### 4.2.3 Second Derivatives and Hessian

Given  $\Omega \subseteq \mathbb{R}^n$ , we say  $f$  is **twice differentiable** at  $\mathbf{x} \in \mathbb{R}^n$  if:

- a)  $\mathbf{x} \in \text{int}(\Omega)$  (the interior condition ensures that  $\mathbf{h}$  can come from every direction), and
- b) there exist  $\mathbf{g} \in \mathbb{R}^n$  (the gradient) and  $\mathbf{H} \in \mathbb{R}^{n \times n}$  (the **Hessian** or **second derivative**) such that

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \mathbf{g}^\top \mathbf{h} + \frac{1}{2} \mathbf{h}^\top \mathbf{H} \mathbf{h} + \psi_{\mathbf{x}}(\mathbf{h}),$$

where

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}, \mathbf{h} \neq \mathbf{0}} \frac{\psi_{\mathbf{x}}(\mathbf{h})}{\|\mathbf{h}\|^2} = 0.$$

If  $\mathbf{H}$  exists, it is uniquely determined; we denote it by  $\nabla^2 f(\mathbf{x})$ .

We say  $f$  is  $\mathcal{C}^2$  at  $\mathbf{x}$  if  $\nabla^2 f(\mathbf{y})$  is continuous for  $\mathbf{y}$  in some open ball around  $\mathbf{x}$ .

#### Theorem

If  $f$  is  $\mathcal{C}^2$  at  $\mathbf{x}$ , then

$$\nabla^2 f(\mathbf{x}) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n}(\mathbf{x}) \end{pmatrix},$$

and this matrix is symmetric:

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x}) = \frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}), \quad \forall i, j = 1, \dots, n.$$

#### Example

(Quadratic, Hessian.)

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} + \mathbf{g}^\top \mathbf{x} + d.$$

Then

$$\nabla f(\mathbf{x}) = \mathbf{H} \mathbf{x} + \mathbf{g}, \quad \nabla^2 f(\mathbf{x}) = \mathbf{H}.$$

### 4.2.4 Taylor's Theorem Variants

Let  $f : \Omega \rightarrow \mathbb{R}$  be  $\mathcal{C}^1$ , and let  $\mathbf{x} \in \text{int}(\Omega)$ . By the **Fundamental Theorem of Calculus**,

$$f(\mathbf{x} + \mathbf{p}) = f(\mathbf{x}) + \int_0^1 \nabla f(\mathbf{x} + \gamma \mathbf{p})^\top \mathbf{p} \, d\gamma.$$

Applying the **Mean Value Theorem for integrals**, we obtain

$$f(\mathbf{x} + \mathbf{p}) = f(\mathbf{x}) + \nabla f(\mathbf{x} + \gamma \mathbf{p})^\top \mathbf{p}$$

for some  $\gamma \in [0, 1]$ .

(Here we are implicitly assuming that  $\mathbf{x} \in \text{int}(\Omega)$ ,  $\mathbf{x} + \mathbf{p} \in \text{int}(\Omega)$ , and  $\mathbf{x} + \gamma\mathbf{p} \in \text{int}(\Omega)$  for all  $\gamma \in [0, 1]$ .)

These results require that  $\mathbf{x}$ ,  $\mathbf{x} + \mathbf{p}$ , and  $\mathbf{x} + \gamma\mathbf{p}$  all lie in  $\text{int}(\Omega)$  for every  $\gamma \in [0, 1]$ .

#### 4.2.5 Convexity

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

$$\forall \mathbf{x} \in C, \forall \mathbf{p} \text{ with } \mathbf{x} + \mathbf{p} \in C, \text{ we also have } \mathbf{x} + \gamma\mathbf{p} \in C \quad \forall \gamma \in [0, 1].$$

Equivalently,  $C$  is convex if and only if

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

#### Example

(i) Some trivial examples:  $\emptyset$ ,  $\{\mathbf{x}_0\}$ ,  $\mathbb{R}^n$ .

(ii) **Affine set.** Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$ ,

$$C = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{b}\}$$

is an affine set.

An equivalent definition:  $C$  is **affine** if and only if

$$\forall \mathbf{x}, \mathbf{y} \in C, \forall \gamma \in \mathbb{R}, \quad (1 - \gamma)\mathbf{x} + \gamma\mathbf{y} \in C,$$

where  $\gamma$  is allowed to be any real number (not just in  $[0, 1]$ ).

Note: Affine sets are always convex.

(iii) Open or closed ball in any norm:

$$\mathbb{B}_{\square}(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\|_{\square} < r\}, \quad \bar{\mathbb{B}}_{\square}(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\|_{\square} \leq r\}.$$

(iv) **Halfspace.** Given  $\mathbf{A} \in \mathbb{R}^n$  and  $b \in \mathbb{R}$ ,

$$H = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}^{\top} \mathbf{x} \leq b\}.$$

(v) **Polyhedron.** A polyhedron is the intersection of a finite number of halfspaces.

Given  $\mathbf{A}_1, \dots, \mathbf{A}_m \in \mathbb{R}^n$  and  $b_1, \dots, b_m \in \mathbb{R}$ , define

$$C = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}_1^{\top} \mathbf{x} \leq b_1, \dots, \mathbf{A}_m^{\top} \mathbf{x} \leq b_m\}.$$

Let

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_1^{\top} \\ \vdots \\ \mathbf{A}_m^{\top} \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}.$$

Then equivalently,

$$C = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} \leq \mathbf{b}\},$$

where the inequality is interpreted componentwise.

Note: For  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^m$ , we write  $\mathbf{u} \leq \mathbf{v}$  if

$$u(i) \leq v(i), \quad \forall i = 1, \dots, m.$$

The same componentwise interpretation holds for  $<$ ,  $>$ , and  $\geq$ .

Note: Affine sets are a special case of polyhedra. Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$ ,

$$\mathbf{Ax} = \mathbf{b} \iff \begin{pmatrix} \mathbf{A} \\ -\mathbf{A} \end{pmatrix} \mathbf{x} \leq \begin{pmatrix} \mathbf{b} \\ -\mathbf{b} \end{pmatrix}.$$

#### 4.2.6 Lipschitz Continuity

Let  $\mathbf{f} : \Omega \rightarrow \mathbb{R}^m$ ,  $\Omega \subseteq \mathbb{R}^n$ . We say  $\mathbf{f}$  is **Lipschitz continuous** with modulus  $L$  if

$$\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \Omega.$$



## LEC 5

5.1  $L$ -smooth Functions

Given  $\mathbf{f} : \Omega \rightarrow \mathbb{R}^m$ ,  $\Omega \subseteq \mathbb{R}^n$ , we say  $\mathbf{f}$  is **Lipschitz continuous** with modulus  $L$  if

$$\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \Omega.$$

Given  $f : \Omega \rightarrow \mathbb{R}$ , we say  $f$  is  **$L$ -smooth** (smooth with modulus  $L$ ) if  $f$  is differentiable and its gradient is  $L$ -Lipschitz continuous, i.e.,

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \Omega.$$

**Lemma**

If  $f$  is  $L$ -smooth and  $\text{dom}(f)$  is open and convex, then

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f).$$

(This looks like a Taylor expansion with a quadratic remainder term.)

**Proof**

$$\begin{aligned} f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) &= \left( \int_0^1 \nabla f(\mathbf{x} + \gamma(\mathbf{y} - \mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) \, d\gamma \right) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \\ &= \int_0^1 [\nabla f(\mathbf{x} + \gamma(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})]^\top (\mathbf{y} - \mathbf{x}) \, d\gamma \\ &\leq \int_0^1 \|\nabla f(\mathbf{x} + \gamma(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\| \cdot \|\mathbf{y} - \mathbf{x}\| \, d\gamma \quad (\text{Cauchy-Schwarz}) \\ &\leq \int_0^1 L\|\mathbf{x} + \gamma(\mathbf{y} - \mathbf{x}) - \mathbf{x}\| \cdot \|\mathbf{y} - \mathbf{x}\| \, d\gamma \quad (L\text{-smoothness}) \\ &= \int_0^1 L\gamma \|\mathbf{y} - \mathbf{x}\|^2 \, d\gamma \\ &= L\|\mathbf{y} - \mathbf{x}\|^2 \int_0^1 \gamma \, d\gamma \\ &= \frac{1}{2} L\|\mathbf{y} - \mathbf{x}\|^2. \end{aligned}$$

Rearranging gives the desired inequality.

## 5.2 Necessary and Sufficient Conditions for Optimality

### 5.2.1 First-order Necessary Condition

#### Theorem

Let  $f : \Omega \rightarrow \mathbb{R}$  be differentiable at a point  $\mathbf{x}^* \in \Omega$ , and assume that  $\mathbf{x}^*$  is an interior point of  $\Omega$  and a local minimizer of  $f$ . Then

$$\nabla f(\mathbf{x}^*) = \mathbf{0}.$$

#### Proof

(Contrapositive.) Assume  $\nabla f(\mathbf{x}^*) \neq \mathbf{0}$ . Let

$$\mathbf{p} := -\nabla f(\mathbf{x}^*).$$

By the definition of differentiability, there exists a function  $\phi_{\mathbf{x}^*}$  such that, for all  $|\alpha|$  sufficiently small,

$$f(\mathbf{x}^* + \alpha\mathbf{p}) = f(\mathbf{x}^*) + \alpha \nabla f(\mathbf{x}^*)^\top \mathbf{p} + \phi_{\mathbf{x}^*}(\alpha\mathbf{p}),$$

and

$$\frac{\phi_{\mathbf{x}^*}(\alpha\mathbf{p})}{\|\alpha\mathbf{p}\|} \rightarrow 0 \quad \Leftrightarrow \quad \frac{\phi_{\mathbf{x}^*}(\alpha\mathbf{p})}{|\alpha|} \rightarrow 0 \quad \text{as } \alpha \rightarrow 0.$$

Since  $\mathbf{p} = -\nabla f(\mathbf{x}^*)$ ,

$$f(\mathbf{x}^* + \alpha\mathbf{p}) = f(\mathbf{x}^*) - \alpha \|\nabla f(\mathbf{x}^*)\|^2 + \phi_{\mathbf{x}^*}(\alpha\mathbf{p}).$$

Because  $\|\mathbf{p}\| > 0$  and  $\phi_{\mathbf{x}^*}(\alpha\mathbf{p})/|\alpha| \rightarrow 0$ , there exists  $\bar{\alpha} > 0$  such that

$$\frac{\|\phi_{\mathbf{x}^*}(\alpha\mathbf{p})\|}{|\alpha|\|\mathbf{p}\|} \leq \frac{\|\mathbf{p}\|}{2}, \quad \forall |\alpha| \leq \bar{\alpha}.$$

For such  $\alpha \in (0, \bar{\alpha}]$ ,

$$\begin{aligned} f(\mathbf{x}^* + \alpha\mathbf{p}) &\leq f(\mathbf{x}^*) - \alpha \|\nabla f(\mathbf{x}^*)\|^2 + \frac{\|\mathbf{p}\|}{2} |\alpha| \|\mathbf{p}\| \\ &= f(\mathbf{x}^*) - \frac{\alpha}{2} \|\mathbf{p}\|^2 \\ &< f(\mathbf{x}^*). \end{aligned}$$

Thus  $f$  strictly decreases in direction  $\mathbf{p}$  away from  $\mathbf{x}^*$ , so  $\mathbf{x}^*$  cannot be a local minimizer.

## 5.2.2 Second-order Necessary Condition

## Theorem

Suppose  $f$  is  $\mathcal{C}^2$  at  $\mathbf{x}^* \in \Omega$ , and  $\mathbf{x}^*$  is a local minimizer. Then

$$\nabla f(\mathbf{x}^*) = \mathbf{0} \quad \text{and} \quad \nabla^2 f(\mathbf{x}^*) \succeq \mathbf{0} \text{ (positive semidefinite).}$$

## Proof

**(Contrapositive idea.)** We already know from Theorem 1 that a local minimizer must satisfy  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ . Now assume  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  but  $\nabla^2 f(\mathbf{x}^*)$  is not positive semidefinite. Then there exists  $\mathbf{p}$  such that

$$\mathbf{p}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{p} < 0.$$

By the second-order Taylor expansion,

$$f(\mathbf{x}^* + \alpha \mathbf{p}) = f(\mathbf{x}^*) + \alpha \nabla f(\mathbf{x}^*)^\top \mathbf{p} + \frac{1}{2} \alpha^2 \mathbf{p}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{p} + \psi_{\mathbf{x}^*}(\alpha \mathbf{p}),$$

and we know

$$\frac{\psi_{\mathbf{x}^*}(\alpha \mathbf{p})}{\alpha^2} \rightarrow 0 \quad \text{as } \alpha \rightarrow 0.$$

Choose  $\bar{\alpha} > 0$  sufficiently small such that

$$\frac{|\psi_{\mathbf{x}^*}(\alpha \mathbf{p})|}{\alpha^2} \leq \frac{1}{4} |\mathbf{p}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{p}|, \quad \forall \alpha \in (0, \bar{\alpha}].$$

Using  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ , we get for all  $\alpha \in (0, \bar{\alpha}]$ ,

$$\begin{aligned} f(\mathbf{x}^* + \alpha \mathbf{p}) &\leq f(\mathbf{x}^*) + \frac{1}{2} \alpha^2 \mathbf{p}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{p} + \frac{1}{4} \alpha^2 |\mathbf{p}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{p}| \\ &= f(\mathbf{x}^*) - \frac{1}{2} \alpha^2 |\mathbf{p}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{p}| + \frac{1}{4} \alpha^2 |\mathbf{p}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{p}| \\ &= f(\mathbf{x}^*) - \frac{\alpha^2}{4} |\mathbf{p}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{p}| \\ &< f(\mathbf{x}^*), \end{aligned}$$

so  $\mathbf{x}^*$  is not a local minimizer.

## Example

**(Quadratic and Least Squares.)**

Let

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} + \mathbf{g}^\top \mathbf{x} + d.$$

Then

$$\nabla f(\mathbf{x}) = \mathbf{H} \mathbf{x} + \mathbf{g},$$

and

$$\nabla f(\mathbf{x}) = \mathbf{0} \iff \mathbf{H} \mathbf{x} + \mathbf{g} = \mathbf{0} \iff \mathbf{H} \mathbf{x} = -\mathbf{g}.$$

Also,

$$\nabla^2 f(\mathbf{x}) = \mathbf{H},$$

so

$\nabla^2 f(\mathbf{x})$  is positive semidefinite  $\iff \mathbf{H}$  is positive semidefinite.

For quadratics, these conditions are **necessary and sufficient** for optimality.

For linear least squares

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|^2,$$

we have

$$\nabla f(\mathbf{x}) = \mathbf{A}^\top (\mathbf{Ax} - \mathbf{y}), \quad \nabla^2 f(\mathbf{x}) = \mathbf{A}^\top \mathbf{A}.$$

### 5.2.3 Second-order Sufficient Condition

#### Theorem

Let  $f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subseteq \mathbb{R}^n$ . Suppose  $\mathbf{x}^* \in \text{int}(\Omega)$ ,  $f$  is  $\mathcal{C}^2$  at  $\mathbf{x}^*$ ,  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ , and

$$\nabla^2 f(\mathbf{x}^*) \succ \mathbf{0}.$$

Then  $\mathbf{x}^*$  is a **strict local minimizer** of  $f$ .

#### Proof

##### Lemma

If  $\mathbf{U} \succ \mathbf{0}$ , then there exists  $r > 0$  such that

$$\mathbf{V} \succ \mathbf{0}, \quad \forall \mathbf{V} \in \mathbb{B}(\mathbf{U}, r),$$

where  $\mathbb{B}(\mathbf{U}, r) \subseteq \mathbb{S}^n$  is the Frobenius-norm ball:

$$\|\mathbf{X}\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^n X(i, j)^2}.$$

Equivalently, the set of positive definite matrices is an open subset of  $\mathbb{S}^n$ .

**Proof****Proof of Lemma.**

- The eigenvalues are continuous functions of the matrix entries.
- If  $\mathbf{U} \succ \mathbf{0}$ , then  $\lambda_{\min}(\mathbf{U}) > 0$ , where

$$\lambda_{\min}(\mathbf{U}) = \min_{\|\mathbf{x}\|=1} \mathbf{x}^\top \mathbf{U} \mathbf{x}.$$

These imply that for sufficiently small  $\Delta \in \mathbb{S}^n$ ,

$$\begin{aligned} \lambda_{\min}(\mathbf{U} + \Delta) &= \min_{\|\mathbf{x}\|=1} \mathbf{x}^\top (\mathbf{U} + \Delta) \mathbf{x} \\ &\geq \min_{\|\mathbf{x}\|=1} (\mathbf{x}^\top \mathbf{U} \mathbf{x} - |\mathbf{x}^\top \Delta \mathbf{x}|) \\ &\geq \min_{\|\mathbf{x}\|=1} \mathbf{x}^\top \mathbf{U} \mathbf{x} - \min_{\|\mathbf{x}\|=1} (\|\mathbf{x}\| \|\Delta\|_F \|\mathbf{x}\|) \\ &= \lambda_{\min}(\mathbf{U}) - \|\Delta\|_F \\ &\geq \frac{\lambda_{\min}(\mathbf{U})}{2} \quad \text{provided } \|\Delta\|_F \leq \frac{\lambda_{\min}(\mathbf{U})}{2}. \end{aligned}$$

Hence  $\lambda_{\min}(\mathbf{U} + \Delta) > 0$ , so  $\mathbf{U} + \Delta \succ \mathbf{0}$ .

**Proof of Theorem.** Because  $f$  is  $\mathcal{C}^2$  at  $\mathbf{x}^*$  and  $\nabla^2 f(\mathbf{x}^*) \succ \mathbf{0}$ , by the lemma there exists  $\rho > 0$  such that

$$\mathbb{B}(\mathbf{x}^*, \rho) \subseteq \Omega \quad \text{and} \quad \nabla^2 f(\mathbf{x}) \succ \mathbf{0}, \quad \forall \mathbf{x} \in \mathbb{B}(\mathbf{x}^*, \rho).$$

Choose any  $\mathbf{y} \in \mathbb{B}(\mathbf{x}^*, \rho) \setminus \{\mathbf{x}^*\}$ . By the second-order Taylor expansion, there exists  $\gamma \in [0, 1]$  such that

$$f(\mathbf{y}) = f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^\top (\mathbf{y} - \mathbf{x}^*) + \frac{1}{2} (\mathbf{y} - \mathbf{x}^*)^\top \nabla^2 f(\mathbf{x}^* + \gamma(\mathbf{y} - \mathbf{x}^*)) (\mathbf{y} - \mathbf{x}^*).$$

Since  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  and

$$\nabla^2 f(\mathbf{x}^* + \gamma(\mathbf{y} - \mathbf{x}^*)) \succ \mathbf{0} \quad (\text{as } \mathbf{x}^* + \gamma(\mathbf{y} - \mathbf{x}^*) \in \mathbb{B}(\mathbf{x}^*, \rho)),$$

we have

$$f(\mathbf{y}) > f(\mathbf{x}^*).$$

Thus  $\mathbf{x}^*$  is a strict local minimizer.

## 5.2.4 Summary of First- and Second-order Conditions

## Summary

Let  $f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subseteq \mathbb{R}^n$ , and  $\mathbf{x}^* \in \Omega$ . Suppose  $f$  is  $\mathcal{C}^2$  at  $\mathbf{x}^*$ . Then we have the chain:

$$\begin{aligned} \nabla f(\mathbf{x}^*) = \mathbf{0}, \quad \nabla^2 f(\mathbf{x}^*) \succ \mathbf{0} &\Rightarrow \mathbf{x}^* \text{ is a strict local minimizer} \\ &\Rightarrow \mathbf{x}^* \text{ is a local minimizer} \\ &\Rightarrow \nabla f(\mathbf{x}^*) = \mathbf{0}, \quad \nabla^2 f(\mathbf{x}^*) \succeq \mathbf{0}. \end{aligned}$$

This chain of implications cannot be reversed in general.

## Example

- $f(x) = x^4$ :

$x^* = 0$  is a strict local minimizer, but  $f''(0) = 0$  (not  $> 0$ ).

- $f(x) = 1$ :

$x^* = 0$  is a local minimizer (in fact, every point is), but not a strict local minimizer.

- $f(x) = x^3$ :

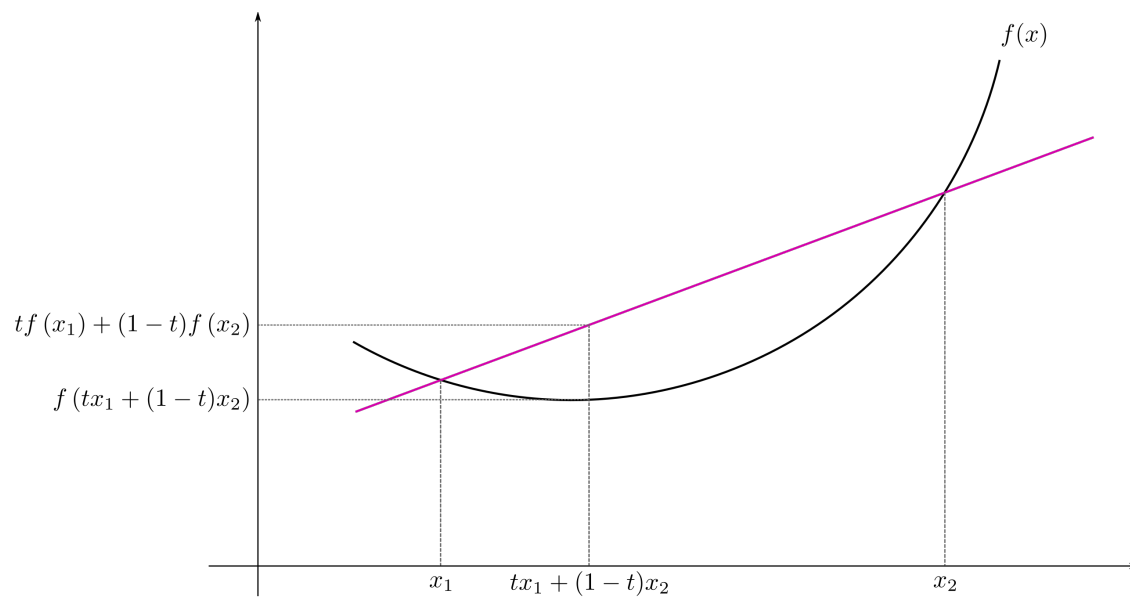
$f'(0) = 0$ ,  $f''(0) = 0$ , but  $x^* = 0$  is not a local minimizer.

### 5.3 Convex Functions

Let  $f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subseteq \mathbb{R}^n$ , and assume  $\Omega = \text{dom}(f)$  is a convex set. We say  $f$  is a **convex function** if

$$\forall \mathbf{x}, \mathbf{y} \in \Omega, \forall \lambda \in [0, 1],$$

$$f((1 - \lambda)\mathbf{x} + \lambda\mathbf{y}) \leq (1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y}).$$



## LEC 6

## 6.1 Convex Functions

## 6.1.1 Definitions

Let  $f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subseteq \mathbb{R}^n$ , and  $\Omega = \text{dom}(f)$ . Assume  $\Omega$  is convex. We say  $f$  is a **convex function** if

$$\forall \mathbf{x}, \mathbf{y} \in \Omega, \forall \lambda \in [0, 1], \quad f((1 - \lambda)\mathbf{x} + \lambda\mathbf{y}) \leq (1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y}).$$

**Extended-value notation.** If  $f : \Omega \rightarrow \mathbb{R}$  is convex, we extend it by setting

$$f(\mathbf{x}) = \infty \quad \text{for } \mathbf{x} \notin \Omega.$$

Thus we can regard

$$f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\},$$

and define

$$\text{dom}(f) = \{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) < \infty\}.$$

The convexity inequality above is then interpreted for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ; it enforces that  $\text{dom}(f)$  is convex.

## 6.1.2 Examples

## Example

1.  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{H}\mathbf{x} + \mathbf{g}^\top \mathbf{x} + d, \quad \mathbf{H} \succeq \mathbf{0}.$
2.  $f(\mathbf{x}) = \|\mathbf{x}\|_{\square}$  (any norm).
3.  $f(x) = e^x.$
4.  $f(x) = \begin{cases} -\ln x, & x > 0, \\ \infty, & x \leq 0, \end{cases}$
5.  $f(x) = \begin{cases} \frac{1}{x}, & x > 0, \\ \infty, & x \leq 0, \end{cases}$
6. If  $g_1, \dots, g_m : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  are convex and  $\alpha_1, \dots, \alpha_m \geq 0$ ,  
then  $\alpha_1 g_1 + \dots + \alpha_m g_m$  is convex.
7. Under the same assumption as in 6,  $\max\{\alpha_1 g_1, \dots, \alpha_m g_m\}$  is convex.  
Example:  $f(x) = \max(x, -x) = |x|$  is convex.



**Theorem**

If  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  is convex and  $\mathbf{x}^* \in \text{dom}(f)$  is a local minimizer, then  $\mathbf{x}^*$  is a global minimizer.

**Proof**

Assume  $\mathbf{x}^*$  is a local minimizer and let  $\mathbf{y} \in \mathbb{R}^n$  be arbitrary. Define

$$\mathbf{y}_i := \left(1 - \frac{1}{i}\right) \mathbf{x}^* + \frac{1}{i} \mathbf{y}, \quad i = 1, 2, 3, \dots$$

Then  $\mathbf{y}_i \rightarrow \mathbf{x}^*$  as  $i \rightarrow \infty$ . Since  $\mathbf{x}^*$  is a local minimizer, there exists  $k$  such that

$$f(\mathbf{y}_i) \geq f(\mathbf{x}^*), \quad \forall i \geq k.$$

In particular,

$$f\left(\left(1 - \frac{1}{k}\right) \mathbf{x}^* + \frac{1}{k} \mathbf{y}\right) \geq f(\mathbf{x}^*).$$

By convexity,

$$f\left(\left(1 - \frac{1}{k}\right) \mathbf{x}^* + \frac{1}{k} \mathbf{y}\right) \leq \left(1 - \frac{1}{k}\right) f(\mathbf{x}^*) + \frac{1}{k} f(\mathbf{y}).$$

Combining,

$$\left(1 - \frac{1}{k}\right) f(\mathbf{x}^*) + \frac{1}{k} f(\mathbf{y}) \geq f(\mathbf{x}^*).$$

Rearranging gives

$$f(\mathbf{y}) \geq f(\mathbf{x}^*).$$

Since  $\mathbf{y}$  was arbitrary,  $\mathbf{x}^*$  is a global minimizer.

**6.2 Differentiable Convex Functions**

Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  is convex and differentiable at  $\mathbf{x} \in \text{intdom}(f)$ . Then for all  $\mathbf{y} \in \mathbb{R}^n$  and all  $\alpha \in [0, 1]$ ,

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

Write the left-hand side using differentiability at  $\mathbf{x}$ :

$$\begin{aligned} f((1 - \alpha)\mathbf{x} + \alpha\mathbf{y}) &= f(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x})) \\ &= f(\mathbf{x}) + \alpha \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \phi(\alpha(\mathbf{y} - \mathbf{x})), \end{aligned}$$

for some remainder term  $\phi$  with

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}, \mathbf{h} \neq \mathbf{0}} \frac{\phi(\mathbf{h})}{\|\mathbf{h}\|} = 0.$$

Combining with convexity,

$$f(\mathbf{x}) + \alpha \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \phi(\alpha(\mathbf{y} - \mathbf{x})) \leq (1 - \alpha)f(\mathbf{x}) + \alpha f(\mathbf{y}).$$

Assuming  $\alpha > 0$ , subtract  $f(\mathbf{x})$  and divide by  $\alpha$ :

$$\nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\phi(\alpha(\mathbf{y} - \mathbf{x}))}{\alpha} \leq f(\mathbf{y}) - f(\mathbf{x}).$$

Taking the limit  $\alpha \downarrow 0$ , the remainder term vanishes and we obtain the

#### Fact

**Subgradient Inequality.** For convex differentiable  $f$  and any  $\mathbf{x} \in \text{intdom}(f)$ ,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{y} \in \mathbb{R}^n.$$

The converse also holds: if this inequality is satisfied for all  $\mathbf{x}, \mathbf{y}$ , then  $f$  is convex.

#### Theorem

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  be convex. Suppose  $\mathbf{x}^* \in \text{intdom}(f)$ ,  $f$  is differentiable at  $\mathbf{x}^*$ , and  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ . Then  $\mathbf{x}^*$  is a global minimizer of  $f$ .

#### Proof

Applying the subgradient inequality at  $\mathbf{x}^*$ ,

$$f(\mathbf{y}) \geq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^\top (\mathbf{y} - \mathbf{x}^*) = f(\mathbf{x}^*), \quad \forall \mathbf{y} \in \mathbb{R}^n,$$

since  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ . Hence  $\mathbf{x}^*$  is a global minimizer.

## 6.3 Descent Methods

### 6.3.1 Descent Directions

Let  $f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subseteq \mathbb{R}^n$ , and  $\mathbf{x} \in \Omega$ . We say  $\mathbf{d} \in \mathbb{R}^n$  is a **descent direction** at  $\mathbf{x}$  if there exists  $\bar{t} > 0$  such that

$$\mathbf{x} + t\mathbf{d} \in \Omega, \quad \forall t \in [0, \bar{t}], \quad \text{and} \quad f(\mathbf{x} + t\mathbf{d}) < f(\mathbf{x}), \quad \forall t \in (0, \bar{t}].$$

#### Theorem

Let  $f : \Omega \rightarrow \mathbb{R}$  with  $\Omega \subseteq \mathbb{R}^n$ . Assume  $\mathbf{x} \in \text{int}(\Omega)$  and  $f$  is differentiable at  $\mathbf{x}$ . Then any  $\mathbf{d}$  satisfying

$$\mathbf{d}^\top \nabla f(\mathbf{x}) < 0$$

is a descent direction at  $\mathbf{x}$ .

#### Proof

Assume  $\nabla f(\mathbf{x}) \neq \mathbf{0}$  and  $\mathbf{d} \neq \mathbf{0}$  (otherwise the statement is vacuous). Differentiability at  $\mathbf{x}$  means

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}, \mathbf{h} \neq \mathbf{0}} \frac{f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top \mathbf{h}}{\|\mathbf{h}\|} = 0.$$

Take  $\mathbf{h} = t\mathbf{d}$  with  $t \downarrow 0$ :

$$\lim_{t \downarrow 0} \frac{f(\mathbf{x} + t\mathbf{d}) - f(\mathbf{x}) - t \nabla f(\mathbf{x})^\top \mathbf{d}}{t} = 0,$$

or equivalently,

$$\lim_{t \downarrow 0} \left( \frac{f(\mathbf{x} + t\mathbf{d}) - f(\mathbf{x})}{t} - \nabla f(\mathbf{x})^\top \mathbf{d} \right) = 0.$$

By assumption,  $\nabla f(\mathbf{x})^\top \mathbf{d} < 0$ . Therefore, for sufficiently small  $t > 0$ ,

$$\frac{f(\mathbf{x} + t\mathbf{d}) - f(\mathbf{x})}{t} < 0,$$

which implies  $f(\mathbf{x} + t\mathbf{d}) < f(\mathbf{x})$  for all sufficiently small  $t > 0$ . Thus  $\mathbf{d}$  is a descent direction.

**Steepest descent direction.** Assume  $\nabla f(\mathbf{x}) \neq \mathbf{0}$ . Among all unit vectors, the direction that minimizes the directional derivative is

$$\arg \min \{ \mathbf{d}^\top \nabla f(\mathbf{x}) : \|\mathbf{d}\| = 1 \} = -\frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|}.$$

This follows from the strong form of the Cauchy–Schwarz inequality:

$$\mathbf{u}^\top \mathbf{v} \leq \|\mathbf{u}\| \|\mathbf{v}\|$$

with equality if and only if either  $\mathbf{u} = \mathbf{0}$  or  $\mathbf{v} = \mathbf{0}$ , or there exists  $\lambda \geq 0$  such that  $\mathbf{u} = \lambda \mathbf{v}$ .

### 6.3.2 Gradient Descent Algorithm (Steepest Descent)

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be differentiable, and choose an initial point  $\mathbf{x}^0 \in \mathbb{R}^n$ .

For  $k = 0, 1, 2, \dots$ :

- select a step size  $\alpha_k > 0$ ,
- update

$$\mathbf{x}^{k+1} := \mathbf{x}^k - \alpha_k \nabla f(\mathbf{x}^k).$$

End.

Here  $\alpha_k$  is the **step size** (learning rate). The procedure used to choose  $\alpha_k$  is called a **line search**.

If  $f$  is  $L$ -smooth, then the constant choice  $\alpha_k \equiv \frac{1}{L}$  always works. Any fixed  $\alpha_k \in (0, \frac{1}{L})$  also leads to convergence but typically more slowly. A common stopping rule is  $\|\nabla f(\mathbf{x}^k)\| \leq \text{tolerance}$ .

#### Theorem

**(Convex case of Gradient Descent)** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be convex and  $L$ -smooth. Suppose there exists  $\mathbf{x}^*$  such that

$$\nabla f(\mathbf{x}^*) = \mathbf{0}.$$

Run gradient descent from  $\mathbf{x}^0$  with step sizes  $\alpha_k \equiv \frac{1}{L}$ . Then

$$f(\mathbf{x}^k) - f(\mathbf{x}^*) \leq \frac{L\|\mathbf{x}^0 - \mathbf{x}^*\|^2}{2k}, \quad \forall k = 1, 2, \dots$$

### Proof

From  $L$ -smoothness,

$$\begin{aligned} f(\mathbf{x}^{k+1}) &\leq f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^\top (\mathbf{x}^{k+1} - \mathbf{x}^k) + \frac{L}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \\ &= f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^\top \left( -\frac{1}{L} \nabla f(\mathbf{x}^k) \right) + \frac{L}{2} \left\| -\frac{1}{L} \nabla f(\mathbf{x}^k) \right\|^2 \\ &= f(\mathbf{x}^k) - \frac{1}{L} \|\nabla f(\mathbf{x}^k)\|^2 + \frac{1}{2L} \|\nabla f(\mathbf{x}^k)\|^2 \\ &= f(\mathbf{x}^k) - \frac{1}{2L} \|\nabla f(\mathbf{x}^k)\|^2. \end{aligned}$$

From the subgradient inequality for convex differentiable  $f$ ,

$$f(\mathbf{x}^*) \geq f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^\top (\mathbf{x}^* - \mathbf{x}^k),$$

so

$$f(\mathbf{x}^k) \leq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^k)^\top (\mathbf{x}^k - \mathbf{x}^*).$$

Combining,

$$f(\mathbf{x}^{k+1}) \leq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^k)^\top (\mathbf{x}^k - \mathbf{x}^*) - \frac{1}{2L} \|\nabla f(\mathbf{x}^k)\|^2.$$

Using the update  $\mathbf{x}^{k+1} = \mathbf{x}^k - \frac{1}{L} \nabla f(\mathbf{x}^k)$ , we compute

$$\begin{aligned} \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 &= \left\| \mathbf{x}^k - \frac{1}{L} \nabla f(\mathbf{x}^k) - \mathbf{x}^* \right\|^2 \\ &= \|\mathbf{x}^k - \mathbf{x}^*\|^2 - \frac{2}{L} \nabla f(\mathbf{x}^k)^\top (\mathbf{x}^k - \mathbf{x}^*) + \frac{1}{L^2} \|\nabla f(\mathbf{x}^k)\|^2. \end{aligned}$$

Rearranging,

$$\nabla f(\mathbf{x}^k)^\top (\mathbf{x}^k - \mathbf{x}^*) = \frac{L}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2) + \frac{1}{2L} \|\nabla f(\mathbf{x}^k)\|^2.$$

Substitute into the bound for  $f(\mathbf{x}^{k+1})$ :

$$\begin{aligned} f(\mathbf{x}^{k+1}) &\leq f(\mathbf{x}^*) + \frac{L}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2) + \frac{1}{2L} \|\nabla f(\mathbf{x}^k)\|^2 - \frac{1}{2L} \|\nabla f(\mathbf{x}^k)\|^2 \\ &= f(\mathbf{x}^*) + \frac{L}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2). \end{aligned}$$

Thus

$$f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*) \leq \frac{L}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2).$$

Summing from  $k = 0$  to  $k = \ell - 1$ ,

$$\sum_{k=0}^{\ell-1} (f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*)) \leq \frac{L}{2} (\|\mathbf{x}^0 - \mathbf{x}^*\|^2 - \|\mathbf{x}^\ell - \mathbf{x}^*\|^2) \leq \frac{L}{2} \|\mathbf{x}^0 - \mathbf{x}^*\|^2.$$

Since  $f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*) \geq 0$  and is nonincreasing, we have

$$\sum_{k=0}^{\ell-1} (f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*)) \geq \ell(f(\mathbf{x}^\ell) - f(\mathbf{x}^*)).$$

Therefore,

$$\ell(f(\mathbf{x}^\ell) - f(\mathbf{x}^*)) \leq \frac{L}{2} \|\mathbf{x}^0 - \mathbf{x}^*\|^2,$$

which gives

$$f(\mathbf{x}^\ell) - f(\mathbf{x}^*) \leq \frac{L\|\mathbf{x}^0 - \mathbf{x}^*\|^2}{2\ell}.$$

Renaming  $\ell$  as  $k$  yields the stated result.

## LEC 7

## 7.1 Descent Methods

7.1.1 Analysis of Gradient Descent for  $L$ -smooth Convex Case

## Theorem

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be convex and  $L$ -smooth, and suppose there exists  $\mathbf{x}^*$  such that

$$\nabla f(\mathbf{x}^*) = \mathbf{0}$$

(so  $\mathbf{x}^*$  is a global minimizer). If gradient descent is initiated at  $\mathbf{x}^0$  with step sizes

$$\alpha_k \equiv \frac{1}{L},$$

then

$$f(\mathbf{x}^k) - f(\mathbf{x}^*) \leq \frac{L\|\mathbf{x}^0 - \mathbf{x}^*\|^2}{2k}, \quad \forall k = 1, 2, \dots$$

Aside:

$$f(\mathbf{x}^l) \leq f(\mathbf{x}^{l-1}) \leq \dots \leq f(\mathbf{x}^0).$$

Any method with this monotone decrease property is called a **descent method**.

The proof was given in the previous lecture.

7.1.2 Analysis of Gradient Descent for  $L$ -smooth Nonconvex Case

For minimizing nonconvex functions, global minimizers are usually hard to find, even for  $n = 1$ . In higher dimensions, even local minimizers can be difficult.

So we often settle for **stationary points**, i.e. seek  $\mathbf{x}^*$  such that

$$\nabla f(\mathbf{x}^*) = \mathbf{0},$$

(although even this may be impossible if  $f$  has no stationary point).

## Example

$$f(x) = \sqrt{x^2 + 1} + x$$

is smooth, convex, and  $L$ -smooth with  $L = 1$ , but it has no stationary point.

If gradient descent is applied to this  $f$ , the sequence  $x^k$  tends to  $-\infty$  and  $f'(x^k) \rightarrow 0$ .

## Theorem

Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $L$ -smooth (not necessarily convex) and bounded below by  $f_{\min}$ . Apply gradient descent with step size  $1/L$  starting at  $\mathbf{x}^0$ . Then:

(i)

$$\min_{0 \leq j \leq k-1} \|\nabla f(\mathbf{x}^j)\| \leq \sqrt{\frac{2L(f(\mathbf{x}^0) - f_{\min})}{k}},$$

(ii)

$$\nabla f(\mathbf{x}^j) \rightarrow \mathbf{0} \quad \text{as } j \rightarrow \infty.$$

**Proof**

Recall from  $L$ -smoothness and the  $GD$  update with step size  $1/L$ :

$$f(\mathbf{x}^{j+1}) \leq f(\mathbf{x}^j) - \frac{1}{2L} \|\nabla f(\mathbf{x}^j)\|^2.$$

Summing from  $j = 0$  to  $k - 1$ ,

$$f(\mathbf{x}^k) \leq f(\mathbf{x}^0) - \frac{1}{2L} \sum_{j=0}^{k-1} \|\nabla f(\mathbf{x}^j)\|^2.$$

Since  $f(\mathbf{x}^k) \geq f_{\min}$ , we obtain

$$f_{\min} - f(\mathbf{x}^0) \leq -\frac{1}{2L} \sum_{j=0}^{k-1} \|\nabla f(\mathbf{x}^j)\|^2,$$

or equivalently

$$\frac{1}{2L} \sum_{j=0}^{k-1} \|\nabla f(\mathbf{x}^j)\|^2 \leq f(\mathbf{x}^0) - f_{\min}.$$

Taking  $k \rightarrow \infty$  shows

$$\sum_{j=0}^{\infty} \|\nabla f(\mathbf{x}^j)\|^2 < \infty,$$

so  $\|\nabla f(\mathbf{x}^j)\| \rightarrow 0$  as  $j \rightarrow \infty$ , proving (2).

For (1), note

$$\sum_{j=0}^{k-1} \|\nabla f(\mathbf{x}^j)\|^2 \geq k \cdot \min_{0 \leq j \leq k-1} \|\nabla f(\mathbf{x}^j)\|^2.$$

Thus

$$\frac{k}{2L} \min_{0 \leq j \leq k-1} \|\nabla f(\mathbf{x}^j)\|^2 \leq f(\mathbf{x}^0) - f_{\min},$$

which implies

$$\min_{0 \leq j \leq k-1} \|\nabla f(\mathbf{x}^j)\| \leq \sqrt{\frac{2L(f(\mathbf{x}^0) - f_{\min})}{k}}.$$

## 7.2 Strongly Convex Functions

### 7.2.1 Definitions

Say  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  is  **$m$ -strongly convex** with modulus  $m \geq 0$  if for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  and all  $\lambda \in [0, 1]$ ,

$$f((1 - \lambda)\mathbf{x} + \lambda\mathbf{y}) \leq (1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y}) - \frac{1}{2}\lambda(1 - \lambda)m\|\mathbf{x} - \mathbf{y}\|^2. \quad (1)$$

Notes:

- $m = 0$  is exactly the usual notion of convexity.
- “Strongly convex” (without specifying  $m$ ) means the above holds for some  $m > 0$ .

### 7.2.2 Differentiable Strongly Convex Functions

#### Theorem

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  be differentiable on  $\text{dom}(f)$ , and assume  $\text{dom}(f)$  is open and convex. Then  $f$  is strongly convex with modulus  $m$  if and only if

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{m}{2}\|\mathbf{y} - \mathbf{x}\|^2. \quad (2)$$

Note that the last term strengthens the usual subgradient inequality.

#### Proof

**Forward direction  $\Rightarrow$ .**

By differentiability, for any  $\mathbf{x}, \mathbf{y}$  and  $\lambda \in [0, 1]$ ,

$$f((1 - \lambda)\mathbf{x} + \lambda\mathbf{y}) = f(\mathbf{x}) + \lambda \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \phi_{\mathbf{x}}(\lambda(\mathbf{y} - \mathbf{x})), \quad (3)$$

where

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}, \mathbf{h} \neq \mathbf{0}} \frac{\phi_{\mathbf{x}}(\mathbf{h})}{\|\mathbf{h}\|} = 0.$$

Combine strong convexity (1) with (3):

$$\begin{aligned} f(\mathbf{x}) + \lambda \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \phi_{\mathbf{x}}(\lambda(\mathbf{y} - \mathbf{x})) \\ \leq (1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y}) - \frac{1}{2}m\lambda(1 - \lambda)\|\mathbf{y} - \mathbf{x}\|^2. \end{aligned}$$

Subtract  $f(\mathbf{x})$  and divide by  $\lambda > 0$ :

$$\nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\phi_{\mathbf{x}}(\lambda(\mathbf{y} - \mathbf{x}))}{\lambda} \leq f(\mathbf{y}) - f(\mathbf{x}) - \frac{1}{2}m(1 - \lambda)\|\mathbf{y} - \mathbf{x}\|^2.$$

Let  $\lambda \downarrow 0$ . The remainder term vanishes and  $(1 - \lambda) \rightarrow 1$ , yielding

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{m}{2}\|\mathbf{y} - \mathbf{x}\|^2,$$



which is (2).

**Backward direction**  $\Leftarrow$ .

Assume (2). Take  $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$  and  $\lambda \in [0, 1]$ . Let

$$\mathbf{z} := (1 - \lambda)\mathbf{x} + \lambda\mathbf{y}.$$

Apply (2) twice, once with  $(\mathbf{z}, \mathbf{x})$  and once with  $(\mathbf{z}, \mathbf{y})$ :

$$\text{a) } f(\mathbf{x}) \geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{z}) + \frac{m}{2} \|\mathbf{x} - \mathbf{z}\|^2,$$

$$\text{b) } f(\mathbf{y}) \geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{z}) + \frac{m}{2} \|\mathbf{y} - \mathbf{z}\|^2.$$

Multiply (a) by  $(1 - \lambda)$  and (b) by  $\lambda$ , then add:

$$\begin{aligned} & (1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y}) \\ & \geq (1 - \lambda) \left[ f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{z}) \right] + \lambda \left[ f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{z}) \right] \\ & \quad + \frac{m}{2} \left[ (1 - \lambda) \|\mathbf{x} - \mathbf{z}\|^2 + \lambda \|\mathbf{y} - \mathbf{z}\|^2 \right]. \end{aligned}$$

Note that  $(1 - \lambda)(\mathbf{x} - \mathbf{z}) + \lambda(\mathbf{y} - \mathbf{z}) = \mathbf{0}$ , and

$$(1 - \lambda) \|\mathbf{x} - \mathbf{z}\|^2 + \lambda \|\mathbf{y} - \mathbf{z}\|^2 = \lambda(1 - \lambda) \|\mathbf{x} - \mathbf{y}\|^2.$$

Hence

$$(1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y}) \geq f(\mathbf{z}) + \frac{m}{2} \lambda(1 - \lambda) \|\mathbf{x} - \mathbf{y}\|^2,$$

which is exactly the strong convexity inequality (1).

### 7.2.3 Second Derivatives and Convexity

#### Theorem

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  be  $\mathcal{C}^2$  on  $\text{dom}(f)$ , and assume  $\text{dom}(f)$  is open and convex. Then  $f$  is strongly convex with modulus  $m$  on  $\text{dom}(f)$  if and only if

$$\nabla^2 f(\mathbf{x}) \succeq m\mathbf{I}, \quad \forall \mathbf{x} \in \text{dom}(f).$$

In particular,

$$f \text{ convex} \iff \nabla^2 f(\mathbf{x}) \text{ is positive semidefinite for all } \mathbf{x}.$$

#### Proof

Fix  $\mathbf{x} \in \text{dom}(f)$  and  $\mathbf{u} \in \mathbb{R}^n$ . Choose  $\alpha \geq 0$  small enough that  $\mathbf{x} + \alpha\mathbf{u} \in \text{dom}(f)$ . By the second-order Taylor expansion, there exists  $\gamma \in [0, 1]$  such that

$$f(\mathbf{x} + \alpha\mathbf{u}) = f(\mathbf{x}) + \alpha \nabla f(\mathbf{x})^\top \mathbf{u} + \frac{1}{2} \alpha^2 \mathbf{u}^\top \nabla^2 f(\mathbf{x} + \gamma\alpha\mathbf{u}) \mathbf{u}. \quad (4)$$

**Forward direction**  $\Rightarrow$ . Assume  $f$  is strongly convex with modulus  $m$ . Then by the characterization (2),

$$f(\mathbf{x} + \alpha \mathbf{u}) \geq f(\mathbf{x}) + \alpha \nabla f(\mathbf{x})^\top \mathbf{u} + \frac{m}{2} \alpha^2 \|\mathbf{u}\|^2.$$

Subtract (4) and cancel the linear terms:

$$\frac{1}{2} \alpha^2 \mathbf{u}^\top \nabla^2 f(\mathbf{x} + \gamma \alpha \mathbf{u}) \mathbf{u} \geq \frac{m}{2} \alpha^2 \|\mathbf{u}\|^2.$$

Dividing by  $\frac{1}{2} \alpha^2$ ,

$$\mathbf{u}^\top \nabla^2 f(\mathbf{x} + \gamma \alpha \mathbf{u}) \mathbf{u} \geq m \|\mathbf{u}\|^2.$$

Let  $\alpha \downarrow 0$ . By continuity of the Hessian,

$$\mathbf{u}^\top \nabla^2 f(\mathbf{x}) \mathbf{u} \geq m \|\mathbf{u}\|^2, \quad \forall \mathbf{u},$$

i.e.  $\nabla^2 f(\mathbf{x}) - m\mathbf{I} \succeq 0$  for all  $\mathbf{x}$ .

**Backward direction**  $\Leftarrow$ . Assume  $\nabla^2 f(\mathbf{x}) \succeq m\mathbf{I}$  for all  $\mathbf{x}$ . Then for any  $\mathbf{x}, \mathbf{z} \in \text{dom}(f)$ , apply (4) with  $\alpha \mathbf{u} := \mathbf{z} - \mathbf{x}$ :

$$f(\mathbf{z}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{z} - \mathbf{x}) + \frac{1}{2} (\mathbf{z} - \mathbf{x})^\top \nabla^2 f(\mathbf{x} + \gamma(\mathbf{z} - \mathbf{x})) (\mathbf{z} - \mathbf{x}).$$

Using  $\nabla^2 f(\cdot) \succeq m\mathbf{I}$ ,

$$\frac{1}{2} (\mathbf{z} - \mathbf{x})^\top \nabla^2 f(\mathbf{x} + \gamma(\mathbf{z} - \mathbf{x})) (\mathbf{z} - \mathbf{x}) \geq \frac{m}{2} \|\mathbf{z} - \mathbf{x}\|^2.$$

Hence

$$f(\mathbf{z}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{z} - \mathbf{x}) + \frac{m}{2} \|\mathbf{z} - \mathbf{x}\|^2,$$

which is the strong convexity condition (2). Thus  $f$  is strongly convex with modulus  $m$ .

### Example

- Quadratic:

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} + \mathbf{g}^\top \mathbf{x} + d, \quad \nabla f(\mathbf{x}) = \mathbf{H} \mathbf{x} + \mathbf{g}, \quad \nabla^2 f(\mathbf{x}) = \mathbf{H}.$$

Then

$$f \text{ convex} \iff \mathbf{H} \succeq \mathbf{0}, \quad f \text{ strongly convex (modulus } m) \iff \mathbf{H} \succeq m\mathbf{I}.$$

- Exponential:

$$f(x) = e^x, \quad f''(x) = e^x > 0 \Rightarrow f \text{ convex}.$$

- Reciprocal:

$$f(x) = \begin{cases} \frac{1}{x}, & x > 0, \\ \infty, & x \leq 0, \end{cases} \quad f''(x) = \frac{2}{x^3} > 0 \Rightarrow f \text{ convex}.$$

- Negative log:

$$f(x) = \begin{cases} -\ln x, & x > 0, \\ \infty, & x \leq 0, \end{cases} \quad f''(x) = \frac{1}{x^2} > 0 \Rightarrow f \text{ convex.}$$

- Non-differentiable example:  $f(x) = \|x\|_{\square}$  is convex, but the second-derivative test does not apply directly.

### 7.2.4 Norm of Matrices

Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , the operator 2-norm is defined as

$$\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})}.$$

#### Fact

- $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F \leq \sqrt{\min(m, n)} \|\mathbf{A}\|_2$ , where

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A(i, j)^2}$$

is the Frobenius norm.

- For a  $\mathcal{C}^2$  function  $f$ ,  $f$  is  $L$ -smooth if and only if

$$\|\nabla^2 f(\mathbf{x})\|_2 \leq L \quad \forall \mathbf{x}.$$

- If  $\mathbf{A} \in \mathbb{S}^n$  (symmetric), then

$$\|\mathbf{A}\|_2 = \max\{|\lambda_{\min}(\mathbf{A})|, |\lambda_{\max}(\mathbf{A})|\}.$$

## LEC 8

### 8.1 Linear Algebra Facts

1. For a matrix  $\mathbf{A}$ , define the operator 2-norm:

$$\|\mathbf{A}\|_2 := \left( \lambda_{\max}(\mathbf{A}^\top \mathbf{A}) \right)^{\frac{1}{2}}.$$

For  $\mathbf{A} \in \mathbb{S}^n$ ,

$$\|\mathbf{A}\|_2 = \max(|\lambda_{\min}(\mathbf{A})|, |\lambda_{\max}(\mathbf{A})|).$$

For  $\mathbf{A} \in \mathbb{S}_+^n$ ,

$$\|\mathbf{A}\|_2 = \lambda_{\max}(\mathbf{A}).$$

2. Recall: eigenvalues of  $\mathbf{A} - \lambda \mathbf{I}$  are the eigenvalues of  $\mathbf{A}$  shifted by  $-\lambda$ .

For  $\mathbf{A} \in \mathbb{S}^n$ ,

$$\mathbf{A} \succeq t\mathbf{I} \Leftrightarrow \lambda_i(\mathbf{A}) \geq t, \quad \forall i = 1, \dots, n,$$

$$\mathbf{A} \preceq t\mathbf{I} \Leftrightarrow \lambda_i(\mathbf{A}) \leq t, \quad \forall i = 1, \dots, n.$$

### 8.2 $m$ -strongly Convex Functions

#### 8.2.1 Two Lemmas

Given  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , with  $f \in \mathcal{C}^2$ ,

$$f \text{ is } L\text{-smooth} \Leftrightarrow -L\mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}, \quad \forall \mathbf{x} \in \mathbb{R}^n,$$

$$f \text{ is convex} \Leftrightarrow \nabla^2 f(\mathbf{x}) \succeq \mathbf{0}, \quad \forall \mathbf{x} \in \mathbb{R}^n,$$

$$f \text{ is } m\text{-strongly convex} \Leftrightarrow \nabla^2 f(\mathbf{x}) \succeq m\mathbf{I}, \quad \forall \mathbf{x} \in \mathbb{R}^n,$$

$$f \text{ is } L\text{-smooth, } m\text{-strongly convex} \Leftrightarrow m\mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

#### Example

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} + \mathbf{g}^\top \mathbf{x} + d, \quad \mathbf{H} \in \mathbb{S}^n.$$

Then

$$f \text{ is } L\text{-smooth} \Leftrightarrow -L\mathbf{I} \preceq \mathbf{H} \preceq L\mathbf{I},$$

$$f \text{ is convex} \Leftrightarrow \mathbf{H} \succeq \mathbf{0},$$

$$f \text{ is } m\text{-strongly convex} \Leftrightarrow \mathbf{H} \succeq m\mathbf{I}.$$

#### Lemma

$$f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\} \text{ is } m\text{-strongly convex} \Leftrightarrow g(\mathbf{x}) := f(\mathbf{x}) - \frac{m}{2} \|\mathbf{x}\|^2 \text{ is convex.}$$

**Proof**

Algebra shows  $\forall \lambda \in \mathbb{R}, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ,

$$\begin{aligned} & (1 - \lambda)g(\mathbf{x}) + \lambda g(\mathbf{y}) - g((1 - \lambda)\mathbf{x} + \lambda\mathbf{y}) \\ &= (1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y}) - f((1 - \lambda)\mathbf{x} + \lambda\mathbf{y}) - \frac{m}{2}\lambda(1 - \lambda)\|\mathbf{x} - \mathbf{y}\|^2. \end{aligned}$$

**Lemma**

If  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  is convex, then  $\exists \mathbf{a} \in \mathbb{R}^n, b \in \mathbb{R}$  such that

$$f(\mathbf{x}) \geq \underbrace{\mathbf{a}^\top \mathbf{x} + b}_{\text{Affine Underestimator}}, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

**Proof**

If  $f$  is differentiable, this is an immediate consequence of the subgradient inequality.

**8.2.2 Two Theorems**

Say  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is *coercive* if

$$\forall s \in \mathbb{R}, \exists r \in \mathbb{R} \text{ such that } \|\mathbf{x}\| \geq r \Rightarrow f(\mathbf{x}) \geq s.$$

Intuitively,  $f(\mathbf{x}) \rightarrow \infty$  as  $\|\mathbf{x}\| \rightarrow \infty$ .

**Facts**

- If  $f, g : \mathbb{R}^n \rightarrow \mathbb{R}, f \geq g$ , and  $g$  is coercive, then  $f$  is coercive.
- $q(t) = \frac{1}{2}at^2 + bt + c : \mathbb{R} \rightarrow \mathbb{R}$  is coercive whenever  $a > 0$ .

**Theorem**

If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuous and coercive, then  $f$  has a minimizer.

**Proof**

Choose  $r$  such that  $f(\mathbf{x}) \geq f(\mathbf{0})$  whenever  $\|\mathbf{x}\| \geq r$  (by coercivity). Observe that  $f$  restricted to  $\overline{B}(\mathbf{0}, r)$  has a minimizer by compactness. Call it  $\mathbf{x}^*$ . Then  $\mathbf{x}^*$  is a minimizer of  $f$  over  $\mathbb{R}^n$  since for  $\mathbf{x} \notin \overline{B}(\mathbf{0}, r)$ ,

$$f(\mathbf{x}) \geq f(\mathbf{0}) \geq f(\mathbf{x}^*).$$

**Theorem**

If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $m$ -strongly convex for  $m > 0$ , then  $f$  has a unique minimizer.

**Proof**

*Uniqueness.* Suppose  $\mathbf{x}_1 \neq \mathbf{x}_2$  are both minimizers. Then

$$f\left(\frac{\mathbf{x}_1 + \mathbf{x}_2}{2}\right) \leq \frac{1}{2}(f(\mathbf{x}_1) + f(\mathbf{x}_2)) - \underbrace{\frac{1}{8}m\|\mathbf{x}_1 - \mathbf{x}_2\|^2}_{>0},$$

contradicting minimality of  $\mathbf{x}_1, \mathbf{x}_2$ .

*Existence.* Let  $g(\mathbf{x}) = f(\mathbf{x}) - \frac{m}{2}\|\mathbf{x}\|^2$ , which is convex by Lemma 1. Let  $\mathbf{a}^\top \mathbf{x} + b$  be an affine underestimate of  $g$  by Lemma 2. Define  $h(\mathbf{x}) := \mathbf{a}^\top \mathbf{x} + b + \frac{m}{2}\|\mathbf{x}\|^2$ , an underestimate of  $f(\mathbf{x})$ . We will show the function  $h$  is coercive, which implies  $f$  is coercive, so  $f$  has a minimizer. Rewrite  $h$  by completing the square:

$$\begin{aligned} h(\mathbf{x}) &= \frac{m}{2} \left\| \mathbf{x} + \frac{\mathbf{a}}{m} \right\|^2 - \frac{1}{2m} \mathbf{a}^\top \mathbf{a} + b, \\ h(\mathbf{x}) - s &= \frac{m}{2} \left\| \mathbf{x} + \frac{\mathbf{a}}{m} \right\|^2 - \frac{1}{2m} \mathbf{a}^\top \mathbf{a} + b - s, \quad \forall s \in \mathbb{R} \\ &\geq \frac{m}{2} \left( \|\mathbf{x}\| - \frac{\|\mathbf{a}\|}{m} \right)^2 - \frac{1}{2m} \mathbf{a}^\top \mathbf{a} + b - s. \end{aligned}$$

This is a univariate quadratic in  $\|\mathbf{x}\|$  with positive leading coefficient, hence coercive. Thus

$$\exists r \text{ such that } \|\mathbf{x}\| \geq r \quad (r > 0) \quad \Rightarrow \quad h(\mathbf{x}) \geq s.$$

**8.3 Analysis of GD for  $m$ -strongly Convex Case****Theorem**

Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  is  $m$ -strongly convex,  $\mathbf{x}^* \in \text{intdom}(f)$ ,  $f$  is differentiable at  $\mathbf{x}^*$ , and  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  (so  $\mathbf{x}^*$  is a minimizer; we cannot deduce the converse since  $f$  may not be differentiable at  $\mathbf{x}^*$ ). Then for all  $\mathbf{x} \in \text{intdom}(f)$  where  $\nabla f(\mathbf{x})$  exists,

$$\begin{aligned} \text{a) } f(\mathbf{x}) - f(\mathbf{x}^*) &\leq \frac{\|\nabla f(\mathbf{x})\|^2}{2m}, \\ \text{b) } \|\mathbf{x} - \mathbf{x}^*\| &\leq \frac{2}{m} \|\nabla f(\mathbf{x})\|. \end{aligned}$$

**Proof**

For all  $\mathbf{y} \in \mathbb{R}^n$ ,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{m}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

Both sides are functions of  $\mathbf{y}$  (with  $\mathbf{x}$  fixed). In general, if  $\phi, \psi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  satisfy  $\phi(\mathbf{y}) \geq \psi(\mathbf{y})$  for all  $\mathbf{y}$ , then

$$\inf_{\mathbf{y}} \phi(\mathbf{y}) \geq \inf_{\mathbf{y}} \psi(\mathbf{y}).$$

Apply this to the preceding inequality:

$$f(\mathbf{x}^*) \geq \min_{\mathbf{y}} \left[ f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{m}{2} \|\mathbf{x} - \mathbf{y}\|^2 \right].$$

The RHS is quadratic in  $\mathbf{y}$ , with quadratic coefficient  $\frac{m}{2}\mathbf{I}$  and linear coefficient  $\nabla f(\mathbf{x}) - m\mathbf{x}$ , so its minimizer  $\mathbf{y}^*$  is

$$\mathbf{y}^* = -\frac{\nabla f(\mathbf{x}) - m\mathbf{x}}{m}.$$

(We regard  $\mathbf{a}^\top \mathbf{x}$  as having coefficient  $\mathbf{a}$ , not  $\mathbf{a}^\top$ .)

The minimal value is

$$\min_{\mathbf{y}} \text{RHS} = f(\mathbf{x}) - \frac{1}{2m} \|\nabla f(\mathbf{x})\|^2,$$

which rearranges to yield (a).

For (b),

$$\begin{aligned} f(\mathbf{x}^*) &\geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{x}^* - \mathbf{x}) + \frac{m}{2} \|\mathbf{x}^* - \mathbf{x}\|^2 \\ &\geq f(\mathbf{x}) - \|\nabla f(\mathbf{x})\| \cdot \|\mathbf{x}^* - \mathbf{x}\| + \frac{m}{2} \|\mathbf{x}^* - \mathbf{x}\|^2 \quad (\text{Cauchy-Schwarz}) \\ 0 \leq f(\mathbf{x}) - f(\mathbf{x}^*) &\leq \|\nabla f(\mathbf{x})\| \cdot \|\mathbf{x}^* - \mathbf{x}\| - \frac{m}{2} \|\mathbf{x}^* - \mathbf{x}\|^2 \\ \frac{m}{2} \|\mathbf{x}^* - \mathbf{x}\|^2 &\leq \|\nabla f(\mathbf{x})\| \cdot \|\mathbf{x}^* - \mathbf{x}\|. \end{aligned}$$

Divide by  $\|\mathbf{x}^* - \mathbf{x}\|$  to obtain (b).

### Theorem

Gradient Descent (GD) with stepsizes  $\alpha_k = \frac{1}{L}$  applied to  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  differentiable and  $m$ -strongly convex satisfies

$$f(\mathbf{x}^k) - f(\mathbf{x}^*) \leq \left(1 - \frac{m}{L}\right)^k \cdot (f(\mathbf{x}^0) - f(\mathbf{x}^*)), \quad \forall k = 1, 2, \dots$$

### Proof

$$\begin{aligned} f(\mathbf{x}^{k+1}) &\leq f(\mathbf{x}^k) - \frac{1}{2L} \|\nabla f(\mathbf{x}^k)\|^2 \\ &\leq f(\mathbf{x}^k) - \frac{m}{L} (f(\mathbf{x}^k) - f(\mathbf{x}^*)). \end{aligned}$$

Subtract  $f(\mathbf{x}^*)$  from both sides and rearrange:

$$f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*) \leq \left(1 - \frac{m}{L}\right) (f(\mathbf{x}^k) - f(\mathbf{x}^*)).$$

Induction on  $k$  proves the theorem:

$$f(\mathbf{x}^k) - f(\mathbf{x}^*) \leq \left(1 - \frac{m}{L}\right)^k (f(\mathbf{x}^0) - f(\mathbf{x}^*)).$$

This is exponential in  $k$ , called *Linear Convergence*. Compared to the convex (not strongly convex) case:

$$f(\mathbf{x}^k) - f(\mathbf{x}^*) \leq \frac{L\|\mathbf{x}^0 - \mathbf{x}^*\|^2}{2k},$$

which is called *Sublinear Convergence*, worse than for  $m$ -strongly convex functions.

## 8.4 Accelerated Gradient Descent and Momentum

*GD* is a *first-order* method, i.e., each iteration requires first derivatives of the objective.

Nemirovsky and Yudin (1982) found a first-order method such that

$$f(\mathbf{x}^k) - f(\mathbf{x}^*) \leq \mathcal{O}\left(\frac{1}{k^2}\right).$$

Furthermore, under some assumptions,  $\mathcal{O}\left(\frac{1}{k^2}\right)$  is the best possible.

The *N-Y algorithm* is difficult to implement. Nesterov (1983) proposed *Accelerated Gradient Descent* (*AGD*) which attains  $\mathcal{O}\left(\frac{1}{k^2}\right)$  (*Nesterov's Fast Gradient Method*).



## LEC 9

## 9.1 Algorithm

An  $L$ -smooth convex function can be minimized with accuracy

$$\mathcal{O}\left(\frac{1}{k^2}\right)$$

after  $k$  iterations.

Choose an initial point  $\mathbf{y}^0 \in \mathbb{R}^n$  arbitrarily. For  $k = 0, 1, 2, \dots$ , update as follows:

$$\begin{aligned}\mathbf{x}^{k+1} &= \mathbf{y}^k - \alpha_k \nabla f(\mathbf{y}^k), \\ \mathbf{y}^{k+1} &= \mathbf{x}^{k+1} + \beta_{k+1}(\mathbf{x}^{k+1} - \mathbf{x}^k).\end{aligned}$$

Here we fix  $\alpha_k \equiv \frac{1}{L}$ , while  $\{\beta_k\}$  is a sequence to be determined.

## 9.2 Convergence Analysis

Assume  $f$  is  $L$ -smooth and convex, and that  $f$  has a minimizer  $\mathbf{x}^*$ . Denote  $f_{\min} := f(\mathbf{x}^*)$ . Define

$$v_k := f(\mathbf{x}^k) - f_{\min} + \frac{L}{2} \left\| \mathbf{x}^k - \mathbf{x}^* - \rho_{k-1}^2 (\mathbf{x}^{k-1} - \mathbf{x}^*) \right\|^2.$$

The sequence  $\{\rho_k\}$  is to be determined. We take  $\mathbf{x}^0 := \mathbf{y}^0$ .

Observe that  $v_k \geq 0$ ; moreover, if  $v_k = 0$  then  $\mathbf{x}^k$  is optimal.

From the definition of  $v_k$ ,

$$v_{k+1} = f(\mathbf{x}^{k+1}) - f_{\min} + \frac{L}{2} \left\| \mathbf{x}^{k+1} - \mathbf{x}^* - \rho_k^2 (\mathbf{x}^k - \mathbf{x}^*) \right\|^2.$$

## Fact

**Recall.** For gradient descent with step size  $1/L$  applied at  $\mathbf{y}^k$ ,

$$f(\mathbf{x}^{k+1}) \leq f(\mathbf{y}^k) - \frac{1}{2L} \|\nabla f(\mathbf{y}^k)\|^2.$$

**Reason.**

$$\begin{aligned}f(\mathbf{x}^{k+1}) &\leq f(\mathbf{y}^k) + \nabla f(\mathbf{y}^k)^\top (\mathbf{x}^{k+1} - \mathbf{y}^k) + \frac{L}{2} \|\mathbf{x}^{k+1} - \mathbf{y}^k\|^2 \\ &= f(\mathbf{y}^k) - \frac{1}{L} \|\nabla f(\mathbf{y}^k)\|^2 + \frac{L}{2} \left\| \frac{1}{L} \nabla f(\mathbf{y}^k) \right\|^2 \\ &= f(\mathbf{y}^k) - \frac{1}{2L} \|\nabla f(\mathbf{y}^k)\|^2.\end{aligned}$$

The orange inequality is exactly the  $L$ -smoothness condition, which is equivalent to the gradient Lipschitz property

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\| \leq L\|\mathbf{y} - \mathbf{x}\|.$$

Intuitively, this restricts the second-order growth of  $f$ . If  $f$  is twice differentiable, it corresponds to bounding the Hessian  $\nabla^2 f$  in operator norm, and the inequality can be seen as a Taylor expansion with a quadratic remainder controlled by  $L$ .

Substituting this into  $v_{k+1}$ , we obtain

$$\begin{aligned} v_{k+1} &\leq f(\mathbf{y}^k) - f_{\min} - \underbrace{\frac{1}{2L}\|\nabla f(\mathbf{y}^k)\|^2}_{\mathcal{A}} + \underbrace{\frac{L}{2}\left\|\mathbf{x}^{k+1} - \mathbf{x}^* - \rho_k^2(\mathbf{x}^k - \mathbf{x}^*)\right\|^2}_{\mathcal{B}} \\ &= f(\mathbf{y}^k) - f_{\min} - \mathcal{A} + \mathcal{B}. \end{aligned}$$

We use an uncommon splitting trick:

$$f(\mathbf{y}^k) - f_{\min} = \rho_k^2(f(\mathbf{y}^k) - f_{\min}) + (1 - \rho_k^2)(f(\mathbf{y}^k) - f_{\min}).$$

Substituting back,

$$v_{k+1} \leq \rho_k^2(f(\mathbf{y}^k) - f_{\min}) + (1 - \rho_k^2)(f(\mathbf{y}^k) - f_{\min}) - \mathcal{A} + \mathcal{B}.$$

#### Fact

**Recall (subgradient inequality).** For convex differentiable  $f$ ,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}).$$

**Reason.** By convexity, from a geometric point of view the graph of  $f$  lies above all its tangents.

Assume  $0 \leq \rho_k \leq 1$  for all  $k$ . By the subgradient inequality, we have

$$\begin{aligned} f(\mathbf{x}^k) &\geq f(\mathbf{y}^k) + \nabla f(\mathbf{y}^k)^\top (\mathbf{x}^k - \mathbf{y}^k), \\ f_{\min} = f(\mathbf{x}^*) &\geq f(\mathbf{y}^k) + \nabla f(\mathbf{y}^k)^\top (\mathbf{x}^* - \mathbf{y}^k). \end{aligned}$$

Plugging these inequalities into the previous relation yields

$$\begin{aligned}
v_{k+1} &\leq \rho_k^2 \left( f(\mathbf{x}^k) - f_{\min} - \nabla f(\mathbf{y}^k)^\top (\mathbf{x}^k - \mathbf{y}^k) \right) + (1 - \rho_k^2) (\nabla f(\mathbf{y}^k)^\top (\mathbf{y}^k - \mathbf{x}^*)) - \mathcal{A} + \mathcal{B} \\
&= \underbrace{\rho_k^2 \left( f(\mathbf{x}^k) - f_{\min} + \frac{L}{2} \|\mathbf{x}^k - \mathbf{x}^* - \rho_{k-1}^2 (\mathbf{x}^{k-1} - \mathbf{x}^*)\|^2 \right)}_{v_k} - \rho_k^2 \nabla f(\mathbf{y}^k)^\top (\mathbf{x}^k - \mathbf{y}^k) \\
&\quad + (1 - \rho_k^2) (\nabla f(\mathbf{y}^k)^\top (\mathbf{y}^k - \mathbf{x}^*)) - \rho_k^2 \cdot \frac{L}{2} \|\mathbf{x}^k - \mathbf{x}^* - \rho_{k-1}^2 (\mathbf{x}^{k-1} - \mathbf{x}^*)\|^2 - \mathcal{A} + \mathcal{B} \\
&= \rho_k^2 v_k - \rho_k^2 \nabla f(\mathbf{y}^k)^\top (\mathbf{x}^k - \mathbf{y}^k) + (1 - \rho_k^2) (\nabla f(\mathbf{y}^k)^\top (\mathbf{y}^k - \mathbf{x}^*)) \\
&\quad - \rho_k^2 \cdot \frac{L}{2} \|\mathbf{x}^k - \mathbf{x}^* - \rho_{k-1}^2 (\mathbf{x}^{k-1} - \mathbf{x}^*)\|^2 - \mathcal{A} + \mathcal{B} \\
&= \rho_k^2 v_k + \underbrace{\nabla f(\mathbf{y}^k)^\top (\mathbf{y}^k - \rho_k^2 \mathbf{x}^k - (1 - \rho_k^2) \mathbf{x}^*)}_{\mathcal{C}} \\
&\quad - \rho_k^2 \cdot \frac{L}{2} \|\mathbf{x}^k - \mathbf{x}^* - \rho_{k-1}^2 (\mathbf{x}^{k-1} - \mathbf{x}^*)\|^2 - \mathcal{A} + \mathcal{B} \\
&= \rho_k^2 v_k - \rho_k^2 \cdot \frac{L}{2} \|\mathbf{x}^k - \mathbf{x}^* - \rho_{k-1}^2 (\mathbf{x}^{k-1} - \mathbf{x}^*)\|^2 - \mathcal{A} + \mathcal{B} + \mathcal{C}.
\end{aligned}$$

Here

$$\begin{aligned}
\mathcal{A} &= \frac{1}{2L} \|\nabla f(\mathbf{y}^k)\|^2, \quad \mathcal{B} = \frac{L}{2} \left\| \mathbf{x}^{k+1} - \mathbf{x}^* - \rho_k^2 (\mathbf{x}^k - \mathbf{x}^*) \right\|^2, \\
\mathcal{C} &= \nabla f(\mathbf{y}^k)^\top (\mathbf{y}^k - \rho_k^2 \mathbf{x}^k - (1 - \rho_k^2) \mathbf{x}^*).
\end{aligned}$$

Using the update

$$\mathbf{x}^{k+1} := \mathbf{y}^k - \frac{1}{L} \nabla f(\mathbf{y}^k),$$

and the identity  $\|\mathbf{a} - \mathbf{b}\|^2 = \|\mathbf{a}\|^2 - 2\mathbf{a}^\top \mathbf{b} + \|\mathbf{b}\|^2$ , we have

$$\begin{aligned}
\underbrace{\frac{L}{2} \left\| \mathbf{x}^{k+1} - \mathbf{x}^* - \rho_k^2 (\mathbf{x}^k - \mathbf{x}^*) \right\|^2}_{\mathcal{B}} &= \frac{L}{2} \left\| \mathbf{y}^k - \frac{1}{L} \nabla f(\mathbf{y}^k) - \mathbf{x}^* - \rho_k^2 (\mathbf{x}^k - \mathbf{x}^*) \right\|^2 \\
&= \frac{L}{2} \left\| \mathbf{y}^k - \mathbf{x}^* - \rho_k^2 (\mathbf{x}^k - \mathbf{x}^*) \right\|^2 \\
&\quad - \underbrace{\nabla f(\mathbf{y}^k)^\top (\mathbf{y}^k - \mathbf{x}^* - \rho_k^2 (\mathbf{x}^k - \mathbf{x}^*))}_{\mathcal{C}} + \underbrace{\frac{1}{2L} \|\nabla f(\mathbf{y}^k)\|^2}_{\mathcal{A}}.
\end{aligned}$$

Thus

$$-\mathcal{A} + \mathcal{B} + \mathcal{C} = \frac{L}{2} \left\| \mathbf{y}^k - \mathbf{x}^* - \rho_k^2 (\mathbf{x}^k - \mathbf{x}^*) \right\|^2.$$

Therefore,

$$\begin{aligned} v_{k+1} &\leq \rho_k^2 v_k + \frac{L}{2} \left\| \mathbf{y}^k - \mathbf{x}^* - \rho_k^2 (\mathbf{x}^k - \mathbf{x}^*) \right\|^2 \\ &\quad - \rho_k^2 \cdot \frac{L}{2} \left\| \mathbf{x}^k - \mathbf{x}^* - \rho_{k-1}^2 (\mathbf{x}^{k-1} - \mathbf{x}^*) \right\|^2 \\ &= \rho_k^2 v_k + \frac{L}{2} \underbrace{\left( \left\| \mathbf{y}^k - \mathbf{x}^* - \rho_k^2 (\mathbf{x}^k - \mathbf{x}^*) \right\|^2 - \rho_k^2 \left\| \mathbf{x}^k - \mathbf{x}^* - \rho_{k-1}^2 (\mathbf{x}^{k-1} - \mathbf{x}^*) \right\|^2 \right)}_{\mathcal{S}}. \end{aligned}$$

We force  $\mathcal{S} = 0$  to specify  $\rho_k$ . Then  $v_{k+1} \leq \rho_k^2 v_k$ , and

$$\mathbf{y}^k - \mathbf{x}^* - \rho_k^2 (\mathbf{x}^k - \mathbf{x}^*) = \rho_k (\mathbf{x}^k - \mathbf{x}^* - \rho_{k-1}^2 (\mathbf{x}^{k-1} - \mathbf{x}^*)).$$

Recalling  $\mathbf{y}^k = (1 + \beta_k) \mathbf{x}^k - \beta_k \mathbf{x}^{k-1}$ , we obtain

$$(1 + \beta_k) \mathbf{x}^k - \beta_k \mathbf{x}^{k-1} - \mathbf{x}^* - \rho_k^2 (\mathbf{x}^k - \mathbf{x}^*) = \rho_k (\mathbf{x}^k - \mathbf{x}^*) - \rho_k \rho_{k-1}^2 (\mathbf{x}^{k-1} - \mathbf{x}^*).$$

Rearranging,

$$(1 + \beta_k - \rho_k^2) (\mathbf{x}^k - \mathbf{x}^*) - \beta_k (\mathbf{x}^{k-1} - \mathbf{x}^*) = \rho_k (\mathbf{x}^k - \mathbf{x}^*) - \rho_k \rho_{k-1}^2 (\mathbf{x}^{k-1} - \mathbf{x}^*).$$

Comparing coefficients, we deduce

$$1 + \beta_k - \rho_k^2 = \rho_k, \quad \beta_k = \rho_k \rho_{k-1}^2.$$

Substituting  $\beta_k$ ,

$$1 + \rho_k \rho_{k-1}^2 - \rho_k^2 = \rho_k,$$

so  $\rho_k$  is a root of the quadratic

$$\rho_k^2 + (1 - \rho_{k-1}^2) \rho_k - 1 = 0.$$

Hence the two roots satisfy

$$\rho_{k_1} + \rho_{k_2} = -(1 - \rho_{k-1}^2), \quad \rho_{k_1} \rho_{k_2} = -1.$$

Since  $v_{k+1} \leq \rho_k^2 v_k$ , by induction we may assume  $\rho_{k-1} \in [0, 1]$ . Thus the positive root lies in  $[0, 1]$ , while the negative root is  $\leq -1$ . Therefore, choosing  $\rho_k$  as the positive root gives  $\rho_k \in [0, 1]$ .

An equivalent way to rewrite the recursion is

$$(1 - \rho_{k-1}^2) \rho_k = 1 - \rho_k^2,$$

that is,

$$\rho_k = \frac{1 - \rho_k^2}{1 - \rho_{k-1}^2}.$$

Since  $v_{k+1} \leq \rho_k^2 v_k$ , by induction

$$v_k \leq \rho_{k-1}^2 \rho_{k-2}^2 \cdots \rho_1^2 v_1.$$

A telescoping argument on the right-hand side yields

$$v_k \leq \frac{(1 - \rho_{k-1}^2)^2}{(1 - \rho_0^2)^2} v_1.$$

Setting  $\rho_0 := 0$  gives

$$v_k \leq (1 - \rho_{k-1}^2)^2 v_1.$$

### Claim

$$1 - \rho_k^2 \leq \frac{2}{k+2}.$$

### Proof

We proceed by induction on  $k$ .

*Base case* ( $k = 0$ ). The inequality holds trivially.

*Induction step* ( $k \geq 1$ ). From the recurrence relation,

$$\rho_k = \frac{1 - \rho_k^2}{1 - \rho_{k-1}^2} \geq \frac{1 - \rho_k^2}{\frac{2}{k+1}} = \frac{k+1}{2} (1 - \rho_k^2).$$

Therefore

$$1 - \rho_k^2 \leq \frac{2}{k+1} \rho_k.$$

Let  $x := 1 - \rho_k^2$ . Then

$$x - \frac{2}{k+1} \sqrt{1-x} \leq 0,$$

and the left-hand side is an increasing function of  $x$ .

Test  $x = \frac{2}{k+2}$ :

$$\frac{2}{k+2} \stackrel{?}{\leq} \frac{2}{k+1} \sqrt{1 - \frac{2}{k+2}}.$$

Squaring both sides and dividing by 4 gives

$$\frac{1}{(k+2)^2} \stackrel{?}{\leq} \frac{1}{(k+1)^2} \left(1 - \frac{2}{k+2}\right),$$

equivalently,

$$(k+1)^2 \stackrel{?}{\leq} (k+2)k.$$

This inequality does not hold, so the condition implies

$$x < \frac{2}{k+2}.$$

Thus

$$1 - \rho_k^2 < \frac{2}{k+2},$$

which completes the induction.

Since  $v_k \leq (1 - \rho_{k-1}^2)^2 v_1$ , the claim implies

$$v_k \leq \frac{4}{(k+1)^2} v_1.$$

Because  $v_k \geq f(\mathbf{x}^k) - f_{\min}$ , we obtain

$$f(\mathbf{x}^k) - f_{\min} \leq \frac{4}{(k+1)^2} v_1,$$

which is the  $\mathcal{O}\left(\frac{1}{k^2}\right)$  convergence rate.

It remains to bound  $v_1$  in terms of the initial point:

$$v_1 = f(\mathbf{x}^1) - f(\mathbf{x}^*) + \frac{L}{2} \|\mathbf{x}^1 - \mathbf{x}^*\|^2.$$

Note that  $\mathbf{x}^1$  is obtained from  $\mathbf{x}^0$  by a plain gradient descent step.

#### Fact

**Recall.** For gradient descent with step size  $1/L$ ,

$$f(\mathbf{x}^{k+1}) \leq f(\mathbf{x}^*) + \frac{L}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2).$$

**Reason.**

$$\begin{aligned} f(\mathbf{x}^k) &\leq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^k)^\top (\mathbf{x}^k - \mathbf{x}^*) \\ f(\mathbf{x}^{k+1}) &\leq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^k)^\top (\mathbf{x}^k - \mathbf{x}^*) - \frac{1}{2L} \|\nabla f(\mathbf{x}^k)\|^2 \\ &\leq f(\mathbf{x}^*) + \frac{L}{2} (2(\mathbf{x}^k - \mathbf{x}^*)^\top \frac{1}{L} \nabla f(\mathbf{x}^k) - \|\frac{1}{L} \nabla f(\mathbf{x}^k)\|^2) \\ &\leq f(\mathbf{x}^*) + \frac{L}{2} (2(\mathbf{x}^k - \mathbf{x}^*)^\top (\mathbf{x}^{k+1} - \mathbf{x}^k) - \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2). \end{aligned}$$

The **purple** inequality is in the form of the law of cosines

$$\|\mathbf{b} - \mathbf{c}\|^2 = \|\mathbf{b}\|^2 + \|\mathbf{c}\|^2 - 2\mathbf{b}^\top \mathbf{c}.$$

Rearranging yields

$$f(\mathbf{x}^{k+1}) \leq f(\mathbf{x}^*) + \frac{L}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2).$$

Substituting  $k = 0$  and rearranging, we obtain

$$v_1 \leq \frac{L}{2} \|\mathbf{x}^0 - \mathbf{x}^*\|^2.$$

Therefore,

$$f(\mathbf{x}^k) - f_{\min} \leq \frac{2L}{(k+1)^2} \|\mathbf{x}^0 - \mathbf{x}^*\|^2,$$

which is the standard bound for Accelerated Gradient Descent *AGD*.

### 9.3 Implementation

A practical implementation of *AGD* is:

$$\rho_0 = 0,$$

$$\mathbf{y}^0 = \text{arbitrary},$$

for  $k = 0, 1, 2, \dots$  do

$$\mathbf{x}^{k+1} := \mathbf{y}^k - \frac{1}{L} \nabla f(\mathbf{y}^k),$$

$$\rho_{k+1} := \text{positive root of } \rho_{k+1}^2 + (1 - \rho_k^2)\rho_{k+1} - 1 = 0,$$

$$\beta_{k+1} := \rho_{k+1}\rho_k^2,$$

$$\mathbf{y}^{k+1} := \mathbf{x}^{k+1} + \beta_{k+1}(\mathbf{x}^{k+1} - \mathbf{x}^k).$$

## LEC 10

## 10.1 AGD in Strongly Convex Cases

AGD for  $L$ -smooth convex functions was discussed in the last lecture.

Here we consider AGD for  $L$ -smooth and  $m$ -strongly convex functions.

We use constant parameters

$$\alpha_k = \frac{1}{L}, \quad \beta_k = \frac{\sqrt{L} - \sqrt{m}}{\sqrt{L} + \sqrt{m}}, \quad \forall k.$$

This implicitly assumes  $L > m$ , where  $L$  is the upper bound on the curvature and  $m$  is the lower bound.

**Theorem**

For an  $L$ -smooth and  $m$ -strongly convex function  $f$  with minimizer  $\mathbf{x}^*$ , accelerated gradient descent with

$$\alpha_k = \frac{1}{L}, \quad \beta_k = \frac{\sqrt{L} - \sqrt{m}}{\sqrt{L} + \sqrt{m}}$$

satisfies

$$f(\mathbf{x}^k) - f(\mathbf{x}^*) \leq \left(1 - \sqrt{\frac{m}{L}}\right)^k \left[ f(\mathbf{x}^0) - f(\mathbf{x}^*) + \frac{m}{2} \|\mathbf{x}^0 - \mathbf{x}^*\|^2 \right].$$

In particular, the method has **linear convergence**.

This rate is faster than standard GD, whose linear rate is

$$\left(1 - \frac{m}{L}\right)^k.$$

Hence, in the strongly convex case, AGD provides a significant improvement over GD.

When  $\frac{m}{L} \ll 1$ , the *convex* (non-strongly-convex) version of AGD may actually perform better for moderate  $k$ . A version of AGD that simultaneously attains both the strongly-convex and merely-convex bounds is well known and can be found in Nesterov's book.

*Note.* Conjugate Gradient is not covered in this lecture.

## 10.2 Aside on Programming

From Problem Set 2:

```
function huber_minimize(A, y, tol, x0, L, tau)

    while ...
        % ← Gradient Descent main loop
        evaluate huber gradient

end
```



We would like to separate the Huber-specific code from the general *GD* algorithm.

*Tendency.*

- The `gradfunc()` represents the gradient of the objective, and it internally depends on `A`, `y`, and `tau`.

### Example

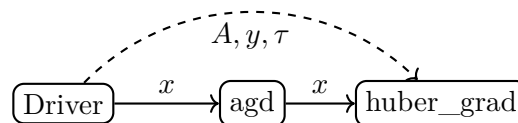
(Huber gradient wrapper.)

```
function g = huber_grad(x)
    % Evaluate Huber gradient
end
```

Then, we pass `huber_grad()` as the first argument of the optimizer:

```
function xopt = agd(gradfunc, x0, tol, L)
```

Now we have a call chain:



However, `huber_grad()` needs access to `A`, `y`, and `tau`, which are defined in `driver()` and not directly passed to it.

*Old solutions.*

(i) **Global variables.**

Store `A`, `y`, and `tau` in global variables known to both `driver()` and `huber_grad()`.

*Risk:* name clashes, hard to debug.

(ii) **Opaque parameter struct.**

Pass an additional argument `params` to `agd()`. The `driver()` packs `A`, `y`, and `tau` inside, and `huber_grad()` unpacks them.

*Risk:* troublesome packing and unpacking, error-prone.

*Modern solutions.*

(i) **Inner function.**

```
function driver
    A = ...;
    y = ...;
    tau = ...;

    function g = huber_grad(x)
        % uses A, y, tau
    end
```

```
xopt = agd(@huber_grad, x0, tol, L);  
end
```

(ii) **Anonymous function (closure / lambda).**

```
function g = huber_grad(x, A, y, tau)  
    % Evaluate Huber gradient  
end  
  
function driver  
    A = ...;  
    y = ...;  
    tau = ...;  
  
    gradfunc = @(x) huber_grad(x, A, y, tau);  
    xopt = agd(gradfunc, x0, tol, L);  
end
```

*Note.* For Problem Set 3, use one of these two modern methods (inner or anonymous functions). Or, besides the two modern methods, can also use *objects* to represent functions in Problem Set 3.

### 10.3 Binary Classification

Given  $N$  data points  $(\mathbf{a}_1, y_1), \dots, (\mathbf{a}_N, y_N)$  with  $y_i \in \{+1, -1\}$  and  $\mathbf{a}_i \in \mathbb{R}^n$  for all  $i$ , we want to learn a vector  $\mathbf{x}$  such that

$$\mathbf{a}_i^\top \mathbf{x} > 0 \text{ for } y_i = +1, \quad \mathbf{a}_i^\top \mathbf{x} < 0 \text{ for } y_i = -1.$$

Such an  $\mathbf{x}$  defines a **linear classifier**.

We can extend this idea to nonlinear classifiers by enriching the feature representation of each  $\mathbf{a}_i$ .

#### Example

**(Feature mapping.)** One approach is to use monomial features (as in Problem Set 1) to construct a higher-dimensional feature vector for each  $\mathbf{a}_i$ , on which a linear classifier is then learned.

To allow the classifier to represent a hyperplane not passing through the origin, we can append a constant 1 to each  $\mathbf{a}_i$ , which is known as the **padding trick**.

#### 10.3.1 Logistic Regression

We say the data is **linearly separable** if there exists  $\mathbf{x}$  such that  $\mathbf{a}_i^\top \mathbf{x} > 0$  for  $y_i = +1$  and  $\mathbf{a}_i^\top \mathbf{x} < 0$  for  $y_i = -1$ .

Define

$$\varphi(t) := \begin{cases} 1, & t > 0, \\ -1, & t < 0. \end{cases}$$

Define

$$L := \prod_{i:y_i=1} \frac{\varphi(\mathbf{a}_i^\top \mathbf{x}) + 1}{2} \prod_{i:y_i=-1} \frac{-\varphi(\mathbf{a}_i^\top \mathbf{x}) + 1}{2}.$$

Then  $L = 1$  if all points are correctly classified, and  $L = 0$  otherwise.

One idea is to find the classifier by maximizing  $L$  as a function of  $\mathbf{x}$ . This fails if the data points are not linearly separable, since  $L$  then cannot achieve 1.

To fix this, redefine  $\varphi$  (since the original one is discontinuous):

$$\varphi(t) := \frac{-1 + e^t}{1 + e^t}.$$

Use the same formula for  $L$ , but with this new  $\varphi$ . (It is easier to work with  $\ln L$ .)

$$\max_{\mathbf{x}} \left[ \sum_{i:y_i=1} \ln(\varphi(\mathbf{a}_i^\top \mathbf{x}) + 1) + \sum_{i:y_i=-1} \ln(-\varphi(\mathbf{a}_i^\top \mathbf{x}) + 1) \right].$$

This is called **Logistic Regression**.

Compute the logarithmic terms:

$$\ln(\varphi(t) + 1) = \ln\left(\frac{-1 + e^t + 1 + e^t}{1 + e^t}\right) = \ln\left(\frac{2e^t}{1 + e^t}\right) = \ln\left(\frac{2}{1 + e^{-t}}\right) = \ln 2 - \ln(1 + e^{-t}),$$

$$\ln(-\varphi(t) + 1) = \ln 2 - \ln(1 + e^t).$$

Dropping additive constants, logistic regression becomes

$$\min_{\mathbf{x}} \left[ \sum_{i:y_i=1} \ln(1 + e^{-\mathbf{a}_i^\top \mathbf{x}}) + \sum_{i:y_i=-1} \ln(1 + e^{\mathbf{a}_i^\top \mathbf{x}}) \right].$$

#### Fact

The logistic regression objective

$$\mathbf{x} \mapsto \sum_{i:y_i=1} \ln(1 + e^{-\mathbf{a}_i^\top \mathbf{x}}) + \sum_{i:y_i=-1} \ln(1 + e^{\mathbf{a}_i^\top \mathbf{x}})$$

is convex and  $L$ -smooth (for some  $L$  depending on the data).

However, if the data is linearly separable, then the problem has *no* minimizer, since any separating  $\mathbf{x}$  can be scaled arbitrarily large. Gradient descent then drives  $\|\mathbf{x}\| \rightarrow \infty$ .

To limit this, add a **regularizing term**:

$$\min_{\mathbf{x}} \frac{1}{N} \left[ \sum_{i:y_i=1} \cdots + \sum_{i:y_i=-1} \cdots \right] + \frac{\gamma}{2} \|\mathbf{x}\|^2.$$

Here  $\gamma > 0$  is the **regularization parameter**, and this is known as  $\ell_2$  regularization.

Later we will cover  $\ell_1$  regularization.

Now we can minimize this regularized objective using *GD* or *AGD*.

But what if  $N = 10^9$ ? Gradient evaluation is very expensive, and not all data points are needed to obtain a good  $\mathbf{x}$ .

### 10.3.2 Stochastic Gradient Descent *SGD*

Rewrite the regularized objective as

$$\mathbb{E}[l(\mathbf{a}, y; \mathbf{x})] + \frac{\gamma}{2} \|\mathbf{x}\|^2,$$

where the expectation is over a random data pair  $(\mathbf{a}, y)$  drawn uniformly from the dataset, and  $l$  is the logistic loss for a single sample.

In *SGD*, instead of using the full gradient of the expectation, we randomly select one data pair  $(\mathbf{a}_i, y_i)$  and compute the stochastic gradient

$$\nabla_{\mathbf{x}} \left[ l(\mathbf{a}_i, y_i; \mathbf{x}) + \frac{\gamma}{2} \|\mathbf{x}\|^2 \right].$$

At each iteration, a new random pair is chosen. This yields much cheaper iterations and scales to very large datasets.

## LEC 11

## 11.1 Stochastic Gradient Descent

Given a dataset  $(\mathbf{a}_i, y_i)$  for  $i = 1, \dots, N$ , we seek  $\mathbf{x}$  such that

$$\mathbf{a}_i^\top \mathbf{x} > 0 \quad \text{when } y_i = 1, \quad \mathbf{a}_i^\top \mathbf{x} < 0 \quad \text{when } y_i = -1.$$

Consider the regularized empirical risk

$$F(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{a}_i, y_i, \mathbf{x}) + \frac{\gamma}{2} \|\mathbf{x}\|^2.$$

**SGD algorithm (logistic regression setting).**

$\mathbf{x}^0$  arbitrary,

for  $k = 0, 1, 2, \dots$

Choose  $i_k \in \{1, \dots, N\}$  uniformly at random,

$$\mathbf{x}^{k+1} := \mathbf{x}^k - \alpha_k \nabla_{\mathbf{x}} \left( \ell(\mathbf{a}_{i_k}, y_{i_k}, \mathbf{x}) + \frac{\gamma}{2} \|\mathbf{x}\|^2 \right) \Big|_{\mathbf{x}=\mathbf{x}^k}.$$

**Finite-sum SGD viewpoint.** Suppose

$$F(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}; \xi_i), \quad \xi_i := (\mathbf{a}_i, y_i).$$

Let

$$g(\mathbf{x}; \xi_i) := \nabla_{\mathbf{x}} f(\mathbf{x}; \xi_i).$$

Then

$$\nabla F(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N g(\mathbf{x}; \xi_i).$$

**SGD algorithm.**

$\mathbf{x}^0$  arbitrary,

for  $k = 0, 1, 2, \dots$

Choose  $i_k \in \{1, \dots, N\}$  uniformly at random,

$$\mathbf{x}^{k+1} := \mathbf{x}^k - \alpha_k g(\mathbf{x}^k; \xi_{i_k}).$$

**Fully general expectation form.** In the most general setting,

$$F(\mathbf{x}) = \mathbb{E}[f(\mathbf{x}; \xi)],$$

where  $\xi$  is drawn from some distribution. On each iteration, we sample  $\xi$  from this distribution and define

$$g(\mathbf{x}; \xi) := \nabla_{\mathbf{x}} f(\mathbf{x}; \xi).$$

Then

$$\nabla F(\mathbf{x}) = \mathbb{E}[g(\mathbf{x}; \xi)].$$

Here  $g(\mathbf{x}; \xi)$  is called a **stochastic gradient**.

11.1.1 Convergence Theorem for *SGD*

## Theorem

**(Bubeck, Theorem 6.3, special case.)**Let  $F : \mathbb{R}^n \rightarrow \mathbb{R}$  be convex and  $L$ -smooth (with  $F$  the expectation or sum of the  $f(\cdot; \xi)$ ).

Assume:

- **Unbiasedness:**

$$\mathbb{E}[g(\mathbf{x}; \xi)] = \nabla F(\mathbf{x}), \quad \forall \mathbf{x}.$$

- **Variance bound:**

$$\mathbb{E}[\|g(\mathbf{x}; \xi) - \nabla F(\mathbf{x})\|^2] \leq \sigma^2, \quad \forall \mathbf{x}.$$

Let  $\mathbf{x}^*$  be a minimizer of  $F$ , and let

$$R := \|\mathbf{x}^0 - \mathbf{x}^*\|, \quad \eta := \frac{R}{\sigma\sqrt{\ell}}.$$

Run *SGD* for  $k = 0, 1, \dots, \ell - 1$  with the constant step size

$$\alpha_k \equiv \frac{1}{L + \frac{1}{\eta}}.$$

Then

$$\mathbb{E}\left[F\left(\frac{1}{\ell} \sum_{k=1}^{\ell} \mathbf{x}^k\right)\right] - F(\mathbf{x}^*) \leq \underbrace{\frac{R\sigma}{\sqrt{\ell}}}_{\text{stochastic part}} + \underbrace{\frac{LR^2}{2\ell}}_{\text{GD-like part (decays faster in } \ell)}, \quad \forall \ell = 1, 2, \dots$$

**Note.** The step size  $\alpha_k$  depends on the total number of iterations  $\ell$ .

## 11.1.2 Jensen's Inequality

## Fact

Let  $F : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  be convex. Suppose  $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^n$  and  $\lambda_1, \dots, \lambda_k \in \mathbb{R}$  satisfy

$$\lambda_i \geq 0, \quad \sum_{i=1}^k \lambda_i = 1.$$

Then

$$F\left(\sum_{i=1}^k \lambda_i \mathbf{x}_i\right) \leq \sum_{i=1}^k \lambda_i F(\mathbf{x}_i).$$

This can be proved by induction on  $k$ ; the case  $k = 2$  is exactly the definition of convexity.

## 11.1.3 Proof of the SGD Theorem

## Proof

We start from  $L$ -smoothness of  $F$ , which gives

$$F(\mathbf{x}^{k+1}) - F(\mathbf{x}^k) \leq \nabla F(\mathbf{x}^k)^\top (\mathbf{x}^{k+1} - \mathbf{x}^k) + \frac{L}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2.$$

Introduce the notation

$$\mathbf{g}_k := g(\mathbf{x}^k; \xi_{i_k}), \quad \mathbf{x}^{k+1} - \mathbf{x}^k = -\frac{1}{L + \frac{1}{\eta}} \mathbf{g}_k.$$

Then

$$\begin{aligned} F(\mathbf{x}^{k+1}) - F(\mathbf{x}^k) &= \nabla F(\mathbf{x}^k)^\top (\mathbf{x}^{k+1} - \mathbf{x}^k) + \frac{L}{2} \left( \frac{1}{L + \frac{1}{\eta}} \right)^2 \|\mathbf{g}_k\|^2 \\ &= \mathbf{g}_k^\top (\mathbf{x}^{k+1} - \mathbf{x}^k) + (\nabla F(\mathbf{x}^k) - \mathbf{g}_k)^\top (\mathbf{x}^{k+1} - \mathbf{x}^k) \\ &\quad + \frac{L}{2} \left( \frac{1}{L + \frac{1}{\eta}} \right)^2 \|\mathbf{g}_k\|^2. \end{aligned}$$

Applying Cauchy–Schwarz gives

$$\begin{aligned} F(\mathbf{x}^{k+1}) - F(\mathbf{x}^k) &\leq \mathbf{g}_k^\top (\mathbf{x}^{k+1} - \mathbf{x}^k) + \|\nabla F(\mathbf{x}^k) - \mathbf{g}_k\| \|\mathbf{x}^{k+1} - \mathbf{x}^k\| \\ &\quad + \frac{L}{2} \left( \frac{1}{L + \frac{1}{\eta}} \right)^2 \|\mathbf{g}_k\|^2. \end{aligned}$$

For any  $\eta > 0$  and  $a, b \in \mathbb{R}$  we have

$$\frac{1}{2} (\sqrt{\eta} a - \frac{1}{\sqrt{\eta}} b)^2 \geq 0 \quad \Rightarrow \quad ab \leq \frac{\eta}{2} a^2 + \frac{1}{2\eta} b^2.$$

Applying this to  $a = \|\nabla F(\mathbf{x}^k) - \mathbf{g}_k\|$  and  $b = \|\mathbf{x}^{k+1} - \mathbf{x}^k\|$ , we obtain

$$\begin{aligned}
 F(\mathbf{x}^{k+1}) - F(\mathbf{x}^k) &\leq \mathbf{g}_k^\top (\mathbf{x}^{k+1} - \mathbf{x}^k) + \frac{\eta}{2} \|\nabla F(\mathbf{x}^k) - \mathbf{g}_k\|^2 + \frac{1}{2\eta} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \\
 &\quad + \frac{L}{2} \left( \frac{1}{L + \frac{1}{\eta}} \right)^2 \|\mathbf{g}_k\|^2 \\
 &= \mathbf{g}_k^\top (\mathbf{x}^{k+1} - \mathbf{x}^k) + \frac{\eta}{2} \|\nabla F(\mathbf{x}^k) - \mathbf{g}_k\|^2 \\
 &\quad + \frac{1}{2} \left( L + \frac{1}{\eta} \right) \left( \frac{1}{L + \frac{1}{\eta}} \right)^2 \|\mathbf{g}_k\|^2 \\
 &= \mathbf{g}_k^\top (\mathbf{x}^{k+1} - \mathbf{x}^k) + \frac{\eta}{2} \|\nabla F(\mathbf{x}^k) - \mathbf{g}_k\|^2 + \frac{1}{2(L + \frac{1}{\eta})} \|\mathbf{g}_k\|^2.
 \end{aligned}$$

**Identity step.** Use the vector identity

$$2(b - a)^\top (a - c) = \|b - c\|^2 - \|a - c\|^2 - \|b - a\|^2.$$

Take

$$b := \mathbf{x}^k, \quad a := \mathbf{x}^{k+1}, \quad c := \mathbf{x}^*,$$

so that

$$\mathbf{x}^k - \mathbf{x}^{k+1} = \frac{1}{L + \frac{1}{\eta}} \mathbf{g}_k.$$

Then the left-hand side is

$$2 \cdot \frac{1}{L + \frac{1}{\eta}} \mathbf{g}_k^\top (\mathbf{x}^{k+1} - \mathbf{x}^*),$$

while the right-hand side is

$$\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 - \left( \frac{1}{L + \frac{1}{\eta}} \right)^2 \|\mathbf{g}_k\|^2.$$

Multiplying both sides by  $\frac{L + \frac{1}{\eta}}{2}$  and rearranging yields

$$\frac{1}{2(L + \frac{1}{\eta})} \|\mathbf{g}_k\|^2 = -\mathbf{g}_k^\top (\mathbf{x}^{k+1} - \mathbf{x}^*) + \frac{L + \frac{1}{\eta}}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2).$$

Substitute this into the previous inequality to eliminate the quadratic term:

$$\begin{aligned}
 F(\mathbf{x}^{k+1}) - F(\mathbf{x}^k) &\leq \mathbf{g}_k^\top (\mathbf{x}^{k+1} - \mathbf{x}^k) + \frac{\eta}{2} \|\nabla F(\mathbf{x}^k) - \mathbf{g}_k\|^2 \\
 &\quad - \mathbf{g}_k^\top (\mathbf{x}^{k+1} - \mathbf{x}^*) + \frac{L + \frac{1}{\eta}}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2) \\
 &= \mathbf{g}_k^\top (\mathbf{x}^* - \mathbf{x}^k) + \frac{\eta}{2} \|\nabla F(\mathbf{x}^k) - \mathbf{g}_k\|^2 \\
 &\quad + \frac{L + \frac{1}{\eta}}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2).
 \end{aligned}$$



Write

$$\mathbf{g}_k^\top (\mathbf{x}^* - \mathbf{x}^k) = \nabla F(\mathbf{x}^k)^\top (\mathbf{x}^* - \mathbf{x}^k) + (\mathbf{g}_k - \nabla F(\mathbf{x}^k))^\top (\mathbf{x}^* - \mathbf{x}^k),$$

so that

$$\begin{aligned} F(\mathbf{x}^{k+1}) - F(\mathbf{x}^k) &\leq \nabla F(\mathbf{x}^k)^\top (\mathbf{x}^* - \mathbf{x}^k) + (\mathbf{g}_k - \nabla F(\mathbf{x}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) \\ &\quad + \frac{\eta}{2} \|\nabla F(\mathbf{x}^k) - \mathbf{g}_k\|^2 \\ &\quad + \frac{L + \frac{1}{\eta}}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2). \end{aligned}$$

By convexity of  $F$ ,

$$\nabla F(\mathbf{x}^k)^\top (\mathbf{x}^* - \mathbf{x}^k) \leq F(\mathbf{x}^*) - F(\mathbf{x}^k),$$

so

$$\begin{aligned} F(\mathbf{x}^{k+1}) - F(\mathbf{x}^k) &\leq F(\mathbf{x}^*) - F(\mathbf{x}^k) + (\mathbf{g}_k - \nabla F(\mathbf{x}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) \\ &\quad + \frac{\eta}{2} \|\nabla F(\mathbf{x}^k) - \mathbf{g}_k\|^2 \\ &\quad + \frac{L + \frac{1}{\eta}}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2). \end{aligned}$$

Rearranging,

$$\begin{aligned} F(\mathbf{x}^{k+1}) - F(\mathbf{x}^*) &\leq (\mathbf{g}_k - \nabla F(\mathbf{x}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) + \frac{\eta}{2} \|\nabla F(\mathbf{x}^k) - \mathbf{g}_k\|^2 \\ &\quad + \frac{L + \frac{1}{\eta}}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2). \end{aligned}$$

**Taking expectations.** Take conditional expectation with respect to the randomness at step  $k$ , conditioning on the history so that  $\mathbf{x}^k$  is fixed. By the unbiasedness and variance assumptions,

$$\mathbb{E}[\mathbf{g}_k \mid \text{history}] = \nabla F(\mathbf{x}^k), \quad \mathbb{E}[\|\nabla F(\mathbf{x}^k) - \mathbf{g}_k\|^2 \mid \text{history}] \leq \sigma^2.$$

Thus

$$\mathbb{E}[(\mathbf{g}_k - \nabla F(\mathbf{x}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) \mid \text{history}] = 0,$$

and we obtain

$$\mathbb{E}[F(\mathbf{x}^{k+1}) \mid \text{history}] - F(\mathbf{x}^*) \leq \frac{\eta}{2} \sigma^2 + \frac{L + \frac{1}{\eta}}{2} (\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \mathbb{E}[\|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 \mid \text{history}]).$$

Now take full expectation and let  $R^2 := \|\mathbf{x}^0 - \mathbf{x}^*\|^2$ :

$$\mathbb{E}[F(\mathbf{x}^{k+1})] - F(\mathbf{x}^*) \leq \frac{\eta}{2} \sigma^2 + \frac{L + \frac{1}{\eta}}{2} (\mathbb{E}\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \mathbb{E}\|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2).$$

Average this inequality over  $k = 0, 1, \dots, \ell - 1$ :

$$\begin{aligned} \frac{1}{\ell} \sum_{k=0}^{\ell-1} (\mathbb{E}[F(\mathbf{x}^{k+1})] - F(\mathbf{x}^*)) &\leq \frac{\eta}{2} \sigma^2 + \frac{L + \frac{1}{\eta}}{2\ell} (\mathbb{E}\|\mathbf{x}^0 - \mathbf{x}^*\|^2 - \mathbb{E}\|\mathbf{x}^\ell - \mathbf{x}^*\|^2) \\ &\leq \frac{\eta}{2} \sigma^2 + \frac{L + \frac{1}{\eta}}{2\ell} R^2. \end{aligned}$$

By convexity of  $F$  and Jensen's inequality,

$$\mathbb{E} \left[ F \left( \frac{1}{\ell} \sum_{k=1}^{\ell} \mathbf{x}^k \right) \right] \leq \frac{1}{\ell} \sum_{k=1}^{\ell} \mathbb{E}[F(\mathbf{x}^k)],$$

so

$$\mathbb{E} \left[ F \left( \frac{1}{\ell} \sum_{k=1}^{\ell} \mathbf{x}^k \right) \right] - F(\mathbf{x}^*) \leq \frac{\eta}{2} \sigma^2 + \frac{L + \frac{1}{\eta}}{2\ell} R^2.$$

Finally, substitute  $\eta = \frac{R}{\sigma\sqrt{\ell}}$ :

$$\begin{aligned} \frac{\eta}{2} \sigma^2 + \frac{L + \frac{1}{\eta}}{2\ell} R^2 &= \frac{1}{2} \cdot \frac{R}{\sigma\sqrt{\ell}} \sigma^2 + \frac{1}{2\ell} \left( L + \frac{\sigma\sqrt{\ell}}{R} \right) R^2 \\ &= \frac{R\sigma}{2\sqrt{\ell}} + \frac{LR^2}{2\ell} + \frac{R\sigma}{2\sqrt{\ell}} \\ &= \frac{R\sigma}{\sqrt{\ell}} + \frac{LR^2}{2\ell}, \end{aligned}$$

which gives the desired bound.

### Remarks.

- The step size  $\alpha_k$  depends on  $\ell$ , the total number of iterations. Why not keep going after iteration  $\ell$ ? The convergence effectively *stalls* after  $\ell$  because each new step introduces a stochastic error of size on the order of

$$\frac{\eta}{2} \sigma^2 = \frac{1}{2} \cdot \frac{R}{\sigma\sqrt{\ell}} \sigma^2 = \frac{R\sigma}{2\sqrt{\ell}}.$$

- This is an *ergodic* bound (on the average iterate  $\frac{1}{\ell} \sum_{k=1}^{\ell} \mathbf{x}^k$ ), rather than a *last-iterate* bound. In practice, people often just keep the last iterate, but averaging is important for the theoretical guarantee.
- By increasing  $\ell$ , we can make the bound smaller; however, the variance term decays only as  $1/\sqrt{\ell}$ , unlike the  $1/\ell$  decay seen in deterministic gradient descent.

## LEC 12

12.1 Review of *SGD*

From last lecture: *SGD* convergence.

- *Ergodic* convergence vs. *last-iterate* convergence.
- Step size  $\alpha = \frac{1}{L + \frac{1}{\eta}}$  depends on the number of steps  $\ell$  (this implies  $\alpha$  cannot be arbitrarily large, otherwise the method diverges; the upper bound is  $\frac{1}{L}$ ).
- Larger  $\alpha$  means faster convergence to a poorer solution.
- Decrease step size: one can use, for example,  $\alpha_k = \frac{\text{const}}{\text{const} + k}$ .
- In practice: use a *stepsize schedule* (e.g.  $\alpha_k$  is a staircase function of iteration  $k$ ).

*SGD* and its variants are the main optimization algorithms used for training modern machine-learning models (including ChatGPT, Deepseek, ...).

12.2 *SGD* Variants12.2.1 *SGD* Variants in Theory

- SVRG,
- SAGA.

Both of these do variance reduction, mainly for convex problems.

12.2.2 *SGD* Variants in Practice

- **Minibatching**: e.g. select 8 random samples instead of 1 and compute the gradient of the sum of 8 terms.
- **Momentum** (similar to *AGD* momentum).
- **Coordinatewise scale**

$$\mathbf{x}^{k+1} := \mathbf{x}^k - \mathbf{D}^k \mathbf{g}^k,$$

where  $\mathbf{D}^k$  is a diagonal matrix adaptively updated. In ordinary *SGD*,  $\mathbf{D}^k = \alpha_k \mathbf{I}$ .

An extreme case is

$$\mathbf{x}^{k+1} := \mathbf{x}^k - \alpha_k \cdot \text{sign}(\mathbf{g}^k).$$

- A well-known *SGD* variant is called **ADAM**, which combines momentum and coordinatewise scaling.

## 12.3 Constrained and Nonsmooth Analysis

### 12.3.1 Motivating Examples

**$\ell_1$ -regularized least squares.** Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{y} \in \mathbb{R}^m$ ,  $\gamma > 0$ ,

$$\ell_1\text{LS} : \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|^2 + \gamma \|\mathbf{x}\|_1,$$

where  $\|\mathbf{x}\|_1 = |x_1| + \dots + |x_n|$ . One can check the three properties of a norm are satisfied.

Compare to ridge regression:

$$\text{RR} : \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|^2 + \frac{\gamma}{2} \|\mathbf{x}\|_2^2.$$

- (i) Both  $\ell_1$ LS and RR encourage shrinkage, i.e. large  $\mathbf{x}$  is penalized.
- (ii)  $\ell_1$ LS has a second effect: *selection*.

A solution  $\mathbf{x}$  to  $\ell_1$ LS has many entries identically equal to 0. Thus it identifies features (columns of  $\mathbf{A}$ ) that are irrelevant to explaining  $\mathbf{y}$ . Closely related to  $\ell_1$ LS is Lasso regression:

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|^2 \quad \text{s.t.} \quad \|\mathbf{x}\|_1 \leq \rho.$$

Observe that

$$g(\mathbf{x}) := \|\mathbf{x}\|_1$$

is a convex function, but not differentiable. So  $\ell_1$ LS is convex and nonsmooth.

Therefore plain *GD* and *AGD* (which need gradients) are not directly applicable.

We can rewrite  $\ell_1$ LS as a smooth constrained problem by introducing auxiliary variables:

$$\min_{\mathbf{x}, \mathbf{t}} \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|^2 + \gamma \sum_{i=1}^n t_i \quad \text{s.t.} \quad t_i \geq x_i, \quad t_i \geq -x_i, \quad \forall i = 1, \dots, n.$$

These constraints are equivalent to  $t_i \geq |x_i|$ , but this latter form is nonsmooth.

To exactly reproduce  $\ell_1$ LS, we should constrain  $t_i = |x_i|$ . It is nevertheless OK to use  $t_i \geq |x_i|$ , because no optimizer of the above problem would have  $t_i > |x_i|$ , due to the form of the objective function.

**Support vector machines.** Binary classification: given  $(\mathbf{a}_1, y_1), \dots, (\mathbf{a}_N, y_N)$ , with  $\mathbf{a}_i \in \mathbb{R}^n$ ,  $y_i \in \{1, -1\}$ , seek  $\mathbf{x} \in \mathbb{R}^n$ ,  $\xi \in \mathbb{R}$  such that

$$\begin{aligned} \mathbf{a}_i^\top \mathbf{x} + \xi &> 0 & \text{if } y_i = 1, \\ \mathbf{a}_i^\top \mathbf{x} + \xi &< 0 & \text{if } y_i = -1. \end{aligned}$$

These are strict linear inequalities on  $(\mathbf{x}, \xi)$ . The solution set of a system of (non-strict) linear inequalities is called a *polyhedron* (generalization of polygons in 2D).

Idea of SVM: find  $(\mathbf{x}, \xi)$  to maximize the width between the +1 and -1 classes.

**Theorem**

Assume  $\|\mathbf{x}\| = 1$ . Given  $\mathbf{a} \in \mathbb{R}^n$ , the distance from  $\mathbf{a}$  to the hyperplane

$$\{\mathbf{v} : \mathbf{v}^\top \mathbf{x} + \xi = 0\}$$

is  $|\mathbf{a}^\top \mathbf{x} + \xi|$ .

**Proof**

Claim that the closest point on the plane to  $\mathbf{a}$  is

$$\mathbf{v}^* := \mathbf{a} - (\mathbf{a}^\top \mathbf{x} + \xi)\mathbf{x}.$$

Check that  $\mathbf{v}^*$  lies on the plane.

$$(\mathbf{v}^*)^\top \mathbf{x} + \xi = \mathbf{a}^\top \mathbf{x} - (\mathbf{a}^\top \mathbf{x} + \xi)\mathbf{x}^\top \mathbf{x} + \xi = \mathbf{a}^\top \mathbf{x} - (\mathbf{a}^\top \mathbf{x} + \xi) \cdot 1 + \xi = 0.$$

Check that the distance from  $\mathbf{a}$  to  $\mathbf{v}^*$  is as claimed:

$$\|\mathbf{a} - \mathbf{v}^*\| = \|\mathbf{a} - \mathbf{a} + (\mathbf{a}^\top \mathbf{x} + \xi)\mathbf{x}\| = |\mathbf{a}^\top \mathbf{x} + \xi| \cdot \|\mathbf{x}\| = |\mathbf{a}^\top \mathbf{x} + \xi|$$

since  $\|\mathbf{x}\| = 1$ .

Check that no other point on the plane is closer to  $\mathbf{a}$  than  $\mathbf{v}^*$ .

Let  $\mathbf{v}^1$  also lie on the plane, so  $(\mathbf{v}^1)^\top \mathbf{x} + \xi = 0$ . Then

$$(\mathbf{v}^1 - \mathbf{v}^*)^\top \mathbf{x} = 0.$$

Let  $\mathbf{p} := \mathbf{v}^1 - \mathbf{v}^*$ . Then

$$\begin{aligned} \|\mathbf{a} - \mathbf{v}^1\|^2 &= \|\mathbf{a} - \mathbf{a} + (\mathbf{a}^\top \mathbf{x} + \xi)\mathbf{x} - \mathbf{p}\|^2 \\ &= |\mathbf{a}^\top \mathbf{x} + \xi|^2 \|\mathbf{x}\|^2 - 2(\mathbf{a}^\top \mathbf{x} + \xi)\mathbf{x}^\top \mathbf{p} + \|\mathbf{p}\|^2 \\ &= |\mathbf{a}^\top \mathbf{x} + \xi|^2 + \|\mathbf{p}\|^2 \end{aligned}$$

because  $\mathbf{x}^\top \mathbf{p} = 0$ . So the distance is minimized when  $\mathbf{p} = \mathbf{0}$ , i.e. when  $\mathbf{v}^1 = \mathbf{v}^*$ .

Thus we can write the SVM optimization problem as

$$\begin{aligned} \text{(SVM-1)} \quad & \max_{\mathbf{x}, \xi} \min_{i=1, \dots, N} |\mathbf{a}_i^\top \mathbf{x} + \xi| \\ \text{s.t.} \quad & \mathbf{a}_i^\top \mathbf{x} + \xi \geq 0 \quad (y_i = 1), \\ & \mathbf{a}_i^\top \mathbf{x} + \xi \leq 0 \quad (y_i = -1), \\ & \|\mathbf{x}\| = 1. \end{aligned}$$

Assume for now that there exists a separating plane  $(\mathbf{x}, \xi)$  for which the above objective is positive. This implies the objective function is positive at the optimizer.

Introduce variables  $\mathbf{t}$  to remove the absolute value:

$$\begin{aligned} \text{(SVM-2)} \quad & \max_{\mathbf{x}, \xi, \mathbf{t}} \min_{i=1, \dots, N} t_i \\ \text{s.t.} \quad & \mathbf{a}_i^\top \mathbf{x} + \xi = t_i \quad (y_i = 1), \\ & \mathbf{a}_i^\top \mathbf{x} + \xi = -t_i \quad (y_i = -1), \\ & \|\mathbf{x}\| = 1, \quad \mathbf{t} \geq \mathbf{0}. \end{aligned}$$

This formulation gets rid of the absolute value in SVM-1.

### Claim

We obtain the same optimizer if we replace  $\|\mathbf{x}\| = 1$  with  $\|\mathbf{x}\| \leq 1$ :

$$\begin{aligned} \text{(SVM-3)} \quad & \max_{\mathbf{x}, \xi, \mathbf{t}} \min_{i=1, \dots, N} t_i \\ \text{s.t.} \quad & \mathbf{a}_i^\top \mathbf{x} + \xi = t_i \quad (y_i = 1), \\ & \mathbf{a}_i^\top \mathbf{x} + \xi = -t_i \quad (y_i = -1), \\ & \|\mathbf{x}\| \leq 1, \quad \mathbf{t} \geq \mathbf{0}. \end{aligned}$$

### Proof

Suppose  $(\mathbf{x}, \xi, \mathbf{t})$  is feasible for SVM-3 with  $\|\mathbf{x}\| < 1$ .

Assume  $\mathbf{x} \neq \mathbf{0}$  (indeed, if  $\mathbf{x} = \mathbf{0}$ , then the constraints force  $\xi = 0$  and hence the objective  $\min_i t_i = 0$ , which cannot be optimal since we assumed the optimal value is positive).

Define the rescaled variables

$$\mathbf{x}^1 := \frac{\mathbf{x}}{\|\mathbf{x}\|}, \quad \xi^1 := \frac{\xi}{\|\mathbf{x}\|}, \quad \mathbf{t}^1 := \frac{\mathbf{t}}{\|\mathbf{x}\|}.$$

One checks that  $(\mathbf{x}^1, \xi^1, \mathbf{t}^1)$  is feasible and

$$\min_i t_i^1 = \frac{1}{\|\mathbf{x}\|} \min_i t_i > \min_i t_i,$$

since  $\|\mathbf{x}\| < 1$ . Thus we obtain a strictly higher objective value, contradicting optimality of a solution with  $\|\mathbf{x}\| < 1$ . Therefore any optimizer of SVM-3 must satisfy  $\|\mathbf{x}\| = 1$ , and the optimizers of SVM-2 and SVM-3 coincide.

## LEC 13

## 13.1 Constrained and Nonsmooth Analysis

## 13.1.1 SVM Example

Binary classification: given  $(\mathbf{a}_1, y_1), \dots, (\mathbf{a}_N, y_N)$  with

$$\mathbf{a}_i \in \mathbb{R}^n, \quad y_i \in \{1, -1\},$$

seek  $\mathbf{x} \in \mathbb{R}^n, \xi \in \mathbb{R}$  such that

$$\begin{aligned} \mathbf{a}_i^\top \mathbf{x} + \xi &> 0 & \text{if } y_i = 1, \\ \mathbf{a}_i^\top \mathbf{x} + \xi &< 0 & \text{if } y_i = -1. \end{aligned}$$

$$(\text{SVM-1}) \quad \max_{\mathbf{x}, \xi} \min_{i=1, \dots, N} |\mathbf{a}_i^\top \mathbf{x} + \xi| \quad \text{s.t.} \quad \mathbf{a}_i^\top \mathbf{x} + \xi \geq 0 \ (y_i = 1), \ \mathbf{a}_i^\top \mathbf{x} + \xi \leq 0 \ (y_i = -1), \ \|\mathbf{x}\| = 1.$$

$$(\text{SVM-2}) \quad \max_{\mathbf{x}, \xi, \mathbf{t}} \min_{i=1, \dots, N} t_i \quad \text{s.t.} \quad \mathbf{a}_i^\top \mathbf{x} + \xi = t_i \ (y_i = 1), \ \mathbf{a}_i^\top \mathbf{x} + \xi = -t_i \ (y_i = -1), \ \|\mathbf{x}\| = 1, \ \mathbf{t} \geq \mathbf{0}.$$

$$(\text{SVM-3}) \quad \max_{\mathbf{x}, \xi, \mathbf{t}} \min_{i=1, \dots, N} t_i \quad \text{s.t.} \quad \mathbf{a}_i^\top \mathbf{x} + \xi = t_i \ (y_i = 1), \ \mathbf{a}_i^\top \mathbf{x} + \xi = -t_i \ (y_i = -1), \ \|\mathbf{x}\| \leq 1, \ \mathbf{t} \geq \mathbf{0}.$$

SVM-2 is nonconvex; SVM-3 is convex.

Note:  $\min(x, y)$  is not differentiable when  $x = y$ .

Let  $\tilde{t} = \min(t_1, \dots, t_N)$ . Then we can write

$$(\text{SVM-4}) \quad \max_{\mathbf{x}, \xi, \tilde{t}} \tilde{t} \quad \text{s.t.} \quad \mathbf{a}_i^\top \mathbf{x} + \xi \geq \tilde{t} \ (y_i = 1), \ \mathbf{a}_i^\top \mathbf{x} + \xi \leq -\tilde{t} \ (y_i = -1), \ \|\mathbf{x}\| \leq 1, \ \tilde{t} \geq 0.$$

This is equivalent to SVM-3: given a feasible point for SVM-4, define

$$t_i = \begin{cases} \mathbf{a}_i^\top \mathbf{x} + \xi, & y_i = 1, \\ -\mathbf{a}_i^\top \mathbf{x} - \xi, & y_i = -1. \end{cases}$$

Any  $\tilde{t} \leq \min(t_1, \dots, t_N)$  is feasible for SVM-4. Among these choices of  $\tilde{t}$ , SVM-4 selects the largest one, i.e.  $\tilde{t} = \min(t_1, \dots, t_N)$ .

Move the quadratic constraint into the objective: change variables

$$\tilde{\mathbf{x}} := \frac{\mathbf{x}}{\tilde{t}}, \quad \tilde{\xi} := \frac{\xi}{\tilde{t}},$$

which is valid since  $\tilde{t} > 0$  at the minimizer of SVM-4. Divide all constraints by  $\tilde{t}$  to obtain SVM-5:

$$(\text{SVM-5}) \quad \max_{\tilde{\mathbf{x}}, \tilde{\xi}, \tilde{t}} \tilde{t} \quad \text{s.t.} \quad \mathbf{a}_i^\top \tilde{\mathbf{x}} + \tilde{\xi} \geq 1 \ (y_i = 1), \ \mathbf{a}_i^\top \tilde{\mathbf{x}} + \tilde{\xi} \leq -1 \ (y_i = -1), \ \|\tilde{\mathbf{x}}\| \leq 1, \ \tilde{t} \geq 0.$$

Observe that

$$\max \tilde{t} \text{ s.t. } \|\tilde{t} \tilde{\mathbf{x}}\| \leq 1, \tilde{t} \geq 0$$

has the same optimizer (in  $\tilde{\mathbf{x}}$ ) as

$$\min \|\tilde{\mathbf{x}}\| \text{ s.t. } \|\tilde{t} \tilde{\mathbf{x}}\| \leq 1, \tilde{t} \geq 0.$$

Thus we arrive at

$$(\text{SVM-6}) \quad \min_{\tilde{\mathbf{x}}, \tilde{\xi}} \|\tilde{\mathbf{x}}\|^2 \quad \text{s.t.} \quad \mathbf{a}_i^\top \tilde{\mathbf{x}} + \tilde{\xi} \geq 1 \ (y_i = 1), \ \mathbf{a}_i^\top \tilde{\mathbf{x}} + \tilde{\xi} \leq -1 \ (y_i = -1).$$

Both SVM-6 and the smooth reformulation of  $\ell_1$ LS (with constraints) are examples of *quadratic programming* (QP).

#### Fact

QP refers to minimizing a quadratic objective function subject to linear equality and/or inequality constraints.

SVM-6 is called **hard-margin SVM**; it requires that the data have an affine linear separator (otherwise the problem is infeasible).

For nonseparable data, we use SVM-7, the **soft-margin SVM**, where  $\gamma > 0$  is a penalty parameter:

$$(\text{SVM-7}) \quad \min_{\mathbf{x}, \xi, \mathbf{s}} \frac{1}{2} \|\tilde{\mathbf{x}}\|^2 + \gamma \sum_{i=1}^N s_i$$

$$\text{s.t.} \quad \mathbf{a}_i^\top \mathbf{x} + \xi \geq 1 - s_i \ (y_i = 1), \quad \mathbf{a}_i^\top \mathbf{x} + \xi \leq -1 + s_i \ (y_i = -1), \quad \mathbf{s} \geq \mathbf{0}.$$

If  $\mathbf{s} = \mathbf{0}$ , we recover SVM-6. If  $s_i > 0$  for some  $i$ , then  $\mathbf{a}_i$  is misclassified. Thus the second term of the objective penalizes misclassification;  $\gamma \rightarrow \infty$  forces  $s_i = 0$  and recovers SVM-6. When  $\gamma \rightarrow 0$ , misclassifications are penalized very little, so the solution prioritizes a wider margin.

Many of the  $s_i$  are expected to be exactly zero at a typical solution. We can eliminate the  $s_i$  to obtain a nonsmooth equivalent problem.

If  $\mathbf{x}, \xi$  are fixed, the best choice for  $s_i$  when  $y_i = 1$  is either

$$s_i = 1 - \mathbf{a}_i^\top \mathbf{x} - \xi \quad \text{or} \quad s_i = 0,$$

so the optimal choice is

$$s_i = \max(0, 1 - \mathbf{a}_i^\top \mathbf{x} - \xi).$$

Write this as  $s_i = \Phi(\mathbf{a}_i^\top \mathbf{x} + \xi)$  where

$$\Phi(t) = \max(0, 1 - t).$$

For  $y_i = -1$ , the best choice is  $s_i = \Phi(-(\mathbf{a}_i^\top \mathbf{x} + \xi))$ . With this notation we obtain SVM-8:

$$(\text{SVM-8}) \quad \min_{\mathbf{x}, \xi} \frac{1}{2} \|\tilde{\mathbf{x}}\|^2 + \gamma \sum_{y_i=1} \Phi(\mathbf{a}_i^\top \mathbf{x} + \xi) + \gamma \sum_{y_i=-1} \Phi(-(\mathbf{a}_i^\top \mathbf{x} + \xi)).$$

Observe that  $\Phi$  is convex, so SVM-8 is a convex, nonsmooth, unconstrained problem. The function  $\Phi$  is called the *hinge loss*.



**Fact**

The logistic loss  $\ln(1 + e^{-t})$  from logistic regression is a smooth approximation of the hinge loss.

Differences between SVM-8 and logistic regression with  $\ell_2$  regularization:

- (i) In logistic regression the penalty parameter multiplies the quadratic term (a notational difference).
- (ii) Logistic regression uses a smooth loss; the hinge loss is nonsmooth.
- (iii) The quadratic term in logistic regression is  $\frac{1}{2}\gamma\|(\mathbf{x}, \xi)\|^2$ , whereas in SVM-8 the quadratic term is independent of  $\xi$ .

**13.1.2 General Case**

Consider the constrained optimization problem

$$\min f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{x} \in \Omega, \quad \Omega \subseteq \mathbb{R}^n.$$

Assume  $\Omega$  is closed and convex.

The **normal cone** to  $\Omega$  at  $\mathbf{x} \in \Omega$  is

$$N_{\Omega}(\mathbf{x}) = \left\{ \mathbf{v} \in \mathbb{R}^n : \mathbf{v}^{\top}(\mathbf{y} - \mathbf{x}) \leq 0, \quad \forall \mathbf{y} \in \Omega \right\}.$$

Geometrically, if  $\mathbf{x}$  lies on the boundary  $\partial\Omega$ , the normal cone  $N_{\Omega}(\mathbf{x})$  consists of all vectors that point “outward” and are normal to  $\Omega$  at  $\mathbf{x}$ . For an interior point,  $N_{\Omega}(\mathbf{x}) = \{\mathbf{0}\}$ . For a nonsmooth boundary point, the normal cone is generated by all limiting normals coming from adjacent smooth patches.

A **closed convex cone**  $C \subseteq \mathbb{R}^n$  satisfies:

- $C$  is closed;
- $\mathbf{0} \in C$ ;
- if  $\mathbf{x} \in C$  and  $\lambda \geq 0$ , then  $\lambda\mathbf{x} \in C$ ;
- if  $\mathbf{x}, \mathbf{y} \in C$ , then  $\mathbf{x} + \mathbf{y} \in C$ .

**Fact**

- Any closed convex cone is a convex set.
- If  $\Omega$  is a closed convex set and  $\mathbf{x} \in \Omega$ , then  $N_{\Omega}(\mathbf{x})$  is a closed convex cone.

**Theorem**

Suppose  $\Omega \subseteq \mathbb{R}^n$  is closed, nonempty, and convex, and  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  is convex and differentiable, with  $\Omega \subseteq \text{int}(\text{dom } f)$ . Then

$$\mathbf{x}^* \in \arg \min \{f(\mathbf{x}) : \mathbf{x} \in \Omega\} \iff \mathbf{x}^* \in \Omega \text{ and } -\nabla f(\mathbf{x}^*) \in N_{\Omega}(\mathbf{x}^*).$$

## Proof

**Backward direction** ( $\Leftarrow$ ).

Assume  $\mathbf{x}^* \in \Omega$  and  $-\nabla f(\mathbf{x}^*) \in N_\Omega(\mathbf{x}^*)$ . For any  $\mathbf{z} \in \Omega$ , the subgradient (gradient) inequality for convex  $f$  gives

$$f(\mathbf{z}) \geq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^\top (\mathbf{z} - \mathbf{x}^*).$$

Since  $-\nabla f(\mathbf{x}^*) \in N_\Omega(\mathbf{x}^*)$  and  $\mathbf{z} \in \Omega$ , we have

$$-\nabla f(\mathbf{x}^*)^\top (\mathbf{z} - \mathbf{x}^*) \leq 0,$$

so

$$f(\mathbf{z}) \geq f(\mathbf{x}^*).$$

Thus  $\mathbf{x}^*$  is a minimizer over  $\Omega$ .

**Forward direction** ( $\Rightarrow$ ).

Assume  $\mathbf{x}^*$  is a minimizer. Take any  $\mathbf{z} \in \Omega$ . For all  $\alpha \in [0, 1]$ , convexity of  $\Omega$  implies  $(1 - \alpha)\mathbf{x}^* + \alpha\mathbf{z} \in \Omega$ , and by optimality,

$$f(\mathbf{x}^*) \leq f((1 - \alpha)\mathbf{x}^* + \alpha\mathbf{z}) = f(\mathbf{x}^* + \alpha(\mathbf{z} - \mathbf{x}^*)).$$

By differentiability, write a first-order expansion with remainder:

$$f(\mathbf{x}^* + \alpha(\mathbf{z} - \mathbf{x}^*)) = f(\mathbf{x}^*) + \alpha \nabla f(\mathbf{x}^*)^\top (\mathbf{z} - \mathbf{x}^*) + \phi(\alpha(\mathbf{z} - \mathbf{x}^*)),$$

where

$$\frac{\phi(\alpha(\mathbf{z} - \mathbf{x}^*))}{\|\alpha(\mathbf{z} - \mathbf{x}^*)\|} \rightarrow 0 \quad \text{as } \alpha \rightarrow 0.$$

Thus

$$0 \leq \alpha \nabla f(\mathbf{x}^*)^\top (\mathbf{z} - \mathbf{x}^*) + \phi(\alpha(\mathbf{z} - \mathbf{x}^*)).$$

Divide by  $\alpha > 0$ :

$$0 \leq \nabla f(\mathbf{x}^*)^\top (\mathbf{z} - \mathbf{x}^*) + \frac{\phi(\alpha(\mathbf{z} - \mathbf{x}^*))}{\alpha}.$$

Let  $\alpha \downarrow 0$ ; the second term tends to 0, so

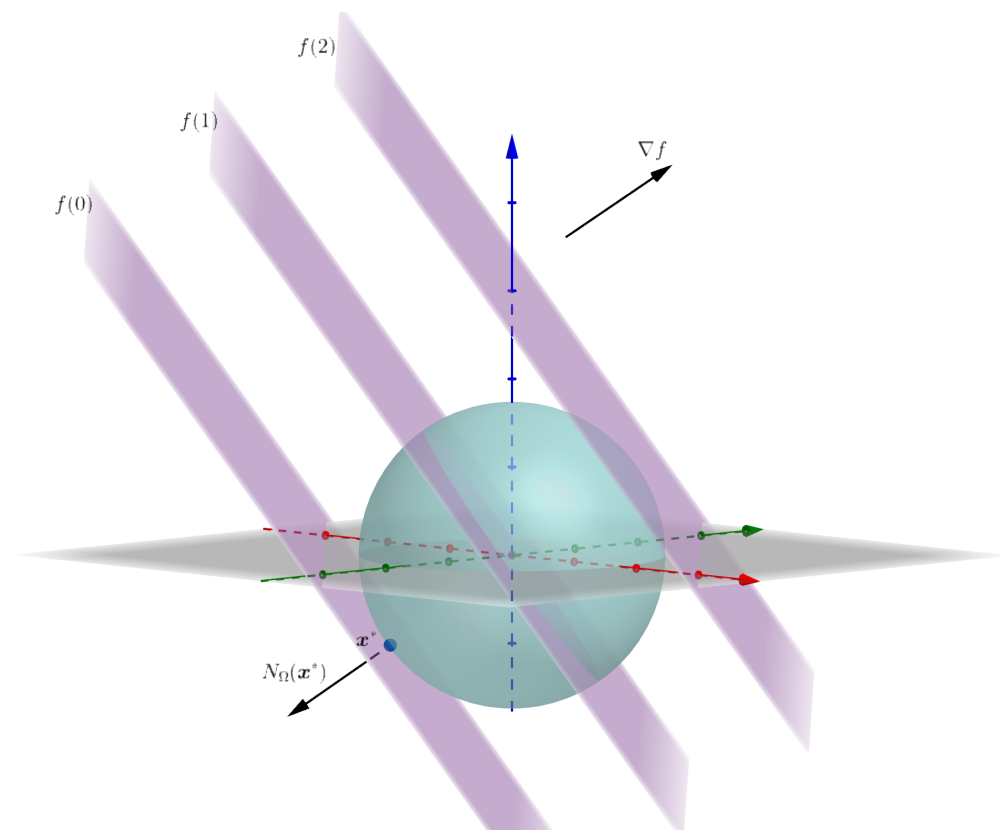
$$\nabla f(\mathbf{x}^*)^\top (\mathbf{z} - \mathbf{x}^*) \geq 0, \quad \forall \mathbf{z} \in \Omega.$$

Equivalently,

$$-\nabla f(\mathbf{x}^*)^\top (\mathbf{z} - \mathbf{x}^*) \leq 0, \quad \forall \mathbf{z} \in \Omega,$$

which means  $-\nabla f(\mathbf{x}^*) \in N_\Omega(\mathbf{x}^*)$ .

Moreover, if  $\lambda \geq 0$ , then  $-\lambda \nabla f(\mathbf{x}^*) \in N_\Omega(\mathbf{x}^*)$  as well, showing that  $N_\Omega(\mathbf{x}^*)$  is a cone.



*(This has a nice geometric interpretation and explains why the Lagrange multiplier method makes sense.)*

## LEC 14

## 14.1 Theorems with Normal Cones

## Theorem

**(From last lecture).** Suppose  $\Omega \subseteq \mathbb{R}^n$  is closed, nonempty, and convex, and  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  is convex and differentiable, with  $\Omega \subseteq \text{intdom}(f)$ . Then

$$\mathbf{x}^* \in \arg \min \{f(\mathbf{x}) : \mathbf{x} \in \Omega\} \iff \mathbf{x}^* \in \Omega \text{ and } -\nabla f(\mathbf{x}^*) \in N_{\Omega}(\mathbf{x}^*).$$

## Theorem

Let  $\Omega = \Omega_1 \cap \cdots \cap \Omega_m$ , where each  $\Omega_i$  is closed and convex (so  $\Omega$  is closed and convex). Then, for all  $\mathbf{x} \in \Omega$ ,

$$N_{\Omega_1}(\mathbf{x}) + \cdots + N_{\Omega_m}(\mathbf{x}) \subseteq N_{\Omega}(\mathbf{x}).$$

Here  $+$  denotes the **Minkowski sum**: given  $A, B \subseteq \mathbb{R}^n$ ,

$$A + B = \{\mathbf{a} + \mathbf{b} : \mathbf{a} \in A, \mathbf{b} \in B\}.$$

## Proof

Let  $\mathbf{v} = \mathbf{v}_1 + \cdots + \mathbf{v}_m$  with  $\mathbf{v}_i \in N_{\Omega_i}(\mathbf{x})$  for all  $i = 1, \dots, m$ .

By definition of the normal cone,

$$\mathbf{v}_i^\top (\mathbf{z} - \mathbf{x}) \leq 0, \quad \forall \mathbf{z} \in \Omega, \forall i = 1, \dots, m.$$

Summing these inequalities over  $i$  gives

$$\mathbf{v}^\top (\mathbf{z} - \mathbf{x}) = \sum_{i=1}^m \mathbf{v}_i^\top (\mathbf{z} - \mathbf{x}) \leq 0, \quad \forall \mathbf{z} \in \Omega.$$

Hence  $\mathbf{v} \in N_{\Omega}(\mathbf{x})$ .

In many generic situations, the inclusion in the theorem is actually an equality:

$$N_{\Omega_1}(\mathbf{x}) + \cdots + N_{\Omega_m}(\mathbf{x}) = N_{\Omega}(\mathbf{x}).$$

A condition that guarantees this equality is called a **constraint qualification** (CQ).

## Fact

**Constraint qualifications ensuring equality.**

- (CQ1) If  $\Omega_1, \dots, \Omega_m$  are all polyhedral sets, then

$$N_{\Omega_1}(\mathbf{x}) + \cdots + N_{\Omega_m}(\mathbf{x}) = N_{\Omega}(\mathbf{x}).$$

- (CQ2, Slater condition) If  $\text{int}(\Omega) \neq \emptyset$ , then the same equality holds.

## Example

**Example where equality fails.**Let  $n = 2$ ,  $m = 2$ , and

$$\Omega_1 = \{\mathbf{x} : x_2 \leq 0\}, \quad \Omega_2 = \{\mathbf{x} : x_2 \geq x_1^2\},$$

and set  $\Omega := \Omega_1 \cap \Omega_2 = \{\mathbf{0}\}$ .

Then

$$N_\Omega(\mathbf{0}) = \mathbb{R}^2.$$

However,

$$N_{\Omega_1}(\mathbf{0}) = \{\lambda(0, 1)^\top : \lambda \geq 0\}, \quad N_{\Omega_2}(\mathbf{0}) = \{\lambda(0, 1)^\top : \lambda \leq 0\}.$$

Thus

$$N_{\Omega_1}(\mathbf{0}) + N_{\Omega_2}(\mathbf{0}) = \{\lambda(0, 1)^\top : \lambda \in \mathbb{R}\} \subsetneq \mathbb{R}^2 = N_\Omega(\mathbf{0}).$$

Neither CQ1 nor CQ2 holds in this example. Intuitively, the normal cone summarizes the local behavior of  $\Omega$  via linear approximation; here the linear approximations of  $\Omega_1$  and  $\Omega_2$  fail to capture the “sharp” intersection at the origin.

## Theorem

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  be convex, and fix  $\beta \in \mathbb{R}$ . Define the level set

$$S = \{\mathbf{y} : f(\mathbf{y}) \leq \beta\}.$$

For any  $\mathbf{x} \in S$ :

- a)  $S$  is convex.
- b) If  $f(\mathbf{x}) < \beta$ , then  $N_S(\mathbf{x}) = \{\mathbf{0}\}$ .
- c) If  $f(\mathbf{x}) = \beta$  and  $\nabla f(\mathbf{x}) \neq \mathbf{0}$ , then

$$N_S(\mathbf{x}) = \{\lambda \nabla f(\mathbf{x}) : \lambda \geq 0\}.$$

## Proof

**(a) Convexity of  $S$ .**Let  $\mathbf{x}_1, \mathbf{x}_2 \in S$  and  $\lambda \in [0, 1]$ . By convexity of  $f$ ,

$$f((1 - \lambda)\mathbf{x}_1 + \lambda\mathbf{x}_2) \leq (1 - \lambda)f(\mathbf{x}_1) + \lambda f(\mathbf{x}_2) \leq (1 - \lambda)\beta + \lambda\beta = \beta.$$

Hence  $(1 - \lambda)\mathbf{x}_1 + \lambda\mathbf{x}_2 \in S$ , so  $S$  is convex.**(b) Interior point:**  $f(\mathbf{x}) < \beta \Rightarrow N_S(\mathbf{x}) = \{\mathbf{0}\}$ .Since  $f$  is convex, it is continuous on the interior of its domain. Thus there exists  $r > 0$  such that

$$f(\mathbf{y}) \leq \beta, \quad \forall \mathbf{y} \in \overline{\mathbb{B}}(\mathbf{x}, r),$$

so  $\overline{\mathbb{B}}(\mathbf{x}, r) \subseteq S$ .Assume, to get a contradiction, that  $\mathbf{v} \in N_S(\mathbf{x}) \setminus \{\mathbf{0}\}$ . Then, by definition of the normal cone,

$$\mathbf{v}^\top (\mathbf{x} + r \frac{\mathbf{v}}{\|\mathbf{v}\|} - \mathbf{x}) \leq 0,$$

i.e.

$$\frac{r \mathbf{v}^\top \mathbf{v}}{\|\mathbf{v}\|} = \frac{r \|\mathbf{v}\|^2}{\|\mathbf{v}\|} = r \|\mathbf{v}\| \leq 0,$$

which is impossible since  $r > 0$  and  $\mathbf{v} \neq \mathbf{0}$ . Hence  $N_S(\mathbf{x}) = \{\mathbf{0}\}$ .

**(c) Boundary point:**  $f(\mathbf{x}) = \beta$ ,  $\nabla f(\mathbf{x}) \neq \mathbf{0}$ .

First we show

$$\{\lambda \nabla f(\mathbf{x}) : \lambda \geq 0\} \subseteq N_S(\mathbf{x}).$$

Let  $\mathbf{y} \in S$ , so  $f(\mathbf{y}) \leq \beta = f(\mathbf{x})$ . By convexity and differentiability,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}).$$

Hence

$$0 \geq f(\mathbf{y}) - f(\mathbf{x}) \geq \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}),$$

so  $\nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \leq 0$  for all  $\mathbf{y} \in S$ , which means  $\nabla f(\mathbf{x}) \in N_S(\mathbf{x})$ , and therefore  $\lambda \nabla f(\mathbf{x}) \in N_S(\mathbf{x})$  for all  $\lambda \geq 0$ .

The converse inclusion

$$N_S(\mathbf{x}) \subseteq \{\lambda \nabla f(\mathbf{x}) : \lambda \geq 0\}$$

is more delicate and is omitted here.

### Theorem

Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$ , and define

$$\Omega = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{b}\}.$$

Then for every  $\mathbf{x} \in \Omega$ ,

$$N_\Omega(\mathbf{x}) = \text{Range}(\mathbf{A}^\top).$$

### Example

#### Geometric intuition.

Fix a point  $\mathbf{x} \in \Omega$ . Any feasible direction  $\mathbf{d}$  at  $\mathbf{x}$  must satisfy

$$\mathbf{A}(\mathbf{x} + \mathbf{d}) = \mathbf{b} \implies \mathbf{A}\mathbf{d} = \mathbf{0},$$

so  $\mathbf{d} \in \text{Null}(\mathbf{A})$ . By the fundamental theorem of linear algebra,

$$\text{Null}(\mathbf{A})^\perp = \text{Range}(\mathbf{A}^\top).$$

On the other hand, the normal cone  $N_\Omega(\mathbf{x})$  consists of all vectors orthogonal to every feasible direction  $\mathbf{d}$ , i.e.

$$N_\Omega(\mathbf{x}) = \text{Null}(\mathbf{A})^\perp = \text{Range}(\mathbf{A}^\top).$$

**Proof**

We prove both inclusions.

( $\subseteq$  **direction**). We first show

$$\text{Range}(\mathbf{A}^\top) \subseteq N_\Omega(\mathbf{x}).$$

Let  $\mathbf{p} \in \text{Range}(\mathbf{A}^\top)$ , so  $\mathbf{p} = \mathbf{A}^\top \mathbf{v}$  for some  $\mathbf{v} \in \mathbb{R}^m$ . For any  $\mathbf{w} \in \Omega$ ,

$$\mathbf{p}^\top(\mathbf{w} - \mathbf{x}) = \mathbf{v}^\top \mathbf{A}(\mathbf{w} - \mathbf{x}) = \mathbf{v}^\top(\mathbf{b} - \mathbf{b}) = 0.$$

Thus  $\mathbf{p}^\top(\mathbf{w} - \mathbf{x}) \leq 0$  for all  $\mathbf{w} \in \Omega$ , so  $\mathbf{p} \in N_\Omega(\mathbf{x})$ .

( $\supseteq$  **direction**). We next show

$$N_\Omega(\mathbf{x}) \subseteq \text{Range}(\mathbf{A}^\top).$$

Equivalently, we show

$$\mathbb{R}^n \setminus \text{Range}(\mathbf{A}^\top) \subseteq \mathbb{R}^n \setminus N_\Omega(\mathbf{x}).$$

Take  $\mathbf{p} \in \mathbb{R}^n \setminus \text{Range}(\mathbf{A}^\top)$ . By the fundamental theorem of linear algebra, we can write

$$\mathbf{p} = \mathbf{p}_1 + \mathbf{p}_2,$$

where

$$\mathbf{p}_1 \in \text{Range}(\mathbf{A}^\top), \quad \mathbf{p}_2 \in \text{Null}(\mathbf{A}), \quad \mathbf{p}_1^\top \mathbf{p}_2 = 0.$$

Since  $\mathbf{p} \notin \text{Range}(\mathbf{A}^\top)$ , we must have  $\mathbf{p}_2 \neq \mathbf{0}$ .

Note that  $\mathbf{x} + \mathbf{p}_2 \in \Omega$  because

$$\mathbf{A}(\mathbf{x} + \mathbf{p}_2) = \mathbf{A}\mathbf{x} + \mathbf{A}\mathbf{p}_2 = \mathbf{b} + \mathbf{0} = \mathbf{b}.$$

Then

$$\mathbf{p}^\top((\mathbf{x} + \mathbf{p}_2) - \mathbf{x}) = \mathbf{p}^\top \mathbf{p}_2 = (\mathbf{p}_1 + \mathbf{p}_2)^\top \mathbf{p}_2 = \mathbf{p}_2^\top \mathbf{p}_2 > 0.$$

Hence  $\mathbf{p}^\top(\mathbf{w} - \mathbf{x}) > 0$  for the feasible point  $\mathbf{w} = \mathbf{x} + \mathbf{p}_2$ , so  $\mathbf{p} \notin N_\Omega(\mathbf{x})$ . Thus  $N_\Omega(\mathbf{x}) \subseteq \text{Range}(\mathbf{A}^\top)$ .

## 14.2 Convex Programming

A **convex program (CP)** has the form

$$\min_{\mathbf{x}} f_0(\mathbf{x}) \quad \text{s.t.} \quad f_1(\mathbf{x}) \leq 0, \dots, f_m(\mathbf{x}) \leq 0, \quad \mathbf{A}\mathbf{x} = \mathbf{b},$$

where  $f_0, \dots, f_m$  are convex functions.

**Example**

**Examples of convex programming.**

- Least squares (LS).
- $\ell_1$ -regularized least squares ( $\ell_1$ LS).
- The SVM formulations SVM-3, SVM-4, SVM-6, SVM-7, SVM-8.

Let

$$\Omega = \{\mathbf{x} : f_1(\mathbf{x}) \leq 0, \dots, f_m(\mathbf{x}) \leq 0, \mathbf{Ax} = \mathbf{b}\},$$

the **feasible region** of convex programming.

We say the problem is **smooth** if  $f_0, \dots, f_m$  are all differentiable on an open set containing  $\Omega$ . (Examples: SVM-4, SVM-6, SVM-7.)

Recall SVM-3:

$$\max_{\mathbf{x}, \xi, \mathbf{t}} \min_{i=1, \dots, N} t_i \quad \text{s.t.} \quad \mathbf{a}_i^\top \mathbf{x} + \xi = t_i \ (y_i = 1), \quad \mathbf{a}_i^\top \mathbf{x} + \xi = -t_i \ (y_i = -1), \quad \|\mathbf{x}\| \leq 1, \ \mathbf{t} \geq \mathbf{0}.$$

Equivalently,

$$\min_{\mathbf{x}, \xi, \mathbf{t}} - \min_{i=1, \dots, N} t_i \quad \text{s.t. the same constraints,}$$

or

$$\min_{\mathbf{x}, \xi, \mathbf{t}} \max_{i=1, \dots, N} (-t_i) \quad \text{s.t. the same constraints.}$$

Here the function  $f(s, t) = \max(s, t)$  is convex but not differentiable when  $s = t$ .

### 14.2.1 KKT Conditions

Let  $\hat{\mathbf{x}} \in \Omega$  and assume

$$\nabla f_1(\hat{\mathbf{x}}) \neq \mathbf{0}, \dots, \nabla f_m(\hat{\mathbf{x}}) \neq \mathbf{0},$$

and that a suitable constraint qualification (polyhedrality, Slater, or similar) holds.

Then:

$$\hat{\mathbf{x}} \text{ is optimal for CP} \iff -\nabla f_0(\hat{\mathbf{x}}) \in N_\Omega(\hat{\mathbf{x}}).$$

Using the normal cone decomposition,

$$N_\Omega(\hat{\mathbf{x}}) = N_{\Omega_1}(\hat{\mathbf{x}}) + \dots + N_{\Omega_{m+1}}(\hat{\mathbf{x}}),$$

where

$$\Omega_i := \{\mathbf{x} : f_i(\mathbf{x}) \leq 0\}, \ i = 1, \dots, m, \quad \Omega_{m+1} := \{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}.$$

Let

$$C := \{i \in \{1, \dots, m\} : f_i(\hat{\mathbf{x}}) = 0\}$$

be the set of **active constraints** at  $\hat{\mathbf{x}}$ , and

$$I := \{1, \dots, m\} \setminus C$$

the set of **inactive constraints**.

Then the optimality condition becomes

$$-\nabla f_0(\hat{\mathbf{x}}) \in \sum_{i \in C} N_{\Omega_i}(\hat{\mathbf{x}}) + N_{\Omega_{m+1}}(\hat{\mathbf{x}}).$$

Using the previous theorem about level sets and the equality-constraint normal cone, this is equivalent to the existence of multipliers  $\lambda_i \geq 0$  for  $i \in C$  and  $\mathbf{v} \in \mathbb{R}^m$  such that

$$-\nabla f_0(\hat{\mathbf{x}}) = \sum_{i \in C} \lambda_i \nabla f_i(\hat{\mathbf{x}}) + \mathbf{A}^\top \mathbf{v}.$$

These are the (primal) **KKT conditions** in this smooth convex setting.



## LEC 15

## 15.1 KKT Conditions

Recall convex programming (CP):

$$\min_{\mathbf{x}} f_0(\mathbf{x}) \quad \text{s.t.} \quad f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \quad \mathbf{Ax} = \mathbf{b},$$

with

$$\Omega := \{\mathbf{x} : f_1(\mathbf{x}) \leq 0, \dots, f_m(\mathbf{x}) \leq 0, \mathbf{Ax} = \mathbf{b}\}.$$

We have the normal-cone optimality condition:

$$\begin{aligned} \hat{\mathbf{x}} \text{ is optimal for CP} &\iff -\nabla f_0(\hat{\mathbf{x}}) \in N_{\Omega}(\hat{\mathbf{x}}), \quad \hat{\mathbf{x}} \in \Omega \\ &\iff -\nabla f_0(\hat{\mathbf{x}}) \in N_{\Omega_1}(\hat{\mathbf{x}}) + \dots + N_{\Omega_{m+1}}(\hat{\mathbf{x}}), \quad \hat{\mathbf{x}} \in \Omega, \end{aligned}$$

where

$$\Omega_i := \{\mathbf{x} : f_i(\mathbf{x}) \leq 0\}, \quad i = 1, \dots, m, \quad \Omega_{m+1} := \{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\},$$

and we assume a suitable **constraint qualification** (CQ).

Let

$$C \subseteq \{1, \dots, m\}$$

be the set of **active constraints** at  $\hat{\mathbf{x}}$ :

$$i \in C \iff f_i(\hat{\mathbf{x}}) = 0.$$

Assume also

$$\nabla f_i(\hat{\mathbf{x}}) \neq \mathbf{0}, \quad i = 1, \dots, m.$$

Then

$$-\nabla f_0(\hat{\mathbf{x}}) \in \sum_{i \in C} N_{\Omega_i}(\hat{\mathbf{x}}) + N_{\Omega_{m+1}}(\hat{\mathbf{x}})$$

is equivalent to the existence of multipliers  $\lambda_i \geq 0$  ( $i \in C$ ) and  $\mathbf{v} \in \mathbb{R}^p$  such that

$$-\nabla f_0(\hat{\mathbf{x}}) = \sum_{i \in C} \lambda_i \nabla f_i(\hat{\mathbf{x}}) + \mathbf{A}^\top \mathbf{v}.$$

The condition above is called the **KKT condition** (Karush–Kuhn–Tucker) for convex programming.

## 15.2 Subdifferential

## 15.2.1 Definition

Recall the (gradient) subgradient inequality: if  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  is convex and differentiable at  $\mathbf{x} \in \text{dom}(f)$ , then for all  $\mathbf{y} \in \mathbb{R}^n$ ,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}).$$

For a nonsmooth convex function (e.g.  $f(x) = \max(x, 0)$  at  $x = 0$ ), the gradient may not exist but we can still have affine lower bounds.

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  be convex and  $\mathbf{x} \in \text{dom}(f)$ . A vector  $\mathbf{g} \in \mathbb{R}^n$  is called a **subgradient** of  $f$  at  $\mathbf{x}$  if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{y} \in \mathbb{R}^n.$$

The set of all subgradients at  $\mathbf{x}$  is the **subdifferential** of  $f$  at  $\mathbf{x}$ , denoted

$$\partial f(\mathbf{x}).$$

### 15.2.2 Relative Interior

For a convex set  $C \subseteq \mathbb{R}^n$ , the (usual) interior is

$$\text{int}(C) = \{\mathbf{x} \in C : \exists r > 0 \text{ s.t. } \mathbb{B}(\mathbf{x}, r) \subseteq C\}.$$

Given  $S \subseteq \mathbb{R}^n$ , its **affine hull** is

$$\text{aff}(S) = \left\{ \mathbf{z} : \exists \mathbf{s}_1, \dots, \mathbf{s}_m \in S, \exists \lambda_1, \dots, \lambda_m \text{ with } \sum_{i=1}^m \lambda_i = 1, \mathbf{z} = \sum_{i=1}^m \lambda_i \mathbf{s}_i \right\}.$$

Moreover, there exist  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$  such that

$$\text{aff}(S) = \{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}.$$

For a convex set  $S \subseteq \mathbb{R}^n$ , the **relative interior** is

$$\text{relint}(S) = \{\mathbf{x} \in S : \exists r > 0 \text{ s.t. } \mathbb{B}(\mathbf{x}, r) \cap \text{aff}(S) \subseteq S\}.$$

Intuitively, this is the interior of  $S$  relative to the (possibly lower-dimensional) affine space  $\text{aff}(S)$ .

For any nonempty convex set  $S \subseteq \mathbb{R}^n$ ,

$$\text{relint}(S) \neq \emptyset.$$

#### Fact

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  be convex. Then:

- (i) For all  $\mathbf{x} \in \text{dom}(f)$ , the subdifferential  $\partial f(\mathbf{x})$  is a closed, convex set.
- (ii) For all  $\mathbf{x} \in \text{relint}(\text{dom}(f))$ , the subdifferential is nonempty:

$$\partial f(\mathbf{x}) \neq \emptyset.$$

- (iii) If  $\mathbf{x} \in \text{int}(\text{dom}(f))$ , then  $\partial f(\mathbf{x})$  is bounded; hence it is convex, compact, and nonempty.
- (iv) If  $\mathbf{x} \in \text{int}(\text{dom}(f))$ , then  $f$  is differentiable at  $\mathbf{x}$  if and only if

$$\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}.$$

**Example****Why use relative interior?**

Consider

$$f(x) = \begin{cases} -\sqrt{1-x^2}, & x \in [-1, 1], \\ \infty, & \text{otherwise.} \end{cases}$$

Then  $\partial f(x) = \emptyset$  at the endpoints  $x = \pm 1$ , even though  $f$  is convex. These points lie on the boundary of  $\text{dom}(f)$ , outside the relative interior.

**Theorem**

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  be convex and  $\mathbf{x} \in \text{dom}(f)$ . Then

$$\mathbf{x} \in \arg \min f \iff \mathbf{0} \in \partial f(\mathbf{x}).$$

**Proof**

We have

$$\mathbf{x} \text{ is a global minimizer} \iff f(\mathbf{y}) \geq f(\mathbf{x}), \forall \mathbf{y} \in \mathbb{R}^n.$$

This is equivalent to

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{0}^\top (\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{y} \in \mathbb{R}^n,$$

which is exactly the condition  $\mathbf{0} \in \partial f(\mathbf{x})$ .

**Terminology.** A convex function  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  is called *proper* if  $\text{dom}(f) \neq \emptyset$ .

### 15.3 Epigraphs

For  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ , the **epigraph** of  $f$  is

$$\text{epi}(f) = \{(\mathbf{x}, y) \in \mathbb{R}^n \times \mathbb{R} : y \geq f(\mathbf{x})\}.$$

We say that  $f$  is **closed** if  $\text{epi}(f)$  is a closed set.

**Example**

Let

$$f(x) = \begin{cases} -\ln x, & x > 0, \\ \infty, & x \leq 0. \end{cases}$$

Here  $\text{dom}(f) = (0, \infty)$  is not closed, but  $\text{epi}(f)$  is closed. Conversely, a closed domain does not guarantee a closed epigraph.

**Fact**

A proper convex function with closed epigraph is often called **lower semicontinuous (lsc)** or a *closed convex function*. If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is (finite-valued) convex, then  $f$  is automatically

proper and closed.

### Theorem

A convex function  $f$  is closed if and only if each level set

$$\{\mathbf{x} : f(\mathbf{x}) \leq \beta\}, \quad \beta \in \mathbb{R},$$

is a closed set.

### Theorem

If  $f, g : \mathbb{R}^n \rightarrow \mathbb{R}$  are closed (finite-valued) convex functions, then their sum  $f + g$  is also closed.

### Theorem

Let  $f_1, f_2 : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  be convex, closed, and proper. Assume

$$\text{relint}(\text{dom}(f_1)) \cap \text{relint}(\text{dom}(f_2)) \neq \emptyset.$$

Then for all  $\mathbf{x} \in \text{dom}(f_1) \cap \text{dom}(f_2)$ ,

$$\partial f_1(\mathbf{x}) + \partial f_2(\mathbf{x}) = \partial(f_1 + f_2)(\mathbf{x}),$$

where  $+$  denotes the Minkowski sum of sets.

### Fact

The relative-interior intersection condition above acts as a type of **constraint qualification** for subdifferentials of sums.

## 15.4 Indicator Function

### Example

#### Indicator of an interval.

Let

$$f(x) = \begin{cases} 0, & x \in [0, 1], \\ \infty, & x \notin [0, 1]. \end{cases}$$

This is the **indicator function** of  $[0, 1]$ , denoted  $\mathbb{I}_{[0,1]}(x)$ .

### Claim

For the interval  $[0, 1]$ ,

$$\partial \mathbb{I}_{[0,1]}(1) = [0, \infty).$$

**Theorem**

Let  $\Omega \subseteq \mathbb{R}^n$  be nonempty, closed, and convex. The indicator function

$$\mathbb{I}_\Omega(\mathbf{x}) = \begin{cases} 0, & \mathbf{x} \in \Omega, \\ \infty, & \mathbf{x} \notin \Omega, \end{cases}$$

is proper, closed, and convex. Moreover, for all  $\mathbf{x} \in \Omega$ ,

$$\partial \mathbb{I}_\Omega(\mathbf{x}) = N_\Omega(\mathbf{x}),$$

the normal cone of  $\Omega$  at  $\mathbf{x}$ .

**Proof**

**Proper:** Since  $\Omega \neq \emptyset$ , we have  $\text{dom}(\mathbb{I}_\Omega) = \Omega \neq \emptyset$ .

**Convex:** For any  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$  and  $\lambda \in [0, 1]$ , convexity follows by checking cases depending on whether  $\mathbf{x}_1$  and  $\mathbf{x}_2$  lie in  $\Omega$ .

**Closed:** The epigraph is

$$\text{epi}(\mathbb{I}_\Omega) = \{(\mathbf{x}, y) : \mathbf{x} \in \Omega, y \geq 0\} = \Omega \times [0, \infty),$$

which is closed because  $\Omega$  is closed.

**Subdifferential equals normal cone.**

*First inclusion:*  $\partial \mathbb{I}_\Omega(\mathbf{x}) \subseteq N_\Omega(\mathbf{x})$ .

Let  $\mathbf{g} \in \partial \mathbb{I}_\Omega(\mathbf{x})$ . Then, by definition,

$$\mathbb{I}_\Omega(\mathbf{y}) \geq \mathbb{I}_\Omega(\mathbf{x}) + \mathbf{g}^\top(\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{y} \in \mathbb{R}^n.$$

For  $\mathbf{y} \in \Omega$ , both sides are finite and

$$0 = \mathbb{I}_\Omega(\mathbf{y}) \geq 0 + \mathbf{g}^\top(\mathbf{y} - \mathbf{x}),$$

so  $\mathbf{g}^\top(\mathbf{y} - \mathbf{x}) \leq 0$  for all  $\mathbf{y} \in \Omega$ , which means  $\mathbf{g} \in N_\Omega(\mathbf{x})$ .

*Second inclusion:*  $N_\Omega(\mathbf{x}) \subseteq \partial \mathbb{I}_\Omega(\mathbf{x})$ .

Take  $\mathbf{v} \in N_\Omega(\mathbf{x})$ , so

$$\mathbf{v}^\top(\mathbf{y} - \mathbf{x}) \leq 0, \quad \forall \mathbf{y} \in \Omega.$$

For  $\mathbf{y} \in \Omega$ ,

$$\mathbb{I}_\Omega(\mathbf{y}) = 0 \geq 0 + \mathbf{v}^\top(\mathbf{y} - \mathbf{x}) = \mathbb{I}_\Omega(\mathbf{x}) + \mathbf{v}^\top(\mathbf{y} - \mathbf{x}).$$

For  $\mathbf{y} \notin \Omega$ , we have  $\mathbb{I}_\Omega(\mathbf{y}) = \infty$ , so the inequality

$$\mathbb{I}_\Omega(\mathbf{y}) \geq \mathbb{I}_\Omega(\mathbf{x}) + \mathbf{v}^\top(\mathbf{y} - \mathbf{x})$$

holds trivially. Thus  $\mathbf{v} \in \partial \mathbb{I}_\Omega(\mathbf{x})$ .

Combining both inclusions gives  $\partial \mathbb{I}_\Omega(\mathbf{x}) = N_\Omega(\mathbf{x})$ .

## 15.5 Rederiving the Normal-Cone Optimality Condition

## Theorem

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  be differentiable on a neighborhood of  $\Omega$ , where  $\Omega \subseteq \mathbb{R}^n$  is closed and convex. Then

$$\mathbf{x} \in \arg \min \{f(\mathbf{x}) : \mathbf{x} \in \Omega\}$$

if and only if

$$-\nabla f(\mathbf{x}) \in N_{\Omega}(\mathbf{x}).$$

## Proof

We can rewrite the constrained problem as unconstrained:

$$\min_{\mathbf{x} \in \Omega} f(\mathbf{x}) \iff \min_{\mathbf{x}} (f(\mathbf{x}) + \mathbb{I}_{\Omega}(\mathbf{x})).$$

Thus

$$\mathbf{x} \in \arg \min \{f(\mathbf{x}) : \mathbf{x} \in \Omega\} \iff \mathbf{x} \in \arg \min \{f(\mathbf{x}) + \mathbb{I}_{\Omega}(\mathbf{x})\}.$$

By the subgradient optimality condition,

$$\mathbf{x} \text{ minimizer} \iff \mathbf{0} \in \partial(f + \mathbb{I}_{\Omega})(\mathbf{x}).$$

Assuming

$$\text{relint}(\text{dom } f) \cap \text{relint}(\text{dom } \mathbb{I}_{\Omega}) \neq \emptyset,$$

we can apply the subdifferential sum rule:

$$\partial(f + \mathbb{I}_{\Omega})(\mathbf{x}) = \partial f(\mathbf{x}) + \partial \mathbb{I}_{\Omega}(\mathbf{x}).$$

Since  $f$  is differentiable at  $\mathbf{x}$ ,

$$\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\},$$

and by the indicator theorem,

$$\partial \mathbb{I}_{\Omega}(\mathbf{x}) = N_{\Omega}(\mathbf{x}).$$

Thus

$$\mathbf{0} \in \partial(f + \mathbb{I}_{\Omega})(\mathbf{x}) \iff \mathbf{0} \in \{\nabla f(\mathbf{x})\} + N_{\Omega}(\mathbf{x}) \iff -\nabla f(\mathbf{x}) \in N_{\Omega}(\mathbf{x}).$$

## LEC 16

## 16.1 Composite Minimization

We consider the problem

$$\min_{\mathbf{x}} F(\mathbf{x}), \quad F(\mathbf{x}) := f(\mathbf{x}) + \psi(\mathbf{x}),$$

where

- $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex and  $L$ -smooth;
- $\psi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  is proper, closed, and convex.

**Example**

- SVM-8 (hinge loss + quadratic regularizer).
- $\ell_1$ -regularized least squares:

$$\ell_1 \text{LS} : \min_{\mathbf{x}} \underbrace{\frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2}_{f(\mathbf{x})} + \underbrace{\gamma \|\mathbf{x}\|_1}_{\psi(\mathbf{x})}.$$

**Theorem****Optimality condition for composite minimization.**

For  $F(\mathbf{x}) = f(\mathbf{x}) + \psi(\mathbf{x})$  as above,

$$\mathbf{x} \in \arg \min F \iff -\nabla f(\mathbf{x}) \in \partial\psi(\mathbf{x}).$$

**Proof**

Since  $F = f + \psi$  with  $f$  differentiable,

$$\partial F(\mathbf{x}) = \partial f(\mathbf{x}) + \partial\psi(\mathbf{x}) = \{\nabla f(\mathbf{x})\} + \partial\psi(\mathbf{x}).$$

By the subgradient optimality condition,

$$\mathbf{x} \in \arg \min F \iff \mathbf{0} \in \partial F(\mathbf{x}) \iff \mathbf{0} \in \{\nabla f(\mathbf{x})\} + \partial\psi(\mathbf{x}),$$

which is equivalent to  $-\nabla f(\mathbf{x}) \in \partial\psi(\mathbf{x})$ .

**Theorem****Uniqueness under strong convexity.**

Suppose, in addition, that  $f$  is  $m$ -strongly convex with  $m > 0$ . Then  $F = f + \psi$  has a unique minimizer.

**Proof**

A sum of a strongly convex function and a convex function is strongly convex. Thus  $F$  is strongly convex, so it is strictly convex and hence has at most one minimizer. Since  $F$  is proper, closed, and coercive (cf. the level-set argument below), a minimizer exists and must be unique.

**Claim**

For any  $\mathbf{x}_0 \in \text{dom}(\psi)$ , the level set

$$S := \{\mathbf{x} \in \mathbb{R}^n : F(\mathbf{x}) \leq F(\mathbf{x}_0)\}$$

is closed.

**Proof**

Both  $f$  and  $\psi$  are closed, so their sum  $F = f + \psi$  is closed. For a closed function, all sublevel sets  $\{\mathbf{x} : F(\mathbf{x}) \leq \beta\}$  are closed. Taking  $\beta = F(\mathbf{x}_0)$  gives that  $S$  is closed.

**Claim**

Choose  $\mathbf{x}_0 \in \text{ri}(\text{dom}(\psi))$  and let  $S$  be the level set above. Then  $S$  is bounded.

**Proof**

Because  $\mathbf{x}_0 \in \text{ri}(\text{dom}(\psi))$  and  $\psi$  is convex, we have  $\partial\psi(\mathbf{x}_0) \neq \emptyset$ . Choose some  $\mathbf{g} \in \partial\psi(\mathbf{x}_0)$ . Then for all  $\mathbf{x}$ ,

$$\psi(\mathbf{x}) \geq \psi(\mathbf{x}_0) + \mathbf{g}^\top(\mathbf{x} - \mathbf{x}_0).$$

Since  $f$  is  $m$ -strongly convex,

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^\top(\mathbf{x} - \mathbf{x}_0) + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_0\|^2.$$

Adding gives

$$F(\mathbf{x}) \geq F(\mathbf{x}_0) + (\nabla f(\mathbf{x}_0) + \mathbf{g})^\top(\mathbf{x} - \mathbf{x}_0) + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_0\|^2.$$

For  $\mathbf{x} \in S$ , we have  $F(\mathbf{x}) - F(\mathbf{x}_0) \leq 0$ , hence

$$0 \geq F(\mathbf{x}) - F(\mathbf{x}_0) \geq (\nabla f(\mathbf{x}_0) + \mathbf{g})^\top(\mathbf{x} - \mathbf{x}_0) + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_0\|^2.$$

Applying Cauchy–Schwarz,

$$(\nabla f(\mathbf{x}_0) + \mathbf{g})^\top(\mathbf{x} - \mathbf{x}_0) \geq -\|\nabla f(\mathbf{x}_0) + \mathbf{g}\| \|\mathbf{x} - \mathbf{x}_0\|.$$

Thus

$$0 \geq -\|\nabla f(\mathbf{x}_0) + \mathbf{g}\| \|\mathbf{x} - \mathbf{x}_0\| + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_0\|^2.$$

Rearranging (and assuming  $\mathbf{x} \neq \mathbf{x}_0$  so the norm is nonzero) gives

$$\frac{m}{2} \|\mathbf{x} - \mathbf{x}_0\| \leq \|\nabla f(\mathbf{x}_0) + \mathbf{g}\|.$$

Hence  $\|\mathbf{x} - \mathbf{x}_0\|$  is uniformly bounded for all  $\mathbf{x} \in S$ . Therefore  $S$  is bounded.



**Claim**

$F$  is bounded below on  $S$  and attains a minimum over  $S$  (and hence over  $\mathbb{R}^n$ ).

**Proof**

For any  $\mathbf{x} \in S$ , using the same inequality as above but dropping the quadratic term yields

$$F(\mathbf{x}) \geq F(\mathbf{x}_0) + (\nabla f(\mathbf{x}_0) + \mathbf{g})^\top (\mathbf{x} - \mathbf{x}_0) \geq F(\mathbf{x}_0) - \|\nabla f(\mathbf{x}_0) + \mathbf{g}\| \|\mathbf{x} - \mathbf{x}_0\|.$$

Using the bound on  $\|\mathbf{x} - \mathbf{x}_0\|$  from Claim 2, we obtain a uniform lower bound:

$$F(\mathbf{x}) \geq F_{\min} \quad \text{for all } \mathbf{x} \in S$$

for some finite  $F_{\min}$ .

Since  $S$  is closed and bounded, it is compact. The epigraph

$$\text{epi}(F) = \{(\mathbf{x}, y) : y \geq F(\mathbf{x})\}$$

is closed, so

$$\text{epi}(F) \cap (S \times [F_{\min}, \infty))$$

is a nonempty compact set. The projection map  $\phi(\mathbf{x}, y) = y$  attains its minimum on this set at some  $(\mathbf{x}^*, F(\mathbf{x}^*))$ . Thus  $F$  attains its minimum on  $S$  at  $\mathbf{x}^*$ .

For any  $\mathbf{x} \notin S$ , we have  $F(\mathbf{x}) > F(\mathbf{x}_0) \geq F(\mathbf{x}^*)$ , so  $\mathbf{x}^*$  is in fact a global minimizer of  $F$  over  $\mathbb{R}^n$ .

**16.2 Proximal Operator**

Let  $\psi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  be proper, closed, and convex. The **proximal operator** of  $\psi$  is defined by

$$\text{prox}_\psi(\mathbf{x}) := \arg \min_{\mathbf{z}} \left\{ \psi(\mathbf{z}) + \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|^2 \right\}.$$

Intuitively,  $\text{prox}_\psi(\mathbf{x})$  finds a point  $\mathbf{z}$  that balances two goals: staying close to the current point  $\mathbf{x}$ , and having a small value of the regularizer  $\psi$ .

**Theorem****Existence and uniqueness of the proximal point.**

For each  $\mathbf{x} \in \mathbb{R}^n$ , the function

$$\mathbf{z} \mapsto \psi(\mathbf{z}) + \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|^2$$

has a unique minimizer, i.e.  $\text{prox}_\psi(\mathbf{x})$  is well-defined and single-valued.

**Proof**

The function  $\mathbf{z} \mapsto \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|^2$  is strongly convex. Adding the convex function  $\psi$  preserves strong convexity, so  $\psi(\mathbf{z}) + \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|^2$  is strongly convex and proper, closed. Hence it has a unique minimizer.

We say that  $\psi$  is *proximable* if there exists an efficient algorithm to compute  $\text{prox}_\psi(\mathbf{x})$ .

### Example

#### Example 1: Scalar $\ell_1$ -prox.

Let  $n = 1$  and

$$\psi(x) = t|x|, \quad t > 0.$$

Define

$$q(z) := t|z| + \frac{1}{2}(z - x)^2.$$

Then  $\text{prox}_{t|\cdot|}(x)$  is the unique minimizer of  $q$ . We have

$$\partial q(z) = t \partial|z| + (z - x).$$

Recall

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

The optimality condition

$$0 \in \partial q(z) \iff 0 \in t \partial|z| + z - x$$

leads to three cases:

$$x - z \in \begin{cases} \{-t\}, & z < 0, \\ [-t, t], & z = 0, \\ \{t\}, & z > 0. \end{cases}$$

- **Case  $z < 0$ :**  $x - z = -t \Rightarrow z = x + t$ . Validity requires  $z < 0 \Rightarrow x < -t$ .
- **Case  $z = 0$ :**  $x - z \in [-t, t] \Rightarrow x \in [-t, t]$ .
- **Case  $z > 0$ :**  $x - z = t \Rightarrow z = x - t$ . Validity requires  $z > 0 \Rightarrow x > t$ .

Summarizing,

$$\text{prox}_{t|\cdot|}(x) = \begin{cases} x + t, & x < -t, \\ 0, & x \in [-t, t], \\ x - t, & x > t. \end{cases}$$

This is the usual scalar *soft-thresholding* operator.

### Example

#### Example 2: Vector $\ell_1$ -prox.

For  $\mathbf{x} \in \mathbb{R}^n$  and  $t > 0$ ,

$$\text{prox}_{t\|\cdot\|_1}(\mathbf{x}) = \arg \min_{\mathbf{z}} \{t\|\mathbf{z}\|_1 + \frac{1}{2}\|\mathbf{z} - \mathbf{x}\|^2\}.$$

Because the function is separable across coordinates, the solution is obtained by applying the

scalar soft thresholding coordinatewise:

$$\text{prox}_{t\|\cdot\|_1}(\mathbf{x}) = \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} \quad \text{with} \quad u_i = \begin{cases} x_i + t, & x_i < -t, \\ 0, & x_i \in [-t, t], \\ x_i - t, & x_i > t, \end{cases} \quad i = 1, \dots, n.$$

### Example

#### Example 3: Prox of the Euclidean norm.

Let

$$\psi(\mathbf{x}) = t\|\mathbf{x}\|, \quad t > 0.$$

We know

$$\partial\|\mathbf{z}\| = \begin{cases} \left\{ \frac{\mathbf{z}}{\|\mathbf{z}\|} \right\}, & \mathbf{z} \neq \mathbf{0}, \\ \overline{\mathbb{B}}(\mathbf{0}, 1), & \mathbf{z} = \mathbf{0}, \end{cases}$$

where  $\overline{\mathbb{B}}(\mathbf{0}, 1)$  is the closed unit ball.

Consider

$$q(\mathbf{z}) := t\|\mathbf{z}\| + \frac{1}{2}\|\mathbf{z} - \mathbf{x}\|^2.$$

Then

$$\mathbf{0} \in \partial q(\mathbf{z}) \iff \mathbf{0} \in t\partial\|\mathbf{z}\| + (\mathbf{z} - \mathbf{x}),$$

i.e.

$$\mathbf{x} - \mathbf{z} \in \begin{cases} \overline{\mathbb{B}}(\mathbf{0}, t), & \mathbf{z} = \mathbf{0} \quad (\text{a}), \\ t \frac{\mathbf{z}}{\|\mathbf{z}\|}, & \mathbf{z} \neq \mathbf{0} \quad (\text{b}). \end{cases}$$

**Case (a):**  $\mathbf{z} = \mathbf{0}$  is valid if and only if  $\mathbf{x} \in \overline{\mathbb{B}}(\mathbf{0}, t)$ , i.e.  $\|\mathbf{x}\| \leq t$ .

**Case (b):** Assume  $\mathbf{z} \neq \mathbf{0}$  and write  $\mathbf{z} = \theta\mathbf{x}$  for some scalar  $\theta$ . Then

$$\mathbf{x} - \theta\mathbf{x} = t \frac{\theta\mathbf{x}}{\|\theta\mathbf{x}\|} = t \frac{\theta}{|\theta|} \frac{\mathbf{x}}{\|\mathbf{x}\|}.$$

Assuming  $\theta > 0$  (we expect  $\mathbf{z}$  to be in the same direction as  $\mathbf{x}$ ), we get

$$1 - \theta = \frac{t}{\|\mathbf{x}\|} \implies \theta = 1 - \frac{t}{\|\mathbf{x}\|}.$$

Thus

$$\mathbf{z} = \left(1 - \frac{t}{\|\mathbf{x}\|}\right) \mathbf{x},$$

which is valid only when  $\mathbf{z} \neq \mathbf{0}$ , i.e.  $\|\mathbf{x}\| > t$ .

Combining both cases:

$$\text{prox}_{t\|\cdot\|}(\mathbf{x}) = \begin{cases} \mathbf{0}, & \|\mathbf{x}\| \leq t, \\ \left(1 - \frac{t}{\|\mathbf{x}\|}\right) \mathbf{x}, & \|\mathbf{x}\| > t. \end{cases}$$

This is often called the *vector soft-thresholding* or *shrinkage* operator.

**Fact**

For any closed, proper, convex function  $\psi$ , the proximal map  $\text{prox}_\psi : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is 1-Lipschitz, i.e.

$$\|\text{prox}_\psi(\mathbf{x}) - \text{prox}_\psi(\mathbf{y})\| \leq \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

In particular, the proximal map is continuous.

**16.3 Proximal Gradient**

We return to the composite setup

$$\min_{\mathbf{x}} F(\mathbf{x}), \quad F(\mathbf{x}) = f(\mathbf{x}) + \psi(\mathbf{x}),$$

with  $f$  convex and  $L$ -smooth and  $\psi$  proper, closed, convex.

**Proximal gradient.**

The *proximal gradient* of  $F$  (with parameter  $L$ ) is

$$\mathbf{G}_L(\mathbf{x}) := L \left( \mathbf{x} - \text{prox}_{\psi/L} \left( \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}) \right) \right).$$

Intuitively, the mapping

$$\mathbf{x} \mapsto \text{prox}_{\psi/L} \left( \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}) \right)$$

performs one gradient step on the smooth part  $f$  followed by one prox step on the (possibly nonsmooth) part  $\psi$ . The vector  $\mathbf{G}_L(\mathbf{x})$  is the scaled difference between the current point and this update, and plays the role of a generalized gradient.

**Example****Special case  $\psi \equiv 0$ .**

If  $\psi \equiv 0$ , then

$$\text{prox}_0(\mathbf{x}) = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|^2 = \mathbf{x}.$$

Hence

$$\mathbf{G}_L(\mathbf{x}) = L \left( \mathbf{x} - \left( \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}) \right) \right) = \nabla f(\mathbf{x}),$$

so the proximal gradient reduces to the usual gradient.

## LEC 17

## 17.1 Proximal Gradient

## 17.1.1 Set-up

Recall the composite problem

$$\min_{\mathbf{x}} F(\mathbf{x}), \quad F(\mathbf{x}) = f(\mathbf{x}) + \psi(\mathbf{x}),$$

where

- $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex and  $L$ -smooth,
- $\psi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  is convex, closed, proper, and proximable.

The *proximal gradient* is defined as

$$\mathbf{G}_L(\mathbf{x}) := L \left( \mathbf{x} - \text{prox}_{\psi/L} \left( \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}) \right) \right).$$

**Theorem**

For  $F$  and  $\mathbf{G}_L$  as above,

$$\mathbf{x} \in \arg \min F \iff \mathbf{G}_L(\mathbf{x}) = \mathbf{0}.$$

**Proof**

$$\begin{aligned} \mathbf{G}_L(\mathbf{x}) = \mathbf{0} &\iff \mathbf{x} = \text{prox}_{\psi/L} \left( \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}) \right) \\ &\iff \mathbf{x} = \arg \min_z \left\{ \frac{1}{L} \psi(z) + \frac{1}{2} \left\| z - \mathbf{x} + \frac{1}{L} \nabla f(\mathbf{x}) \right\|^2 \right\} \\ &\iff \mathbf{0} \in \left[ \frac{1}{L} \partial \psi(z) + \left\{ z - \mathbf{x} + \frac{1}{L} \nabla f(\mathbf{x}) \right\} \right]_{z=\mathbf{x}} \\ &\iff \mathbf{0} \in \frac{1}{L} \partial \psi(\mathbf{x}) + \frac{1}{L} \{ \nabla f(\mathbf{x}) \} \\ &\iff \mathbf{0} \in \frac{1}{L} (\partial \psi(\mathbf{x}) + \{ \nabla f(\mathbf{x}) \}) \\ &\iff \mathbf{0} \in \partial \psi(\mathbf{x}) + \{ \nabla f(\mathbf{x}) \} = \partial F(\mathbf{x}) \\ &\iff \mathbf{x} \in \arg \min F. \end{aligned}$$

## 17.1.2 Proximal Gradient Descent

The **Proximal Gradient Descent (PGD)** algorithm is

$$\begin{aligned} \mathbf{x}^0 &= \text{arbitrary}, \\ \mathbf{x}^{k+1} &= \mathbf{x}^k - \frac{1}{L} \mathbf{G}_L(\mathbf{x}^k). \end{aligned}$$

Under our assumptions, this iteration converges to a minimizer of  $F$  (assuming a minimizer exists). The proof is omitted in class.

### Accelerated Proximal Gradient Descent (*APGD*).

- Nesterov (2004),
- Beck & Teboulle (2013): **FISTA**.

*APGD* is obtained from *AGD* by replacing the gradient with the proximal gradient operator.

- *AGD* update:

$$\mathbf{x}^{k+1} = \mathbf{y}^k - \frac{1}{L} \nabla f(\mathbf{y}^k).$$

- *APGD* update:

$$\mathbf{x}^{k+1} = \mathbf{y}^k - \frac{1}{L} \mathbf{G}_L(\mathbf{y}^k).$$

More explicitly,

$$\begin{aligned} \mathbf{x}^{k+1} &= \mathbf{y}^k - \frac{1}{L} \mathbf{G}_L(\mathbf{y}^k) \\ &= \mathbf{y}^k - \frac{1}{L} \left[ L \left( \mathbf{y}^k - \text{prox}_{\psi/L} \left( \mathbf{y}^k - \frac{1}{L} \nabla f(\mathbf{y}^k) \right) \right) \right] \\ &= \mathbf{y}^k - \mathbf{y}^k + \text{prox}_{\psi/L} \left( \mathbf{y}^k - \frac{1}{L} \nabla f(\mathbf{y}^k) \right) \\ &= \text{prox}_{\psi/L} \left( \mathbf{y}^k - \frac{1}{L} \nabla f(\mathbf{y}^k) \right). \end{aligned}$$

This view is called **operator splitting**, or a **forward–backward step**:

- forward step:  $\mathbf{y}^k - \frac{1}{L} \nabla f(\mathbf{y}^k)$  (gradient step on  $f$ ),
- backward step:  $\text{prox}_{\psi/L}(\cdot)$  (prox step on  $\psi$ ).

The terminology comes from an analogy with forward/backward Euler methods in ODEs.

#### Theorem

##### Convergence rates for *PGD* & *APGD*.

Let  $F = f + \psi$  with  $f$  convex and  $L$ -smooth, and  $\psi$  convex, closed, proper.

- ***PGD* (sublinear rate)**. For a minimizer  $\mathbf{x}^*$  of  $F$ ,

$$F(\mathbf{x}^k) - F(\mathbf{x}^*) \leq \frac{L \|\mathbf{x}^0 - \mathbf{x}^*\|^2}{2k}.$$

- If  $f$  is additionally  $m$ -strongly convex with  $m > 0$  (so  $\mathbf{x}^*$  is unique), then *PGD* has a **linear** rate:

$$\|\mathbf{x}^k - \mathbf{x}^*\| \leq \left( \frac{L - m}{L + m} \right)^k \|\mathbf{x}^0 - \mathbf{x}^*\|.$$

- ***APGD* (accelerated rate)**. Without strong convexity, *APGD* attains

$$F(\mathbf{x}^k) - F(\mathbf{x}^*) = \mathcal{O}\left(\frac{1}{k^2}\right).$$

- If  $f$  is also  $m$ -strongly convex,  $APGD$  enjoys a linear rate with factor

$$\left( \frac{\sqrt{L} - \sqrt{m}}{\sqrt{L} + \sqrt{m}} \right)^k.$$

For details, see Beck, *First-Order Methods in Optimization* (2017).

A common termination criterion for  $PGD$  /  $APGD$  is

$$\|G_L(\mathbf{x}^k)\| \leq \text{tol}.$$

### Alternative view of one $PGD$ step.

We can write

$$\begin{aligned} \mathbf{x}^{k+1} &= \text{prox}_{\psi/L} \left( \mathbf{x}^k - \frac{1}{L} \nabla f(\mathbf{x}^k) \right) \\ &= \arg \min_z \left\{ \frac{1}{L} \psi(z) + \frac{1}{2} \left\| \mathbf{x}^k - \frac{1}{L} \nabla f(\mathbf{x}^k) - z \right\|^2 \right\}. \end{aligned}$$

Expanding the square,

$$\begin{aligned} \mathbf{x}^{k+1} &= \arg \min_z \left\{ \frac{1}{L} \psi(z) + \frac{1}{2} \|\mathbf{x}^k - z\|^2 - \frac{1}{L} \nabla f(\mathbf{x}^k)^\top (\mathbf{x}^k - z) + \frac{1}{2L^2} \|\nabla f(\mathbf{x}^k)\|^2 \right\} \\ &= \arg \min_z \left\{ \psi(z) + \frac{L}{2} \|\mathbf{x}^k - z\|^2 - \nabla f(\mathbf{x}^k)^\top (\mathbf{x}^k - z) \right\} \\ &= \arg \min_z \left\{ \underbrace{\psi(z) + f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^\top (z - \mathbf{x}^k) + \frac{L}{2} \|z - \mathbf{x}^k\|^2}_{\tilde{f}(z)} \right\}. \end{aligned}$$

Here  $\tilde{f}(z)$  is a quadratic upper model of  $f$  at  $\mathbf{x}^k$ : we replace  $f(z)$  by its first-order Taylor expansion at  $\mathbf{x}^k$  plus a quadratic term  $\frac{L}{2} \|z - \mathbf{x}^k\|^2$ .

### Theorem

For all  $z \in \mathbb{R}^n$ ,

$$\tilde{f}(z) \geq f(z).$$

This inequality is simply the  $L$ -smoothness (descent lemma) of  $f$ .

Why not use the *true* Hessian  $\nabla^2 f(\mathbf{x}^k)$  as the quadratic term? In general, the resulting subproblem

$$\min_z \psi(z) + f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^\top (z - \mathbf{x}^k) + \frac{1}{2} (z - \mathbf{x}^k)^\top \nabla^2 f(\mathbf{x}^k) (z - \mathbf{x}^k)$$

does not admit an efficient closed-form solution, whereas the  $\frac{L}{2} \|z - \mathbf{x}^k\|^2$  model does for many popular regularizers  $\psi$ .

## 17.2 Projected Gradient

### 17.2.1 Definition

Projected Gradient is the special case of *PGD* where we enforce a simple constraint  $\mathbf{x} \in \Omega$  via an indicator function.

Let  $\Omega \subseteq \mathbb{R}^n$  be closed, convex, nonempty, and define

$$\psi(\mathbf{x}) = \mathbb{I}_{\Omega}(\mathbf{x}) = \begin{cases} 0, & \mathbf{x} \in \Omega, \\ \infty, & \mathbf{x} \notin \Omega. \end{cases}$$

Then

$$\text{prox}_{t\mathbb{I}_{\Omega}}(\mathbf{x}) = \arg \min_{\mathbf{z}} \left\{ t\mathbb{I}_{\Omega}(\mathbf{z}) + \frac{1}{2}\|\mathbf{z} - \mathbf{x}\|^2 \right\} = \arg \min_{\mathbf{z} \in \Omega} \frac{1}{2}\|\mathbf{z} - \mathbf{x}\|^2.$$

The *projection* onto  $\Omega$  is

$$\text{proj}_{\Omega}(\mathbf{x}) := \arg \min_{\mathbf{z} \in \Omega} \|\mathbf{z} - \mathbf{x}\|.$$

In this setting, *PGD* is simply gradient descent followed by projection, and *APGD* also applies provided we can compute  $\text{proj}_{\Omega}$  efficiently.

### 17.2.2 Examples

#### Example

##### Example 1: Projection onto a ball.

Let  $\mathbb{B}(\mathbf{0}, r) = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| \leq r\}$ . Then

$$\text{proj}_{\mathbb{B}(\mathbf{0}, r)}(\mathbf{x}) = \begin{cases} \mathbf{x}, & \|\mathbf{x}\| \leq r, \\ r \frac{\mathbf{x}}{\|\mathbf{x}\|}, & \|\mathbf{x}\| > r. \end{cases}$$

#### Example

##### Example 2: Projection onto an affine subspace.

Let

$$\Omega = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{b}\}, \quad \mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{b} \in \mathbb{R}^m.$$

##### Precomputation.

- Find some  $\mathbf{x}_0$  such that  $\mathbf{A}\mathbf{x}_0 = \mathbf{b}$ .
- Compute a matrix  $\mathbf{U} \in \mathbb{R}^{n \times p}$  whose columns form an orthonormal basis of  $\text{Null}(\mathbf{A})$ , so  $\mathbf{U}^{\top}\mathbf{U} = \mathbf{I}_p$ .

If  $\text{rank}(\mathbf{A}) = m$  (full row rank), we may use a **QR factorization** of  $\mathbf{A}^{\top}$ ; then  $p = n - m$  by the rank-nullity theorem. If  $\text{rank}(\mathbf{A}) < m$ , one can use an **SVD** of  $\mathbf{A}$ ; then  $p > n - m$ . We assume  $\mathbf{b} \in \text{Range}(\mathbf{A})$  so that  $\Omega \neq \emptyset$ .



Note that

$$\mathbf{Ax} = \mathbf{b} \iff \mathbf{Ax} = \mathbf{Ax}_0 \iff \mathbf{A}(\mathbf{x} - \mathbf{x}_0) = \mathbf{0} \iff \mathbf{x} - \mathbf{x}_0 \in \text{Null}(\mathbf{A}) \iff \mathbf{x} - \mathbf{x}_0 = \mathbf{U}\mathbf{w}$$

for some  $\mathbf{w} \in \mathbb{R}^p$ .

### Theorem

With  $\Omega, \mathbf{A}, \mathbf{b}, \mathbf{U}, \mathbf{x}_0$  as above,

$$\text{proj}_{\Omega}(\mathbf{y}) = \mathbf{x}_0 + \mathbf{U}\mathbf{U}^{\top}(\mathbf{y} - \mathbf{x}_0).$$

### Proof

We want

$$\min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2 : \mathbf{Ax} = \mathbf{b} \right\}.$$

Using the parameterization  $\mathbf{x} = \mathbf{x}_0 + \mathbf{U}\mathbf{w}$ , this is equivalent to

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{x}_0 + \mathbf{U}\mathbf{w} - \mathbf{y}\|^2.$$

Expanding,

$$\begin{aligned} \frac{1}{2} \|\mathbf{x}_0 + \mathbf{U}\mathbf{w} - \mathbf{y}\|^2 &= \frac{1}{2} \mathbf{w}^{\top} \mathbf{U}^{\top} \mathbf{U} \mathbf{w} + (\mathbf{x}_0 - \mathbf{y})^{\top} \mathbf{U} \mathbf{w} + C_1 \\ &= \frac{1}{2} \mathbf{w}^{\top} \mathbf{w} + (\mathbf{x}_0 - \mathbf{y})^{\top} \mathbf{U} \mathbf{w} + C_1 \\ &= \frac{1}{2} \|\mathbf{w} + \mathbf{U}^{\top}(\mathbf{x}_0 - \mathbf{y})\|^2 + C_2, \end{aligned}$$

where  $C_1, C_2$  are constants independent of  $\mathbf{w}$ .

This is minimized by

$$\mathbf{w}^* = \mathbf{U}^{\top}(\mathbf{y} - \mathbf{x}_0).$$

Hence

$$\mathbf{x}^* = \mathbf{x}_0 + \mathbf{U}\mathbf{w}^* = \mathbf{x}_0 + \mathbf{U}\mathbf{U}^{\top}(\mathbf{y} - \mathbf{x}_0),$$

which is exactly  $\text{proj}_{\Omega}(\mathbf{y})$ .

## LEC 18

## 18.1 Nonlinear Programming

Consider general nonlinear programming (NLP):

$$\min_{\mathbf{x}} f_0(\mathbf{x}) \quad \text{S.T.} \quad f_1(\mathbf{x}) \leq 0, \dots, f_m(\mathbf{x}) \leq 0, \quad h_1(\mathbf{x}) = 0, \dots, h_p(\mathbf{x}) = 0,$$

where  $f_0, \dots, f_m, h_1, \dots, h_p : \mathbb{R}^n \rightarrow \mathbb{R}$ .

Let the **feasible region** be

$$\Omega := \{\mathbf{x} : f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m; \quad h_i(\mathbf{x}) = 0, \quad i = 1, \dots, p\}.$$

Define the **Lagrangian**

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{i=1}^p \mu_i h_i(\mathbf{x}),$$

where  $\boldsymbol{\lambda} \in \mathbb{R}^m$ ,  $\boldsymbol{\mu} \in \mathbb{R}^p$  are the **multipliers**.

**Theorem**

$$\inf_{\mathbf{x}} \{f_0(\mathbf{x}) : \mathbf{x} \in \Omega\} = \inf_{\mathbf{x}} \sup_{\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\mu}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}).$$

**Note.**  $\inf$  is like  $\min$ , except the minimum need not be attained (e.g.  $\inf\{e^x : x \in \mathbb{R}\} = 0$ , and  $\inf \emptyset = \infty$ ). Similarly,  $\sup$  is analogous to  $\max$ .

**Proof**

Fix  $\mathbf{x} \in \mathbb{R}^n$ .

**Case 1:  $\mathbf{x} \notin \Omega$  (infeasible).**

Then either  $f_i(\mathbf{x}) > 0$  for some  $i$ , or  $h_i(\mathbf{x}) \neq 0$  for some  $i$  (or both).

If  $f_i(\mathbf{x}) > 0$  for some  $i$ , set all multipliers to 0 except  $\lambda_i \rightarrow \infty$ . Then

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \rightarrow \infty,$$

so

$$\sup_{\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\mu}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \infty,$$

which can never be the value of the outer infimum.

Similarly, if  $h_i(\mathbf{x}) \neq 0$  for some  $i$ , fix all multipliers at 0 except  $\mu_i = t \cdot \text{sgn}(h_i(\mathbf{x}))$ ,  $t \rightarrow \infty$ .

Again  $\sup_{\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\mu}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \infty$ .

So for  $\mathbf{x} \notin \Omega$ , the inner sup is  $\infty$  and cannot contribute to the outer inf.

**Case 2:  $\mathbf{x} \in \Omega$  (feasible).**

Then  $f_i(\mathbf{x}) \leq 0$  for all  $i = 1, \dots, m$  and  $h_i(\mathbf{x}) = 0$  for all  $i = 1, \dots, p$ . For any  $\boldsymbol{\lambda} \geq \mathbf{0}$ , any  $\boldsymbol{\mu}$ ,

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{i=1}^p \mu_i h_i(\mathbf{x}) \leq f_0(\mathbf{x}),$$

since  $\sum \lambda_i f_i(\mathbf{x}) \leq 0$  and  $\sum \mu_i h_i(\mathbf{x}) = 0$ . In fact,

$$\sup_{\lambda \geq 0, \mu} \mathcal{L}(\mathbf{x}, \lambda, \mu) = f_0(\mathbf{x}),$$

attained by taking  $\lambda = \mathbf{0}$  (and any  $\mu$ ).

Therefore, when we take the outer infimum over  $\mathbf{x}$ , we see:

$$\sup_{\lambda, \mu} \mathcal{L}(\mathbf{x}, \lambda, \mu) = \begin{cases} \infty, & \mathbf{x} \notin \Omega, \\ f_0(\mathbf{x}), & \mathbf{x} \in \Omega. \end{cases}$$

Thus

$$\inf_{\mathbf{x}} \sup_{\lambda \geq 0, \mu} \mathcal{L}(\mathbf{x}, \lambda, \mu) = \inf_{\mathbf{x} \in \Omega} f_0(\mathbf{x}),$$

as claimed.

### Theorem

(Weak Duality.)

$$\inf_{\mathbf{x}} \sup_{\lambda \geq 0, \mu} \mathcal{L}(\mathbf{x}, \lambda, \mu) \geq \sup_{\lambda \geq 0, \mu} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \mu).$$

### Proof

Fix  $\mathbf{x}_1 \in \mathbb{R}^n$  and  $\lambda_1 \geq \mathbf{0}, \mu_1$ . Then

$$\mathcal{L}(\mathbf{x}_1, \lambda_1, \mu_1) \geq \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda_1, \mu_1).$$

Since  $\lambda_1, \mu_1$  are arbitrary, we can take sup over them:

$$\sup_{\lambda_1 \geq 0, \mu_1} \mathcal{L}(\mathbf{x}_1, \lambda_1, \mu_1) \geq \sup_{\lambda_1 \geq 0, \mu_1} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda_1, \mu_1).$$

The left-hand side holds for all  $\mathbf{x}_1$ , so taking inf over  $\mathbf{x}_1$  preserves the inequality:

$$\inf_{\mathbf{x}_1} \sup_{\lambda_1 \geq 0, \mu_1} \mathcal{L}(\mathbf{x}_1, \lambda_1, \mu_1) \geq \sup_{\lambda_1 \geq 0, \mu_1} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda_1, \mu_1).$$

**Game interpretation.** Think of a two-player game:

- the inf-player chooses  $\mathbf{x}$  (tries to minimize the payoff),
- the sup-player chooses  $(\lambda, \mu)$  (tries to maximize the payoff).

Weak duality says: it is better for the inf-player (i.e. lower cost) if the sup-player must announce their strategy first.

#### 18.1.1 Dual Objective and Dual Problem

Define the **dual objective**

$$g(\lambda, \mu) := \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \mu).$$

The **dual optimization problem** is

$$\max_{\lambda, \mu} g(\lambda, \mu) \quad \text{S.T.} \quad \lambda \geq \mathbf{0}.$$

Weak duality can be rewritten as

$$\sup_{\lambda \geq \mathbf{0}, \mu} g(\lambda, \mu) \leq \inf_{\mathbf{x}} \{f_0(\mathbf{x}) : \mathbf{x} \in \Omega\}.$$

Equivalently:

- For any primal feasible  $\mathbf{x} \in \Omega$ ,

$$\sup_{\lambda \geq \mathbf{0}, \mu} g(\lambda, \mu) \leq f_0(\mathbf{x}).$$

- For any dual feasible  $(\lambda, \mu)$  with  $\lambda \geq \mathbf{0}$ ,

$$\inf_{\mathbf{x}} \{f_0(\mathbf{x}) : \mathbf{x} \in \Omega\} \geq g(\lambda, \mu).$$

Weak duality requires no assumptions, so it is very general but often too weak by itself.

#### Theorem

**(Strong Duality.)**

$$\sup_{\lambda \geq \mathbf{0}, \mu} g(\lambda, \mu) = \inf_{\mathbf{x}} \{f_0(\mathbf{x}) : \mathbf{x} \in \Omega\}.$$

Equivalently,

$$\inf_{\mathbf{x}} \sup_{\lambda \geq \mathbf{0}, \mu} \mathcal{L}(\mathbf{x}, \lambda, \mu) = \sup_{\lambda \geq \mathbf{0}, \mu} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \mu).$$

**Note.** The equality above does *not* require that the infimum/supremum are attained. Some authors define strong duality to include the existence of optimizers:

$$\exists \mathbf{x}^* \in \Omega, \lambda^* \geq \mathbf{0}, \mu^* \quad \text{s.t.} \quad g(\lambda^*, \mu^*) = f_0(\mathbf{x}^*).$$

#### Theorem

**(KKT-type characterization.)**

Strong duality is attained at some  $\mathbf{x}^* \in \Omega$ ,  $\lambda^* \geq \mathbf{0}$ ,  $\mu^*$  if and only if all four conditions hold:

(i) **Primal feasibility:**  $\mathbf{x}^* \in \Omega$ ;

(ii) **Dual feasibility:**  $\lambda^* \geq \mathbf{0}$ ;

(iii) **Complementarity:**

$$\lambda_i^* f_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m;$$

(iv) **Stationarity:**

$$\mathbf{x}^* \in \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda^*, \mu^*).$$

These are exactly the (convex) KKT conditions.

## Proof

( $\Leftarrow$ ) Assume (1)–(4) hold.

Using (1) and (3),

$$\sum_{i=1}^m \lambda_i^* f_i(\mathbf{x}^*) = 0, \quad \sum_{i=1}^p \mu_i^* h_i(\mathbf{x}^*) = 0.$$

Thus

$$\mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = f_0(\mathbf{x}^*) + 0 + 0 = f_0(\mathbf{x}^*).$$

By (4),

$$\mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = g(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*).$$

Hence  $g(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = f_0(\mathbf{x}^*)$ , so strong duality holds and is attained at  $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ .

( $\Rightarrow$ ) Suppose strong duality holds and is attained at  $\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*$ , i.e.

$$g(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = f_0(\mathbf{x}^*), \quad \mathbf{x}^* \in \Omega, \quad \boldsymbol{\lambda}^* \geq \mathbf{0}.$$

Then (1) and (2) hold by assumption.

From the strong duality equality and weak duality chain,

$$\begin{aligned} f_0(\mathbf{x}^*) &= \inf_{\mathbf{x}} \sup_{\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\mu}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \\ &\geq \sup_{\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\mu}} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \\ &= g(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \\ &\leq \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \\ &\leq \sup_{\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\mu}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}, \boldsymbol{\mu}) \\ &\leq f_0(\mathbf{x}^*), \end{aligned}$$

so all inequalities are in fact equalities. In particular,

$$f_0(\mathbf{x}^*) = \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = f_0(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* f_i(\mathbf{x}^*) + \sum_{i=1}^p \mu_i^* h_i(\mathbf{x}^*).$$

Hence

$$\sum_{i=1}^m \lambda_i^* f_i(\mathbf{x}^*) + \sum_{i=1}^p \mu_i^* h_i(\mathbf{x}^*) = 0.$$

But  $f_i(\mathbf{x}^*) \leq 0$  and  $\lambda_i^* \geq 0$ ; therefore each term  $\lambda_i^* f_i(\mathbf{x}^*) \leq 0$ , and the sum being zero forces

$$\lambda_i^* f_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m$$

(complementarity, (3)). Also,

$$f_0(\mathbf{x}^*) = \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*),$$

so  $\mathbf{x}^*$  is a minimizer of  $\mathcal{L}(\cdot, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ , i.e. (4) holds.

**Convex specialization.**

Suppose  $f_0, \dots, f_m$  are convex, and the equality constraints are affine:

$$\begin{pmatrix} h_1(\mathbf{x}) \\ \vdots \\ h_p(\mathbf{x}) \end{pmatrix} = \mathbf{A}\mathbf{x} - \mathbf{b}.$$

Then for fixed  $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$  with  $\boldsymbol{\lambda}^* \geq \mathbf{0}$ , the Lagrangian  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$  is convex in  $\mathbf{x}$ .

Condition (4) is equivalent to

$$\mathbf{0} \in \partial_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*).$$

If moreover  $f_0, \dots, f_m$  are differentiable, then

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = \nabla f_0(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla f_i(\mathbf{x}^*) + \mathbf{A}^\top \boldsymbol{\mu}^*,$$

since

$$\nabla_{\mathbf{x}} (\boldsymbol{\mu}^\top (\mathbf{A}\mathbf{x} - \mathbf{b})) = \mathbf{A}^\top \boldsymbol{\mu}.$$

In earlier KKT formulations, the sum over  $i$  was only over *active* constraints. Here we have a multiplier  $\lambda_i^*$  for every inequality constraint, but complementarity (3) gives  $\lambda_i^* = 0$  whenever  $f_i(\mathbf{x}^*) < 0$  (inactive).

## 18.2 Kernel Machines

Consider **ridge regression**:

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 + \frac{1}{2} \gamma \|\mathbf{x}\|^2,$$

with  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\gamma > 0$ .

The closed-form solution is

$$\mathbf{x}^* = (\mathbf{A}^\top \mathbf{A} + \gamma \mathbf{I})^{-1} \mathbf{A}^\top \mathbf{y}.$$

Often one wants to enrich the feature set: a linear model may not fit well (e.g. a parabola is needed). Starting from  $p$  original features, including all quadratic monomials would yield  $\frac{p(p+1)}{2}$  additional feature columns.

In general, we might want  $N$  new features with  $N \gg m$ . This leads naturally to kernel methods, where the feature space can be extremely high-dimensional (or infinite), but the optimization is done via inner products (kernels) rather than explicit feature vectors.

## LEC 19

## 19.1 Kernel Ridge Regression

## 19.1.1 Ridge Regression

Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , consider

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 + \frac{1}{2} \gamma \|\mathbf{x}\|^2.$$

## 19.1.2 Kernel Ridge Regression

Now “lift” to a higher-dimensional feature matrix  $\tilde{\mathbf{A}} \in \mathbb{R}^{m \times N}$  with  $N \gg \max(m, n)$  and solve

$$\tilde{\mathbf{x}} := (\tilde{\mathbf{A}}^\top \tilde{\mathbf{A}} + \gamma \mathbf{I}_N)^{-1} \tilde{\mathbf{A}}^\top \mathbf{y}.$$

## 19.2 Kernel Trick

## 19.2.1 Part 1

## Theorem

$$\underbrace{(\tilde{\mathbf{A}}^\top \tilde{\mathbf{A}} + \gamma \mathbf{I}_N)^{-1}}_{N \times N} \tilde{\mathbf{A}}^\top = \tilde{\mathbf{A}}^\top \underbrace{(\tilde{\mathbf{A}} \tilde{\mathbf{A}}^\top + \gamma \mathbf{I}_m)^{-1}}_{m \times m}.$$

## Proof

$$\begin{aligned} (\tilde{\mathbf{A}}^\top \tilde{\mathbf{A}} + \gamma \mathbf{I}_N) \tilde{\mathbf{A}}^\top (\tilde{\mathbf{A}} \tilde{\mathbf{A}}^\top + \gamma \mathbf{I}_m)^{-1} &= (\tilde{\mathbf{A}}^\top \tilde{\mathbf{A}} \tilde{\mathbf{A}}^\top + \gamma \tilde{\mathbf{A}}^\top) (\tilde{\mathbf{A}} \tilde{\mathbf{A}}^\top + \gamma \mathbf{I}_m)^{-1} \\ &= \tilde{\mathbf{A}}^\top (\tilde{\mathbf{A}} \tilde{\mathbf{A}}^\top + \gamma \mathbf{I}_m) (\tilde{\mathbf{A}} \tilde{\mathbf{A}}^\top + \gamma \mathbf{I}_m)^{-1} \\ &= \tilde{\mathbf{A}}^\top. \end{aligned}$$

Thus the claimed identity holds.

## 19.2.2 Part 2

We do not need to store  $\tilde{\mathbf{x}} \in \mathbb{R}^N$  explicitly. By the theorem,

$$\tilde{\mathbf{x}} = \tilde{\mathbf{A}}^\top (\tilde{\mathbf{A}} \tilde{\mathbf{A}}^\top + \gamma \mathbf{I}_m)^{-1} \mathbf{y} = \tilde{\mathbf{A}}^\top \mathbf{w},$$

where

$$\mathbf{w} := (\tilde{\mathbf{A}} \tilde{\mathbf{A}}^\top + \gamma \mathbf{I}_m)^{-1} \mathbf{y} \in \mathbb{R}^m.$$

So we store only  $\mathbf{w}$ . For a new instance  $\tilde{\mathbf{a}} \in \mathbb{R}^N$ ,

$$\tilde{\mathbf{a}}^\top \tilde{\mathbf{x}} = \tilde{\mathbf{a}}^\top \tilde{\mathbf{A}}^\top \mathbf{w},$$

so classification can be done using matrix–vector products involving  $\tilde{\mathbf{A}}$  and the new  $\tilde{\mathbf{a}}$ , and  $\mathbf{w}$ .

### 19.2.3 Part 3

A **kernel** is the matrix

$$\mathbf{K} := \tilde{\mathbf{A}}\tilde{\mathbf{A}}^\top \in \mathbb{S}^m,$$

often called the Gram matrix or inner-product matrix.

For some feature families it is possible to compute an approximation  $\hat{\mathbf{K}}$  much faster than explicitly forming  $\tilde{\mathbf{A}}$  and then  $\tilde{\mathbf{A}}\tilde{\mathbf{A}}^\top$ .

Example: all monomials up to degree  $d$  in the original features.

#### Example

##### (Polynomial features.)

Let  $n = 2$ , degree  $d = 2$ . Then  $N = 6$  features:

$$1, x, y, x^2, y^2, xy.$$

For  $\mathbf{K}(i, j)$ , suppose row  $i$  of  $\mathbf{A}$  is  $(v_1, v_2)$  and row  $j$  is  $(w_1, w_2)$ . Then the corresponding rows of  $\tilde{\mathbf{A}}$  are

$$\begin{aligned}\text{row } i &\mapsto (1, v_1, v_2, v_1^2, v_1v_2, v_2^2), \\ \text{row } j &\mapsto (1, w_1, w_2, w_1^2, w_1w_2, w_2^2).\end{aligned}$$

Then  $\mathbf{K}(i, j)$  is the inner product of these two 6-vectors.

For a more convenient kernel, consider

$$(1 + v_1w_1 + v_2w_2)^2.$$

Expanding gives

$$1 + 2v_1w_1 + 2v_2w_2 + v_1^2w_1^2 + v_2^2w_2^2 + 2v_1w_1v_2w_2,$$

which is the inner product of

$$(1, \sqrt{2}v_1, \sqrt{2}v_2, v_1^2, \sqrt{2}v_1v_2, v_2^2)$$

and

$$(1, \sqrt{2}w_1, \sqrt{2}w_2, w_1^2, \sqrt{2}w_1w_2, w_2^2).$$

Call this inner product  $\hat{\mathbf{K}}(i, j)$ . Then  $\hat{\mathbf{K}}$  behaves like  $\mathbf{K}$  (up to scaling factors of  $\sqrt{2}$ , which only rescale  $\tilde{\mathbf{x}}$ ) but can be computed without explicitly forming the rows of  $\tilde{\mathbf{A}}$ .

*Conclusion.* The kernel trick lets us effectively add many columns (features) without working explicitly with an  $m \times N$  matrix.

Part 3 also applies to forming  $\tilde{\mathbf{A}}\tilde{\mathbf{a}}$  for a new data point  $\tilde{\mathbf{a}}$ : kernels can be defined for many feature families beyond polynomials.



### 19.3 Kernel Machines

Let the  $i$ -th row of  $\mathbf{A}$  be  $\boldsymbol{\xi}_i^\top \in \mathbb{R}^n$ . After lifting by a feature map

$$\Phi : \mathbb{R}^n \rightarrow \mathbb{R}^N, \quad \Phi(\boldsymbol{\xi}_i) = \begin{bmatrix} \Phi_1(\boldsymbol{\xi}_i) \\ \vdots \\ \Phi_N(\boldsymbol{\xi}_i) \end{bmatrix},$$

we have

$$\mathbf{A} = \begin{bmatrix} \boldsymbol{\xi}_1^\top \\ \vdots \\ \boldsymbol{\xi}_m^\top \end{bmatrix} \xrightarrow{\text{Map}} \tilde{\mathbf{A}} = \begin{bmatrix} \Phi^\top(\boldsymbol{\xi}_1) \\ \vdots \\ \Phi^\top(\boldsymbol{\xi}_m) \end{bmatrix} = \begin{bmatrix} \Phi_1(\boldsymbol{\xi}_1) & \cdots & \Phi_N(\boldsymbol{\xi}_1) \\ \vdots & \ddots & \vdots \\ \Phi_1(\boldsymbol{\xi}_m) & \cdots & \Phi_N(\boldsymbol{\xi}_m) \end{bmatrix}.$$

Then the kernel (Gram matrix)  $\mathbf{K}$  is

$$\mathbf{K} = \begin{bmatrix} \Phi^\top(\boldsymbol{\xi}_1)\Phi(\boldsymbol{\xi}_1) & \cdots & \Phi^\top(\boldsymbol{\xi}_1)\Phi(\boldsymbol{\xi}_m) \\ \vdots & \ddots & \vdots \\ \Phi^\top(\boldsymbol{\xi}_m)\Phi(\boldsymbol{\xi}_1) & \cdots & \Phi^\top(\boldsymbol{\xi}_m)\Phi(\boldsymbol{\xi}_m) \end{bmatrix} \in \mathbb{S}^m.$$

$\mathbf{K}$  is also called a *Gram matrix* or *inner-product matrix*.

#### Example

(Nonlinear separator via kernel SVM.)

Recall soft-margin SVM (SVM-7):

$$\min_{\mathbf{x}, \xi, \mathbf{s}} \frac{1}{2} \|\mathbf{x}\|^2 + \gamma \sum_{i=1}^m s_i \quad \text{S.T.} \quad \begin{cases} \mathbf{a}_i^\top \mathbf{x} + \xi \geq 1 - s_i, & y_i = 1, \\ \mathbf{a}_i^\top \mathbf{x} + \xi \leq -1 + s_i, & y_i = -1, \\ \mathbf{s} \geq \mathbf{0}, \end{cases}$$

with  $\mathbf{a}_i \in \mathbb{R}^n$ .

After mapping  $\mathbf{a}_i \mapsto \Phi(\mathbf{a}_i) \in \mathbb{R}^N$ , the *kernel SVM* seeks  $\tilde{\mathbf{x}} \in \mathbb{R}^N$ ,  $\xi \in \mathbb{R}$ ,  $\mathbf{s} \in \mathbb{R}^m$  solving

$$\min_{\tilde{\mathbf{x}}, \xi, \mathbf{s}} \frac{1}{2} \|\tilde{\mathbf{x}}\|^2 + \gamma \sum_{i=1}^m s_i \quad \text{S.T.} \quad \begin{cases} \Phi^\top(\mathbf{a}_i)\tilde{\mathbf{x}} + \xi \geq 1 - s_i, & y_i = 1, \\ \Phi^\top(\mathbf{a}_i)\tilde{\mathbf{x}} + \xi \leq -1 + s_i, & y_i = -1, \\ \mathbf{s} \geq \mathbf{0}. \end{cases}$$

#### Lagrangian formulation.

Introduce multipliers  $\boldsymbol{\lambda} \geq \mathbf{0}$  for  $y_i = 1$  constraints,  $\boldsymbol{\lambda}' \geq \mathbf{0}$  for  $y_i = -1$  constraints, and  $\boldsymbol{\pi} \geq \mathbf{0}$  for  $s_i \geq 0$ . The Lagrangian is

$$\min_{\tilde{\mathbf{x}}, \xi, \mathbf{s}} \max_{\substack{\boldsymbol{\lambda} \geq \mathbf{0} \\ \boldsymbol{\lambda}' \geq \mathbf{0} \\ \boldsymbol{\pi} \geq \mathbf{0}}} \mathcal{L}(\tilde{\mathbf{x}}, \xi, \mathbf{s}; \boldsymbol{\lambda}, \boldsymbol{\lambda}', \boldsymbol{\pi}),$$

where

$$\begin{aligned}\mathcal{L}(\tilde{\mathbf{x}}, \xi, \mathbf{s}; \boldsymbol{\lambda}, \boldsymbol{\lambda}', \boldsymbol{\pi}) = & \underbrace{\frac{1}{2} \|\tilde{\mathbf{x}}\|^2 + \gamma \sum_{i=1}^m s_i}_{f_0(\tilde{\mathbf{x}}, \mathbf{s})} \\ & + \sum_{y_i=1} \lambda_i (1 - s_i - \boldsymbol{\Phi}^\top(\mathbf{a}_i) \tilde{\mathbf{x}} - \xi) \\ & + \sum_{y_i=-1} \lambda'_i (1 - s_i + \boldsymbol{\Phi}^\top(\mathbf{a}_i) \tilde{\mathbf{x}} + \xi) + \sum_{i=1}^m (-s_i) \pi_i.\end{aligned}$$

### Quadratic programming (QP).

A QP is an optimization problem of the form

$$\min_{\mathbf{x}} q(\mathbf{x}) \quad \text{S.T.} \quad A_1 \mathbf{x} = \mathbf{b}_1, \quad A_2 \mathbf{x} \leq \mathbf{b}_2,$$

where  $q$  is quadratic with positive semidefinite Hessian. Soft-margin SVM (SVM-7) and  $\ell_1$ -LS are QPs.

#### Fact

If a QP has at least one feasible point and is bounded below, then it has an optimizer and strong duality holds.

For SVM-7 the objective is bounded below by 0, and there is an obvious feasible point (e.g.  $\tilde{\mathbf{x}} = \mathbf{0}$ ,  $\xi = 0$ ,  $\mathbf{s} = \mathbf{e}$ ), so strong duality applies.

Thus the dual problem is

$$\max_{\substack{\boldsymbol{\lambda} \geq \mathbf{0} \\ \boldsymbol{\lambda}' \geq \mathbf{0} \\ \boldsymbol{\pi} \geq \mathbf{0}}} \min_{\tilde{\mathbf{x}}, \xi, \mathbf{s}} \mathcal{L}(\tilde{\mathbf{x}}, \xi, \mathbf{s}; \boldsymbol{\lambda}, \boldsymbol{\lambda}', \boldsymbol{\pi}).$$

The inner minimization in each  $s_i$  is

$$\min_{s_i} \gamma s_i - \begin{cases} \lambda_i s_i, & y_i = 1, \\ \lambda'_i s_i, & y_i = -1 \end{cases} - \pi_i s_i,$$

which is  $-\infty$  unless

$$\gamma - \begin{cases} \lambda_i, \\ \lambda'_i \end{cases} - \pi_i = 0.$$

These are **hidden constraints** on dual variables; they must hold at any optimal dual solution. Under these, the  $s_i$  disappear.

We then minimize over  $\xi$ ; the Lagrangian is linear in  $\xi$ , so the minimum is finite only if

$$\sum_{y_i=1} \lambda_i = \sum_{y_i=-1} \lambda'_i,$$

another hidden constraint. With both hidden constraints satisfied, the inner problem reduces to

$$\min_{\tilde{\mathbf{x}}} \frac{1}{2} \|\tilde{\mathbf{x}}\|^2 + \mathbf{p}^\top \tilde{\mathbf{x}} + c,$$

where

$$\mathbf{p} = \sum_{y_i=-1} \lambda'_i \Phi(\mathbf{a}_i) - \sum_{y_i=1} \lambda_i \Phi(\mathbf{a}_i), \quad c = \sum_{y_i=1} \lambda_i + \sum_{y_i=-1} \lambda'_i.$$

The minimizer is

$$\tilde{\mathbf{x}}^* = -\mathbf{p},$$

and the optimal inner value is

$$-\frac{1}{2} \|\mathbf{p}\|^2 + c.$$

Thus the dual becomes

$$\begin{aligned} \max_{\substack{\lambda \geq 0 \\ \lambda' \geq 0 \\ \pi \geq 0}} & -\frac{1}{2} \left\| \sum_{y_i=1} \lambda_i \Phi(\mathbf{a}_i) - \sum_{y_i=-1} \lambda'_i \Phi(\mathbf{a}_i) \right\|^2 + \sum_{y_i=1} \lambda_i + \sum_{y_i=-1} \lambda'_i \\ \text{s.t.} & \begin{cases} \lambda \geq \mathbf{0}, \lambda' \geq \mathbf{0}, \pi \geq \mathbf{0}, \\ \lambda_i + \pi_i = \gamma, \quad \lambda'_i + \pi_i = \gamma, \quad \forall i, \\ \sum_{y_i=1} \lambda_i = \sum_{y_i=-1} \lambda'_i. \end{cases} \end{aligned}$$

Eliminating  $\pi$  via  $\pi_i = \gamma - \lambda_i = \gamma - \lambda'_i$  gives the simpler dual:

$$\begin{aligned} \max_{\lambda, \lambda'} & -\frac{1}{2} \left\| \sum_{y_i=1} \lambda_i \Phi(\mathbf{a}_i) - \sum_{y_i=-1} \lambda'_i \Phi(\mathbf{a}_i) \right\|^2 + \sum_{y_i=1} \lambda_i + \sum_{y_i=-1} \lambda'_i \\ \text{s.t.} & \lambda \geq \mathbf{0}, \quad \lambda' \geq \mathbf{0}, \quad \lambda \leq \gamma \mathbf{e}, \quad \lambda' \leq \gamma \mathbf{e}, \\ & \sum_{y_i=1} \lambda_i = \sum_{y_i=-1} \lambda'_i. \end{aligned}$$

Let  $\mathbf{K}$  denote the quadratic coefficients in  $(\lambda, \lambda')$  in the dual objective. In terms of feature vectors,

$$\mathbf{K} = \begin{bmatrix} \Phi^\top(\mathbf{a}_1)\Phi(\mathbf{a}_1) & \cdots & \pm \Phi^\top(\mathbf{a}_1)\Phi(\mathbf{a}_m) \\ \vdots & \ddots & \vdots \\ \pm \Phi^\top(\mathbf{a}_m)\Phi(\mathbf{a}_1) & \cdots & \Phi^\top(\mathbf{a}_m)\Phi(\mathbf{a}_m) \end{bmatrix},$$

with signs depending on the labels  $y_i$  (same labels  $\Rightarrow +$ , different labels  $\Rightarrow -$ ).

## LEC 20

## 20.1 Singular Value Decomposition

Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , we want a matrix  $\tilde{\mathbf{A}}$  such that  $\tilde{\mathbf{A}} \approx \mathbf{A}$  and  $\text{rank}(\tilde{\mathbf{A}}) = r \ll \min(m, n)$ .

**Theorem****(SVD.)**

For any matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , there exist orthogonal  $\mathbf{U} \in \mathbb{R}^{m \times m}$ ,  $\mathbf{V} \in \mathbb{R}^{n \times n}$  and a diagonal  $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$  such that

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top.$$

The diagonal entries of  $\mathbf{\Sigma}$ , denoted  $\sigma_1, \dots, \sigma_{\min(m, n)}$ , are the **singular values** of  $\mathbf{A}$  and satisfy

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min(m, n)} \geq 0.$$

The singular values are uniquely determined by  $\mathbf{A}$ .

Properties of the SVD:

1. If  $\sigma_1 \geq \dots \geq \sigma_r \geq 0$  and  $\sigma_{r+1} = \dots = \sigma_{\min(m, n)} = 0$ , then  $\text{rank}(\mathbf{A}) = r$ .
2.  $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top \iff \mathbf{A}^\top = \mathbf{V} \mathbf{\Sigma}^\top \mathbf{U}^\top$ .
3. If  $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$ , then

$$\mathbf{A}^\top \mathbf{A} = \mathbf{V} \mathbf{\Sigma}^\top \mathbf{\Sigma} \mathbf{V}^\top,$$

where  $\mathbf{\Sigma}^\top \mathbf{\Sigma}$  is square diagonal. Thus the singular values of  $\mathbf{A}$  are the square roots of the eigenvalues of  $\mathbf{A}^\top \mathbf{A}$ .

4. If  $\text{rank}(\mathbf{A}) = r$ , then

$$\mathbf{A} = \mathbf{U}(:, 1:r) \mathbf{\Sigma}(1:r, 1:r) \mathbf{V}(:, 1:r)^\top = \sum_{i=1}^r \mathbf{U}(:, i) \sigma_i \mathbf{V}(:, i)^\top,$$

a sum of  $r$  rank-1 matrices.

5. If  $\mathbf{B} := \mathbf{Q} \mathbf{A} \mathbf{Z}^\top$  with  $\mathbf{Q}, \mathbf{Z}$  orthogonal and  $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$  is an SVD, then

$$\mathbf{B} = (\mathbf{Q} \mathbf{U}) \mathbf{\Sigma} (\mathbf{V} \mathbf{Z})^\top$$

is an SVD of  $\mathbf{B}$ . Hence  $\mathbf{A}$  and  $\mathbf{B}$  have the same singular values.

6. If  $\mathbf{D}$  is diagonal, then its singular values are the absolute values of its diagonal entries, sorted in nonincreasing order.

**Theorem****(Eckart–Young.)**

Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $0 \leq \rho \leq \min(m, n)$ . Then

$$\min\{\|\mathbf{A} - \mathbf{X}\|_F : \mathbf{X} \in \mathbb{R}^{m \times n}, \text{rank}(\mathbf{X}) \leq \rho\} = \sqrt{\sigma_{\rho+1}^2 + \dots + \sigma_{\min(m, n)}^2},$$

and a minimizer is

$$\mathbf{X} = \mathbf{U}(:, 1:\rho) \mathbf{\Sigma}(1:\rho, 1:\rho) \mathbf{V}(:, 1:\rho)^\top,$$

where  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  is an SVD.

## 20.2 Latent Semantic Indexing

(Deerwester et al., 1990)

Given a corpus of documents, form the **term–document matrix**  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , where  $m$  is the number of terms in the language and  $n$  is the number of documents. The entry  $\mathbf{A}(i, j)$  is the number of occurrences of term  $i$  in document  $j$  (after normalization).

Compute the SVD of  $\mathbf{A}$  and keep the rank- $\rho$  approximation

$$\hat{\mathbf{A}} = \mathbf{U}(:, 1:\rho) \mathbf{\Sigma}(1:\rho, 1:\rho) \mathbf{V}(:, 1:\rho)^\top,$$

where  $\rho$  is an estimate of the number of topics in the corpus.

Given a query vector  $\mathbf{q} \in \mathbb{R}^m$ , the vector

$$\mathbf{q}^\top \hat{\mathbf{A}} \in \mathbb{R}^n$$

gives scores that can be interpreted as the relevance of each document to the query.

## 20.3 Matrix Completion Problem

Let  $\mathbf{M} \in (\mathbb{R} \cup \{?\})^{m \times n}$  be a partially observed matrix. We want to fill in the missing entries “?” with numbers.

**Application:** recommender systems.

Rows of  $\mathbf{M}$  correspond to users, columns to products, and  $\mathbf{M}(i, j)$  is the rating user  $i$  gives to product  $j$ . Most entries are missing; the vendor wants to predict missing entries to recommend products.

A reasonable model: assume the completed  $\mathbf{M}$  has low rank.

Intuition: each user can be modeled as a linear combination of a small number of *prototypical users* (unknown). If ratings are linear combinations of the ratings given by prototypical users, then the completed  $\mathbf{M}$  will have rank equal to the number of prototypical users.

Mathematical formulation:

Find  $\mathbf{X} \in \mathbb{R}^{m \times n}$  such that  $\mathbf{X}(\Omega) = \mathbf{M}(\Omega)$ , where

$$\Omega = \{(i, j) : \mathbf{M}(i, j) \neq ?\},$$

and subject to this constraint,  $\text{rank}(\mathbf{X})$  is minimized. This problem is **NP-hard**.

Candès and Recht showed that under certain assumptions on  $\mathbf{M}$  and  $\Omega$ , the matrix  $\mathbf{X}$  solving this NP-hard problem is also a solution of the convex optimization problem

$$\min \|\mathbf{X}\|_* \quad \text{s.t.} \quad \mathbf{X}(\Omega) = \mathbf{M}(\Omega),$$

where  $\|\mathbf{X}\|_*$  is the **nuclear norm** of  $\mathbf{X}$ ,

$$\|\mathbf{X}\|_* = \sigma_1(\mathbf{X}) + \cdots + \sigma_{\min(m,n)}(\mathbf{X}).$$

We check that this is (indeed) a norm.

A norm  $\|\cdot\|$  must satisfy:

1.  $\|\mathbf{X}\| \geq 0$  and  $\|\mathbf{X}\| = 0 \iff \mathbf{X} = \mathbf{0}$ .
2.  $\|\lambda\mathbf{X}\| = |\lambda| \|\mathbf{X}\|$  for all scalars  $\lambda$ .
3.  $\|\mathbf{X} + \mathbf{Y}\| \leq \|\mathbf{X}\| + \|\mathbf{Y}\|$  (triangle inequality).

### Theorem

Let  $\|\cdot\|_\square$  be a norm on  $\mathbb{R}^m$  and  $\|\cdot\|_\triangle$  a norm on  $\mathbb{R}^n$ . The  $\triangle \rightarrow \square$  **operator norm** on  $\mathbb{R}^{m \times n}$  is defined by

$$\|\mathbf{A}\|_{\triangle \rightarrow \square} := \sup\{\|\mathbf{A}\mathbf{v}\|_\square : \|\mathbf{v}\|_\triangle \leq 1\}.$$

This is a norm on  $\mathbb{R}^{m \times n}$  (an **induced norm**).

### Theorem

For  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , the largest singular value  $\sigma_1(\mathbf{A})$  is the operator norm induced by the Euclidean norm  $\|\cdot\|_2$ , i.e.

$$\|\mathbf{A}\|_{2 \rightarrow 2} = \sigma_1(\mathbf{A}).$$

### Proof

We show that

$$\sup_{\|\mathbf{x}\|_2 \leq 1} \|\mathbf{A}\mathbf{x}\|_2 = \sigma_1(\mathbf{A}).$$

First,  $\|\mathbf{A}\mathbf{x}\|_2 \leq \sigma_1(\mathbf{A})$  for all  $\|\mathbf{x}\|_2 \leq 1$ .

Let  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  be an SVD. For any  $\mathbf{x}$  with  $\|\mathbf{x}\|_2 \leq 1$ ,

$$\|\mathbf{A}\mathbf{x}\|_2 = \|\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top \mathbf{x}\|_2 = \|\mathbf{\Sigma}\mathbf{V}^\top \mathbf{x}\|_2,$$

since orthogonal matrices preserve the 2-norm. Let  $\mathbf{w} = \mathbf{V}^\top \mathbf{x}$ . Then  $\|\mathbf{w}\|_2 = \|\mathbf{x}\|_2 \leq 1$ , and if  $m \geq n$ ,

$$\|\mathbf{\Sigma}\mathbf{w}\|_2 = \left(\sigma_1^2 w_1^2 + \cdots + \sigma_n^2 w_n^2\right)^{1/2} \leq \sigma_1 \left(w_1^2 + \cdots + w_n^2\right)^{1/2} \leq \sigma_1.$$

Thus  $\|\mathbf{A}\mathbf{x}\|_2 \leq \sigma_1$  for all  $\|\mathbf{x}\|_2 \leq 1$ , so  $\|\mathbf{A}\|_2 \leq \sigma_1(\mathbf{A})$ .

For the reverse inequality, take

$$\mathbf{x} = \mathbf{V}\mathbf{e}_1, \quad \mathbf{e}_1 = (1, 0, \dots, 0)^\top.$$

Then  $\|\mathbf{x}\|_2 = 1$  and

$$\|\mathbf{A}\mathbf{x}\|_2 = \|\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top \mathbf{V}\mathbf{e}_1\|_2 = \|\mathbf{U}\mathbf{\Sigma}\mathbf{e}_1\|_2 = \|\mathbf{\Sigma}\mathbf{e}_1\|_2 = \|(\sigma_1, 0, \dots, 0)^\top\|_2 = \sigma_1.$$

Hence  $\|\mathbf{A}\|_2 \geq \sigma_1(\mathbf{A})$ , and combining both directions gives  $\|\mathbf{A}\|_2 = \sigma_1(\mathbf{A})$ .

**Theorem**

Let  $\|\cdot\|_{\square}$  be a norm on  $\mathbb{R}^m$  (or on  $\mathbb{R}^{m \times n}$ ). The function

$$\mathbf{x} \mapsto \sup\{\langle \mathbf{x}, \mathbf{y} \rangle : \|\mathbf{y}\|_{\square} \leq 1\}$$

is the **dual norm** of  $\|\cdot\|_{\square}$  and is itself a norm.

**Proof**

This is the operator norm of the linear map  $\mathbf{y} \mapsto \langle \mathbf{x}, \mathbf{y} \rangle$ , so it satisfies the norm axioms.

For vectors:

1. The dual norm of the Euclidean norm is itself (Cauchy–Schwarz inequality).
2. The dual norm of  $\|\mathbf{x}\|_1$  is  $\|\mathbf{x}\|_{\infty}$ .

**Theorem**

For matrices, the dual norm of  $\|\cdot\|_2$  (spectral norm) is the nuclear norm  $\|\cdot\|_*$ .

**Proof**

See Problem Set 6.

Back to Candès and Recht: why does

$$\min \|\mathbf{X}\|_* \quad \text{s.t.} \quad \mathbf{X}(\Omega) = \mathbf{M}(\Omega)$$

typically yield a low-rank  $\mathbf{X}$ ?

**Theorem**

Let  $\mathbf{A} \in \mathbb{R}^{m \times n} \setminus \{\mathbf{0}\}$  with  $m \geq n$ . Consider

$$\min \|\mathbf{X}\|_* \quad \text{s.t.} \quad \langle \mathbf{A}, \mathbf{X} \rangle = 1.$$

Then an optimal solution is

$$\hat{\mathbf{X}} = \mathbf{U} \begin{bmatrix} \frac{1}{\sigma_1(\mathbf{A})} & & & \\ & 0 & & \\ & & \ddots & \\ & & & 0 \end{bmatrix} \mathbf{V}^{\top},$$

where  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^{\top}$  is an SVD of  $\mathbf{A}$ .

## Proof

First, check that  $\hat{\mathbf{X}}$  is feasible:

$$\begin{aligned}
 \langle \mathbf{A}, \hat{\mathbf{X}} \rangle &= \left\langle \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top, \mathbf{U} \begin{bmatrix} \frac{1}{\sigma_1(\mathbf{A})} & & \\ & 0 & \\ & & \ddots \\ & & & 0 \end{bmatrix} \mathbf{V}^\top \right\rangle \\
 &= \text{Tr} \left( \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top \mathbf{V} \begin{bmatrix} \frac{1}{\sigma_1(\mathbf{A})} & & \\ & 0 & \\ & & \ddots \\ & & & 0 \end{bmatrix}^\top \mathbf{U}^\top \right) \\
 &= \text{Tr} \left( \mathbf{\Sigma} \begin{bmatrix} \frac{1}{\sigma_1(\mathbf{A})} & & \\ & 0 & \\ & & \ddots \\ & & & 0 \end{bmatrix}^\top \right) \quad (\text{cyclic invariance of trace, and } \mathbf{U}^\top \mathbf{U} = \mathbf{I}) \\
 &= \text{Tr} \begin{bmatrix} 1 & & \\ & 0 & \\ & & \ddots \\ & & & 0 \end{bmatrix} = 1.
 \end{aligned}$$

Next, note that  $\|\hat{\mathbf{X}}\|_2 = \frac{1}{\sigma_1(\mathbf{A})}$ .

For any other feasible  $\mathbf{X}$  (i.e.  $\langle \mathbf{A}, \mathbf{X} \rangle = 1$ ), the nuclear norm can be written using the dual norm relationship:

$$\|\mathbf{X}\|_* = \sup \{ \langle \mathbf{X}, \mathbf{Y} \rangle : \|\mathbf{Y}\|_2 \leq 1 \}.$$

Take

$$\bar{\mathbf{Y}} = \frac{\mathbf{A}}{\sigma_1(\mathbf{A})} = \frac{\mathbf{A}}{\|\mathbf{A}\|_2},$$

so  $\|\bar{\mathbf{Y}}\|_2 = 1$ . Then

$$\|\mathbf{X}\|_* \geq \langle \mathbf{X}, \bar{\mathbf{Y}} \rangle = \left\langle \mathbf{X}, \frac{\mathbf{A}}{\sigma_1(\mathbf{A})} \right\rangle = \frac{1}{\sigma_1(\mathbf{A})} \langle \mathbf{X}, \mathbf{A} \rangle = \frac{1}{\sigma_1(\mathbf{A})}.$$

Thus any feasible  $\mathbf{X}$  has nuclear norm at least  $1/\sigma_1(\mathbf{A})$ , which is exactly  $\|\hat{\mathbf{X}}\|_*$ . Therefore  $\hat{\mathbf{X}}$  is optimal.



## LEC 21

## 21.1 Matrix Completion Problem

Given partially specified  $\mathbf{M} \in (\mathbb{R} \cup \{?\})^{m \times n}$ , the goal is to fill in the missing entries to obtain a completion with as low rank as possible.

**Theorem**

**Candès–Recht Relaxation (CR Relaxation).**

$$\min_{\mathbf{X}} \|\mathbf{X}\|_* \quad \text{subject to} \quad \mathbf{X}(\Omega) = \mathbf{M}(\Omega),$$

where

$$\Omega = \{(i, j) \in \{1, \dots, m\} \times \{1, \dots, n\} : \mathbf{M}(i, j) \neq ?\}.$$

Extension of CR Relaxation to noisy  $\mathbf{M}$ :

$$\min_{\mathbf{X}} \frac{1}{2} \sum_{(i,j) \in \Omega} (\mathbf{X}(i, j) - \mathbf{M}(i, j))^2 + \gamma \|\mathbf{X}\|_*.$$

Norms are proper, closed, convex functions and are typically nonsmooth, so we are in a good setting for proximal / first-order methods.

We will use *APGD* and therefore need prox for the nuclear norm  $\|\cdot\|_*$ .

**Theorem**

Let  $\mathbf{X} \in \mathbb{R}^{m \times n}$  with  $m \geq n$ , and let  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  be its SVD with singular values  $\sigma_1, \dots, \sigma_n$  on the diagonal of  $\mathbf{\Sigma}$ . Then

$$\text{prox}_{t\|\cdot\|_*}(\mathbf{X}) = \mathbf{U} \text{diag}(\max(\sigma_1 - t, 0), \dots, \max(\sigma_n - t, 0)) \mathbf{V}^\top.$$

That is, the proximal operator for  $\|\cdot\|_*$  performs soft-thresholding on the singular values.

**Proof**

For matrices, the proximal operator is

$$\text{prox}_{t\|\cdot\|_*}(\mathbf{X}) = \arg \min_{\mathbf{Z}} \left\{ t\|\mathbf{Z}\|_* + \frac{1}{2}\|\mathbf{Z} - \mathbf{X}\|_F^2 \right\}.$$

Use unitary invariance of both  $\|\cdot\|_F$  and  $\|\cdot\|_*$ . Write  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  with  $\mathbf{U}, \mathbf{V}$  orthogonal, and set

$$\mathbf{W} := \mathbf{U}^\top \mathbf{Z} \mathbf{V}.$$

Then

$$\|\mathbf{Z} - \mathbf{X}\|_F = \|\mathbf{U}^\top \mathbf{Z} \mathbf{V} - \mathbf{\Sigma}\|_F = \|\mathbf{W} - \mathbf{\Sigma}\|_F, \quad \|\mathbf{Z}\|_* = \|\mathbf{U}^\top \mathbf{Z} \mathbf{V}\|_* = \|\mathbf{W}\|_*.$$

Hence

$$\arg \min_{\mathbf{Z}} \left\{ t\|\mathbf{Z}\|_* + \frac{1}{2}\|\mathbf{Z} - \mathbf{X}\|_F^2 \right\} = \mathbf{U} \left( \arg \min_{\mathbf{W}} \left\{ t\|\mathbf{W}\|_* + \frac{1}{2}\|\mathbf{W} - \mathbf{\Sigma}\|_F^2 \right\} \right) \mathbf{V}^\top.$$

So it suffices to solve

$$\min_{\mathbf{W}} t\|\mathbf{W}\|_* + \frac{1}{2}\|\mathbf{W} - \mathbf{\Sigma}\|_{\text{F}}^2.$$

### Lemma

For any  $\mathbf{W} \in \mathbb{R}^{m \times n}$  with  $m \geq n$ ,

$$\|\mathbf{W}\|_* \geq |\mathbf{W}(1,1)| + \cdots + |\mathbf{W}(n,n)|.$$

### Proof

#### Proof of Lemma.

By the dual characterization of the nuclear norm,

$$\|\mathbf{W}\|_* = \sup\{\langle \mathbf{W}, \mathbf{Y} \rangle : \|\mathbf{Y}\|_2 \leq 1\}.$$

Choose

$$\hat{\mathbf{Y}} = \text{diag}(\pm 1, \dots, \pm 1),$$

with signs chosen so that the  $(i, i)$  entry of  $\hat{\mathbf{Y}}$  has the same sign as  $\mathbf{W}(i, i)$ . Then  $\|\hat{\mathbf{Y}}\|_2 = 1$  and

$$\|\mathbf{W}\|_* \geq \langle \mathbf{W}, \hat{\mathbf{Y}} \rangle = \sum_{i=1}^n \mathbf{W}(i, i) \hat{\mathbf{Y}}(i, i) = \sum_{i=1}^n |\mathbf{W}(i, i)|.$$

Now consider the objective

$$\Phi(\mathbf{W}) := t\|\mathbf{W}\|_* + \frac{1}{2}\|\mathbf{W} - \mathbf{\Sigma}\|_{\text{F}}^2.$$

Write  $\mathbf{W}_0$  for the diagonal matrix formed from the diagonal of  $\mathbf{W}$ . From the lemma,  $\|\mathbf{W}\|_* \geq \|\mathbf{W}_0\|_*$ , and since  $\mathbf{\Sigma}$  is diagonal,

$$\|\mathbf{W} - \mathbf{\Sigma}\|_{\text{F}}^2 = \|\mathbf{W}_0 - \mathbf{\Sigma}\|_{\text{F}}^2 + (\text{nonnegative off-diagonal terms}),$$

so  $\|\mathbf{W} - \mathbf{\Sigma}\|_{\text{F}} \geq \|\mathbf{W}_0 - \mathbf{\Sigma}\|_{\text{F}}$  with equality only when  $\mathbf{W}$  is diagonal. Thus  $\Phi(\mathbf{W}) \geq \Phi(\mathbf{W}_0)$ , and any minimizer of  $\Phi$  must be diagonal.

Let  $\mathbf{W} = \text{diag}(w_1, \dots, w_n)$  and  $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_n)$ . Then

$$\|\mathbf{W}\|_* = \sum_{i=1}^n |w_i|, \quad \|\mathbf{W} - \mathbf{\Sigma}\|_{\text{F}}^2 = \sum_{i=1}^n (w_i - \sigma_i)^2.$$

So the problem decouples across  $i$ :

$$\min_{\mathbf{w} \in \mathbb{R}^n} \left\{ t\|\mathbf{w}\|_1 + \frac{1}{2}\|\mathbf{w} - \boldsymbol{\sigma}\|_2^2 \right\}, \quad \boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_n)^\top.$$

This is exactly the proximal operator of  $t\|\cdot\|_1$  at  $\boldsymbol{\sigma}$ , whose solution is componentwise soft-thresholding:

$$w_i = \begin{cases} \sigma_i - t, & \sigma_i \geq t, \\ 0, & \sigma_i < t. \end{cases}$$

Hence

$$\text{prox}_{t\|\cdot\|_*}(\mathbf{X}) = \mathbf{U} \text{diag}(\max(\sigma_i - t, 0))_{i=1}^n \mathbf{V}^\top,$$

which completes the proof.

## 21.2 Economy-sized SVD

Given  $\mathbf{X} \in \mathbb{R}^{m \times n}$  with  $m \geq n$ ,

$$\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top = \mathbf{U}(:, 1:n) \mathbf{\Sigma}(1:n, 1:n) \mathbf{V}(:, 1:n)^\top.$$

This is the **economy-sized SVD**.

In practice, for the nuclear-norm proximal operator we only need the nonzero singular values and corresponding singular vectors, so we use the economy-sized SVD, especially when  $m \gg n$ .

## 21.3 Rank–Sparsity Decomposition

Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , seek  $\mathbf{L}, \mathbf{M} \in \mathbb{R}^{m \times n}$  such that

$$\mathbf{A} = \mathbf{L} + \mathbf{M},$$

with

- $\mathbf{L}$  sparse (most entries zero),
- $\mathbf{M}$  low rank ( $\text{rank}(\mathbf{M}) \ll \min(m, n)$ ).

*Application:* background identification / removal from video frames (Problem Set 6). Think of columns of  $\mathbf{A}$  as vectorized frames: background is nearly low rank, moving objects are sparse.

If  $\mathbf{A}$  itself were simultaneously very sparse and very low rank, the decomposition would be ill-posed. So we assume  $\mathbf{L}$  is sparse and high rank, while  $\mathbf{M}$  is dense and low rank.

### Theorem

**Convex relaxation** (Candès et al. 2011; Chandrasekaran et al. 2011).

$$\min_{\mathbf{L}, \mathbf{M}} \|\mathbf{L}\|_{\ell_1} + \gamma \|\mathbf{M}\|_* \quad \text{subject to} \quad \mathbf{L} + \mathbf{M} = \mathbf{A},$$

where

$$\|\mathbf{L}\|_{\ell_1} = \sum_{i=1}^m \sum_{j=1}^n |\mathbf{L}(i, j)|, \quad \gamma > 0.$$

We solve this with **ADMM** (Alternating Direction Method of Multipliers).

### 21.3.1 ADMM for two blocks

General form:

$$\min_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}) + g(\mathbf{y}) \quad \text{subject to} \quad \mathbf{Ax} + \mathbf{By} = \mathbf{c},$$

where  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{p \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{p \times m}$ ,  $\mathbf{c} \in \mathbb{R}^p$ .

**Augmented Lagrangian:**

$$\mathcal{L}_A(\mathbf{x}, \mathbf{y}; \boldsymbol{\mu}) = f(\mathbf{x}) + g(\mathbf{y}) + \boldsymbol{\mu}^\top (\mathbf{Ax} + \mathbf{By} - \mathbf{c}) + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{By} - \mathbf{c}\|_2^2,$$

where  $\rho > 0$ .

**ADMM iteration:**

$$\begin{aligned} &\mathbf{y}^0 \text{ arbitrary, } \boldsymbol{\mu}^0 \text{ arbitrary,} \\ &\text{for } k = 0, 1, 2, \dots \\ &\quad \mathbf{x}^{k+1} := \arg \min_{\mathbf{x}} \mathcal{L}_A(\mathbf{x}, \mathbf{y}^k; \boldsymbol{\mu}^k), \\ &\quad \mathbf{y}^{k+1} := \arg \min_{\mathbf{y}} \mathcal{L}_A(\mathbf{x}^{k+1}, \mathbf{y}; \boldsymbol{\mu}^k), \\ &\quad \boldsymbol{\mu}^{k+1} := \boldsymbol{\mu}^k + \rho (\mathbf{Ax}^{k+1} + \mathbf{By}^{k+1} - \mathbf{c}). \end{aligned}$$

### 21.3.2 Specialization to Rank–Sparsity

Here the constraint is  $\mathbf{L} + \mathbf{M} = \mathbf{A}$ ; the augmented Lagrangian is

$$\mathcal{L}_A(\mathbf{L}, \mathbf{M}; \boldsymbol{\Lambda}) = \|\mathbf{L}\|_{\ell_1} + \gamma \|\mathbf{M}\|_* + \langle \boldsymbol{\Lambda}, \mathbf{L} + \mathbf{M} - \mathbf{A} \rangle + \frac{\rho}{2} \|\mathbf{L} + \mathbf{M} - \mathbf{A}\|_F^2,$$

where  $\boldsymbol{\Lambda}$  is the matrix of Lagrange multipliers and  $\langle \cdot, \cdot \rangle$  is the Frobenius inner product.

**Step 1:  $\mathbf{L}$ -update.**

$$\begin{aligned} \mathbf{L}^{k+1} &:= \arg \min_{\mathbf{L}} \mathcal{L}_A(\mathbf{L}, \mathbf{M}^k; \boldsymbol{\Lambda}^k) \\ &= \arg \min_{\mathbf{L}} \left\{ \|\mathbf{L}\|_{\ell_1} + \langle \boldsymbol{\Lambda}^k, \mathbf{L} \rangle + \frac{\rho}{2} \|\mathbf{L} + \mathbf{M}^k - \mathbf{A}\|_F^2 + (\text{constants}) \right\}. \end{aligned}$$

Complete the square:

$$\mathbf{L}^{k+1} = \arg \min_{\mathbf{L}} \left\{ \|\mathbf{L}\|_{\ell_1} + \frac{\rho}{2} \left\| \mathbf{L} + \mathbf{M}^k - \mathbf{A} + \frac{1}{\rho} \boldsymbol{\Lambda}^k \right\|_F^2 + (\text{constants}) \right\}.$$

Thus

$$\mathbf{L}^{k+1} = \text{prox}_{\frac{1}{\rho} \|\cdot\|_{\ell_1}} \left( -\mathbf{M}^k + \mathbf{A} - \frac{1}{\rho} \boldsymbol{\Lambda}^k \right),$$

i.e. entrywise soft-thresholding with parameter  $1/\rho$ .

**Step 2:  $\mathbf{M}$ -update.**

$$\begin{aligned} \mathbf{M}^{k+1} &:= \arg \min_{\mathbf{M}} \mathcal{L}_A(\mathbf{L}^{k+1}, \mathbf{M}; \boldsymbol{\Lambda}^k) \\ &= \arg \min_{\mathbf{M}} \left\{ \gamma \|\mathbf{M}\|_* + \langle \boldsymbol{\Lambda}^k, \mathbf{M} \rangle + \frac{\rho}{2} \|\mathbf{L}^{k+1} + \mathbf{M} - \mathbf{A}\|_F^2 \right\} \\ &= \arg \min_{\mathbf{M}} \left\{ \gamma \|\mathbf{M}\|_* + \frac{\rho}{2} \left\| \mathbf{M} + \mathbf{L}^{k+1} - \mathbf{A} + \frac{1}{\rho} \boldsymbol{\Lambda}^k \right\|_F^2 \right\}. \end{aligned}$$

So

$$\mathbf{M}^{k+1} = \text{prox}_{\frac{\gamma}{\rho}\|\cdot\|_*} \left( -\mathbf{L}^{k+1} + \mathbf{A} - \frac{1}{\rho}\mathbf{\Lambda}^k \right),$$

i.e. singular value soft-thresholding (using economy-sized SVD).

**Step 3: Multiplier update.**

$$\mathbf{\Lambda}^{k+1} := \mathbf{\Lambda}^k + \rho(\mathbf{L}^{k+1} + \mathbf{M}^{k+1} - \mathbf{A}).$$

### 21.3.3 ADMM Motivation: Method of Multipliers

Consider

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{h}(\mathbf{x}) = \mathbf{0},$$

with  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ .

The augmented Lagrangian is

$$\mathcal{L}_A(\mathbf{x}; \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\mu}^\top \mathbf{h}(\mathbf{x}) + \frac{1}{2}\rho \|\mathbf{h}(\mathbf{x})\|_2^2.$$

Given  $\hat{\boldsymbol{\mu}}$ , define  $\hat{\mathbf{x}}$  by

$$\nabla_{\mathbf{x}} \mathcal{L}_A(\hat{\mathbf{x}}; \hat{\boldsymbol{\mu}}) = \mathbf{0},$$

that is,

$$\nabla f(\hat{\mathbf{x}}) + \mathbf{J}(\hat{\mathbf{x}})^\top \hat{\boldsymbol{\mu}} + \rho \mathbf{J}(\hat{\mathbf{x}})^\top \mathbf{h}(\hat{\mathbf{x}}) = \mathbf{0},$$

where

$$\mathbf{J}(\hat{\mathbf{x}}) = \begin{pmatrix} \nabla h_1(\hat{\mathbf{x}})^\top \\ \vdots \\ \nabla h_p(\hat{\mathbf{x}})^\top \end{pmatrix} \in \mathbb{R}^{p \times n}$$

is the Jacobian.

Compare with the KKT conditions:

$$\begin{aligned} \mathbf{h}(\mathbf{x}^*) &= \mathbf{0}, \\ \nabla f(\mathbf{x}^*) + \mathbf{J}(\mathbf{x}^*)^\top \boldsymbol{\mu}^* &= \mathbf{0}. \end{aligned}$$

If we define

$$\boldsymbol{\mu}^{\text{new}} := \hat{\boldsymbol{\mu}} + \rho \mathbf{h}(\hat{\mathbf{x}}),$$

the stationarity equation above resembles the KKT stationarity equation with updated multiplier.

Thus the **method of multipliers** (augmented Lagrangian method) updates

$$\begin{aligned} \mathbf{x}^{k+1} &:= \arg \min_{\mathbf{x}} \mathcal{L}_A(\mathbf{x}; \boldsymbol{\mu}^k), \\ \boldsymbol{\mu}^{k+1} &:= \boldsymbol{\mu}^k + \rho \mathbf{h}(\mathbf{x}^{k+1}). \end{aligned}$$

ADMM can be viewed as applying this idea but approximating the primal minimization by alternating minimization over blocks of variables ( $\mathbf{x}$ -block,  $\mathbf{y}$ -block), which leads to the simple, proximal-based updates used in rank-sparsity decomposition.

## LEC 22

## 22.1 CVX

## 22.1.1 Semidefinite Programming

CVX is a modeling language for convex optimization.

It accepts a high-level description of an optimization problem, and translates it to an input suitable for a solver; invokes solver; returns optimal values to user's code.

The solvers in CVX are for **semidefinite programming (SDP)**.

**SDP**

$$\min \langle \mathbf{C}, \mathbf{X} \rangle \quad \text{S.T.} \quad \langle \mathbf{A}_1, \mathbf{X} \rangle = b_1, \dots, \langle \mathbf{A}_m, \mathbf{X} \rangle = b_m, \mathbf{X} \succeq 0$$

where  $\mathbf{X} \in \mathbb{S}^n$ , given  $\mathbf{C}, \mathbf{A}_1, \dots, \mathbf{A}_m \in \mathbb{S}^n$ ,  $b_1, \dots, b_m \in \mathbb{R}$ .

**Fact**

**Fact.** All of the convex programs considered this semester ( $\ell_1$  LS, SVM, Huber, Candès-Recht, Rank-Sparsity) are special cases of SDP.

CVX is equipped with rules for translating many convex problems to SDP.

## 22.1.2 Examples

**Example**

(Setting up Candès-Recht relaxation of matrix completion in CVX.)

Assume we have  $\mathbf{M} \in \mathbb{R}^{m \times n}$ , say  $\mathbf{M}(i, j) = \text{NaN}$  for missing entries. So  $\mathbf{M}, m, n$  are already program variables.

```
cvx_begin
variable X(m, n)
minimize norm_nuc(X)
subject to
for i = 1 : m
    for j = 1 : n
        if ~ isnan(M(i, j))
            X(i, j) == M(i, j)
        end
    end
end
end
cvx_end
```

This code segment is inserted inside a MATLAB function. The variables  $\mathbf{M}$ ,  $m$ , and  $n$  are defined

before the segment. The variable  $\mathbf{X}$  will be defined upon completion.

### Example

(Noisy C-R relaxation.)

$\mathbf{M}, m, n$  as above,  $\gamma$  also defined.

```
[is, js] = find(~isnan(M))
nk = length(is)
cvx_begin
variable X(m, n)
variable z(nk)
minimize (quad_form(z, eye(nk)) / 2 + gamma * norm_nuc(X))
subject to
for k = 1 : nk
    z(k) == X(is(k), js(k)) - M(is(k), js(k))
end
cvx_end
```

CVX formulation of  $\min \frac{1}{2} \|\mathbf{z}\|^2 + \gamma \|\mathbf{X}\|_*$  subject to  $z[s(i, j)] = \mathbf{X}(i, j) - \mathbf{M}(i, j)$ ,  $\forall (i, j) \in \Omega$ , where  $s(i, j)$  is sequential numbering  $(i, j) \in \Omega$ .

CVX implements **Disciplined Convex Programming (DCP)**. CVX determines convexity using syntactic rules.

### Example

Consider

$$\min \langle \text{some function} \rangle + \gamma \sum_{i=1}^n t_i \quad \text{S.T.} \quad t_i \geq x_i, t_i \geq -x_i, \forall i = 1 : n.$$

OK in CVX: write as above.

NOT OK in CVX:  $t_i == \text{abs}(x_i)$ .

### 22.1.3 Second-order Methods

CVX standard solvers are **Interior-point Methods** (CO 663 / CO 666).

Interior-point methods are second-order methods. On each iteration, solve large system of linear equations. Obtain high accuracy, fast convergence, but each iteration is costly in terms of space and time.

The methods in this class are all first-order methods, that is, **gradient-based**, no matrices (except matrices native to problem like C-R relaxation). First-order methods have slower convergence, lower accuracy, but much faster iterations.

## 22.2 Nonconvex Optimization

### 22.2.1 Solvers

#### 1. Convex relaxation

- (i) C-R relaxation of matrix completion
- (ii) Relaxation of rank-sparsity
- (iii) Compressive sensing

Nonconvex problem:

$$\min \text{nnz}(\mathbf{x}) \quad \text{S.T.} \quad \mathbf{Ax} = \mathbf{b}$$

where  $\text{nnz}(\cdot)$  is the function of the number of nonzero entries. This is NP-hard.

Convex relaxation:

$$\min \|\mathbf{x}\|_1 \quad \text{S.T.} \quad \mathbf{Ax} = \mathbf{b}.$$

**2. Alternating algorithms** Regression: given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{y} \in \mathbb{R}^m$ , want

$$\mathbf{x} := \arg \min \{\|\mathbf{Ax} - \mathbf{y}\|\}.$$

Robust Regression:  $\ell$  entries of  $\mathbf{y}$  are arbitrarily corrupted. Still want  $\mathbf{x}$  fits the noncorrupted entries.

**AMRR-Alternating Minimization** for robust regression.

Initialize: guess  $\ell$  corrupted entries randomly, for example. Let  $S^0$  be the guess of uncorrupted entries, so

$$S^0 \subseteq \{1, \dots, m\}, |S^0| = m - \ell.$$

Now solve:

$$\mathbf{x}^0 := \arg \min \|\mathbf{A}(S^0, :)\mathbf{x} - \mathbf{y}(S^0)\|.$$

Let  $S^1 :=$  indices of the  $m - \ell$  smallest entries of  $|r_1^0|, \dots, |r_m^0|$  where  $\mathbf{r}^0 = \mathbf{Ax}^0 - \mathbf{y}$ .

Continue alternating: find  $\mathbf{x}^0, S^1, \mathbf{x}^1, S^2, \mathbf{x}^2, \dots$  until convergence.

#### Theorem

(Bhatia, Jain, Kar) AMRR yields an  $\mathbf{x}^k$  such that

$$\|\mathbf{x}^k - \mathbf{x}^*\| \leq \varepsilon$$

in  $\mathcal{O}(\log \|\mathbf{q}^*\|/\varepsilon)$  iterations under some fairly strong assumptions on  $\mathbf{A}$ . Here  $\mathbf{q}^*$  is the residual vector of  $\mathbf{x}^*$ .

From experience, AMRR performs well for a problem arising in an image application, where assumptions on  $\mathbf{A}$  may not hold.

**Nonnegative Matrix Factorization (NMF):**



(Gillis Book) Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{A} \geq 0$  (entrywise), given  $r \leq \min(m, n)$ . Seek  $\mathbf{W}, \mathbf{H}$  such that  $\mathbf{W} \in \mathbb{R}^{m \times r}$ ,  $\mathbf{H} \in \mathbb{R}^{r \times n}$ ,  $\mathbf{W} \geq 0, \mathbf{H} \geq 0, \mathbf{A} \approx \mathbf{WH}$ .

#### Example

(A term-doc matrix.) Interpretation: each document is a nonnegative sum of topics, so more interpretable than SVD.

Other applications of NMF: **Blind Sound Source Separation, Hyperspectral Imaging, Microarray Experiments.**

**3. *GD* / *SGD* / Conjugate Gradient** (Gradient Descent / Stochastic Gradient Descent / Conjugate Gradient.)

## LEC 23

## 23.1 Alternating algorithms

## 23.1.1 Nonnegative Matrix Factorization

## Theorem

Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , given  $r \leq \min(m, n)$  (some variants: NMF figures out  $r$ ), want  $\mathbf{W} \in \mathbb{R}^{m \times r}$ ,  $\mathbf{H} \in \mathbb{R}^{r \times n}$ ,  $\mathbf{W} \geq \mathbf{0}$ ,  $\mathbf{H} \geq \mathbf{0}$ ,  $\mathbf{A} \approx \mathbf{WH}$ .

## Claim

Geometry of NMF:

Say each column of  $\mathbf{A}$  normalized in 1-norm (sum of entries of each column equals 1)

In this case, we can assume that columns of  $\mathbf{W}$ ,  $\mathbf{H}$  are similarly normalized.

Why?

Say  $\mathbf{A} \approx \mathbf{WH}$ .

Rescale columns of  $\mathbf{W}$  to have 1-norm to 1, that is, define  $\bar{\mathbf{W}} = \mathbf{W}\mathbf{D}$ ,  $\mathbf{D}$  is diagonal,  $\mathbf{D} \succ \mathbf{0}$ , entries chosen so that  $\underbrace{\mathbf{e}^\top}_{\mathbb{R}^m} \bar{\mathbf{W}} = \underbrace{\mathbf{e}^\top}_{\mathbb{R}^r}$ .

Let  $\bar{\mathbf{H}} = \mathbf{D}^{-1}\mathbf{H}$  so that  $\mathbf{e}^\top \bar{\mathbf{W}} \bar{\mathbf{H}} = \mathbf{e}^\top \bar{\mathbf{H}}$ .

We know  $\bar{\mathbf{W}} \bar{\mathbf{H}} = \mathbf{WH} \approx \mathbf{A}$ , and  $\underbrace{\mathbf{e}^\top}_{\mathbb{R}^m} \mathbf{A} = \underbrace{\mathbf{e}^\top}_{\mathbb{R}^n}$  by assumption.

So may as well take columns of  $\bar{\mathbf{H}}$  to have 1-norm equal to 1. ( $\mathbf{e}^\top \bar{\mathbf{H}} \approx \mathbf{e}^\top \mathbf{A} = \mathbf{e}^\top$ )

Under assumption that columns of  $\mathbf{W}$ ,  $\mathbf{H}$ ,  $\mathbf{A}$  are all normalized, NMF looks as follows:

$$\Delta^m = m\text{-dimensional standard simplex} = \{\mathbf{x} \in \mathbb{R}^m : \mathbf{e}^\top \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}\}.$$

Columns of  $\mathbf{A}$  lie in  $\Delta^m$  as do columns of  $\mathbf{W}$ .

Columns of  $\mathbf{H}$  are coefficients of convex combinations.

$\mathbf{A} = \mathbf{WH}$  means each column of  $\mathbf{A}$  is a convex combination of columns of  $\mathbf{W}$ .

When posed as optimization, for example

$$\min \|\mathbf{A} - \mathbf{WH}\|_F \quad \text{s.t.} \quad \mathbf{W} \in \mathbb{R}^{m \times r}, \mathbf{H} \in \mathbb{R}^{r \times n}, \mathbf{W} \geq \mathbf{0}, \mathbf{H} \geq \mathbf{0}$$

is NP-hard.

Special case: columns of  $\mathbf{W}$  appear as columns of  $\mathbf{A}$ .

In context of term-doc matrices, for each topic there is a *pure document* only about that topic, and it is comprehensive.

This is called *Separable Case* in literature.

Gillis and Vavasis proposed algorithm **SPA** to solve separable with noise.

For nonseparable case: **Alternating Nonnegative Least Squares (ANLS)**.

If  $\mathbf{W}$  is given, optimal  $\mathbf{H}$  solves:

$$\arg \min_{\mathbf{H}} \{ \|\mathbf{A} - \mathbf{W}\mathbf{H}\|_{\text{F}} : \mathbf{H} \in \mathbb{R}^{r \times n}, \mathbf{H} \geq \mathbf{0} \}$$

is convex in  $\mathbf{H}$ .

In fact, this is quadratic programming.

Even better: decouples into  $n$  problems of size  $r$ , for  $k = 1 : n$

$$\arg \min_{\mathbf{H}(:,k)} \{ \|\mathbf{A}(:,k) - \mathbf{W}\mathbf{H}(:,k)\|^2 : \mathbf{H}(:,k) \geq \mathbf{0} \}.$$

ANLS: alternately optimize for  $\mathbf{H}$  as above, then fix  $\mathbf{H}$  and optimize for  $\mathbf{W}$ .

Optimizing for  $\mathbf{W}$  yields  $m$  quadratic programming problems each of  $r$  variables.

This could be slow.

Another faster algorithm in practice: **Hierarchical Alternating Least Squares (HALS)**.

Consider row  $\ell$  of  $\mathbf{H}$ .

Fix all other entries of  $\mathbf{H}$ .

Optimal choice for row:

$$\arg \min_{\mathbf{H}(\ell,:) \geq \mathbf{0}} \left\| \mathbf{A} - \sum_{k \neq \ell} \mathbf{W}(:,k) \mathbf{H}(k,:) - \mathbf{W}(:,\ell) \mathbf{H}(\ell,:) \right\|_{\text{F}}^2.$$

Simpler notation:  $\mathbf{h} := \mathbf{H}(\ell,:)^{\top}$ ,  $\mathbf{w} := \mathbf{W}(:,\ell)$ ,  $\mathbf{B}^{\ell} = \mathbf{A} - \sum_{k \neq \ell} \mathbf{W}(:,k) \mathbf{H}(k,:)$ ,

$$\arg \min_{\mathbf{h} \geq \mathbf{0}} \left\{ \|\mathbf{B}^{\ell} - \underbrace{\mathbf{w}\mathbf{h}^{\top}}_{\mathbb{R}^{m \times n}}\|_{\text{F}}^2 \right\}.$$

The problem separates over  $j$ :

$$\arg \min_{h_j \geq 0} \left\| \mathbf{B}^{\ell}(:,j) - \mathbf{w}h_j \right\|^2 = \arg \min_{h_j \geq 0} \left\{ \mathbf{w}^{\top} \mathbf{w} h_j^2 - 2\mathbf{w}^{\top} \mathbf{B}^{\ell}(:,j) h_j + \|\mathbf{B}^{\ell}(:,j)\|^2 \right\}.$$

This is a parabola:

$$h_j = \max \left\{ 0, \frac{\mathbf{w}^{\top} \mathbf{B}^{\ell}(:,j)}{\mathbf{w}^{\top} \mathbf{w}} \right\}.$$

So HALS applies this formula to update every entry of  $\mathbf{H}$ , then every entry of  $\mathbf{W}$ , then back to  $\mathbf{H}$ , etc.

Both HALS and ANLS have a descent property:  $\|\mathbf{A} - \mathbf{W}\mathbf{H}\|_{\text{F}}$  decreases on every iteration. To make HALS efficient, must compute  $\mathbf{B}^{\ell}$  via updating previous values.

## 23.2 Gradient Descent and Related Methods

### 23.2.1 GD

Suppose  $GD$  applied to  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $L$ -smooth, not necessarily convex.

Assume  $f$  bounded below by  $f_{\min}$ .

Then after  $k$  iterations of  $GD$  with stepsize  $\frac{1}{L}$ ,

$$\min_{0 \leq j \leq k-1} \|\nabla f(\mathbf{x}^j)\| \leq \sqrt{\frac{2L(f(\mathbf{x}^0) - f_{\min})}{k}}$$

also,

$$\lim_{k \rightarrow \infty} \nabla f(\mathbf{x}^k) = \mathbf{0}.$$

What if  $L$  not known?

Can use **Backtrack Line-Search** (Problem Set 3):

$$\mathbf{x}^{k+1} := \mathbf{x}^k - 2^{\ell_k} \nabla f(\mathbf{x}^k)$$

where  $\ell_k \in \mathbb{Z}$  determined adaptively (Armijo's Condition).

Assuming  $f$  is  $L$ -smooth, previous theorem holds with worse constants for Backtrack Line-Search.

Flaw with Backtrack Line-Search: each iteration is more expensive than using step-size  $\frac{1}{L}$ , because of the need to evaluate  $f(\mathbf{x}^k - 2^\ell \nabla f(\mathbf{x}^k))$  for multiple values of  $\ell$  for each  $k$ .

Published Line-Search Algorithms are more efficient than *Backtrack Line-Search*.

In convex case, can improve on  $GD$  via  $AGD$ . But  $AGD$  is not usable in nonconvex cases.

### 23.2.2 CGD

**Nonlinear Conjugate Gradient Descent (CGD)** often outperforms  $GD$ .

(Linear) Conjugate Gradient (Hestenes and Stiefel 1952) solves

$$\min \underbrace{\frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{b}^\top \mathbf{x}}_{f(\mathbf{x})}, \quad \mathbf{A} \in \mathbb{S}^n, \mathbf{A} \succ 0$$

which is same as solving  $\mathbf{A} \mathbf{x} = \mathbf{b}$ .

Given  $\mathbf{x}_0 \in \mathbb{R}^n$  arbitrarily,

(1)  $\mathbf{g}_0 := \mathbf{A} \mathbf{x}_0 - \mathbf{b}$  ( =  $\nabla f(\mathbf{x}_0)$  )

(2)  $\mathbf{p}_0 := -\mathbf{g}_0$

for  $k = 0, 1, 2, \dots$

(3)  $\alpha_k := \frac{-\mathbf{g}_k^\top \mathbf{p}_k}{\mathbf{p}_k^\top \mathbf{A} \mathbf{p}_k}$

(4)  $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$

(5)  $\mathbf{g}_{k+1} := \mathbf{g}_k + \alpha_k \mathbf{A} \mathbf{p}_k$  ( =  $\mathbf{A} \mathbf{x}_{k+1} - \mathbf{b} = \nabla f(\mathbf{x}_{k+1})$  )

$$(6) \beta_{k+1} := \frac{\mathbf{g}_{k+1}^\top \mathbf{g}_{k+1}}{\mathbf{g}_k^\top \mathbf{g}_k}$$
$$(7) \mathbf{p}_{k+1} := -\mathbf{g}_{k+1} + \beta_{k+1} \mathbf{p}_k.$$

Step (3) is called **Exact Line-Search**

$$\alpha_k := \arg \min_{\alpha} f(\mathbf{x}_k + \alpha \mathbf{p}_k).$$

Exact Line-Search is generally not used in optimization because of little benefit for expensive computation (use Backtrack Line-Search instead), however in current settings it's OK.

Fact: in exact arithmetic,  $CG$  returns the exact minimizer of  $f$  after  $\leq n$  iterations.

Thanks to success of  $CG$ , there are many attempts to extend to an arbitrary  $f$ .

**(Fletcher-Reeves 1960s.)**

- replace  $\mathbf{g}_k$  with  $\nabla f(\mathbf{x}_k)$
- use inexact line-search for  $\alpha_k$
- many formulas for  $\beta_k$
- occasional restart: set  $\beta_{k+1} := 0$

Say  $\mu$  is the modulus of strong convexity,  $L$  is the modulus of smoothness.

(Nonlinear)  $CGD$

Quadratics

$$\left(1 - \sqrt{\frac{\mu}{L}}\right)^k \text{ (better constant)}$$

No guarantees for nonquadratics

No hyperparameters

(Convex Functions)  $AGD$

Convex functions

$$\left(1 - \sqrt{\frac{\mu}{L}}\right)^k$$

Guaranteed for all convex functions

Needs  $\mu, L$

## LEC 24

## 24.1 Motivation

A software package to embed high-dimensional data in a lower dimension.

Example usage: if embed into 2 dimensions, then obtain a visual representation of data.

Premise: realistic high dimensional data actually lies on a lower dimensional manifold (a priori unknown). Goal is to identify the manifold.

If low dimensional manifold is a linear subspace, then can be found via SVD. If it is an affine subspace, then translate data by its mean and apply SVD. Both are called **Principal Components Analysis (PCA)**.

If manifold is curved, need something like UMAP.

## 24.2 UMAP Algorithm

Given  $n$  vectors in  $\mathbb{R}^d$ ,

- (i) Compute  $k$ -nearest-neighbor graph on data. For each data item, find its closest neighbor, second-closest neighbor, etc up to  $k$ . This yields a directed graph with distances on arcs.
- (ii) Assign weights to edges based on distances.
- (iii) Symmetrize the graph (yielding an undirected graph).
- (iv) Use **Spectral Method** to produce an embedding.
- (v) *SGD* to improve spectral embedding.

## Step 4 Spectral Embedding

Let  $\mathbf{p}_i \in \mathbb{R}^d, \forall i = 1 : n$ . The **Euclidean Distance Matrix (EDM)** is an  $n \times n$  symmetric matrix whose  $(i, j)$  entry is  $\|\mathbf{p}_i - \mathbf{p}_j\|^2$ .

The **Gram Matrix** is also  $n \times n$  symmetric,  $(i, j)$  entry is  $\mathbf{p}_i^\top \mathbf{p}_j$ . Given Gram Matrix, it's straightforward to obtain EDM:

$$\|\mathbf{p}_i - \mathbf{p}_j\|^2 = \mathbf{p}_i^\top \mathbf{p}_i - 2\mathbf{p}_i^\top \mathbf{p}_j + \mathbf{p}_j^\top \mathbf{p}_j.$$

To convert an EDM to a Gram Matrix is more complicated. Must assume:  $\sum_{i=1}^n \mathbf{p}_i = \mathbf{0}$ .

Add up one row of EDM, say row  $i$ :

$$\sum_{j=1}^n \|\mathbf{p}_i - \mathbf{p}_j\|^2 = n\mathbf{p}_i^\top \mathbf{p}_i + \sum_{j=1}^n \mathbf{p}_j^\top \mathbf{p}_j.$$

Sum this over  $i$ , letting  $\rho = \sum_{j=1}^n \mathbf{p}_j^\top \mathbf{p}_j$ , obtain  $2n\rho$ .

So from sum of EDM entries, we obtain  $\rho$ . Once  $\rho$  is known, we obtain  $\mathbf{p}_i^\top \mathbf{p}_i$  for all  $i$ , then obtain  $\mathbf{p}_i^\top \mathbf{p}_j$  for all  $i, j$  by

$$\|\mathbf{p}_i - \mathbf{p}_j\|^2 = \mathbf{p}_i^\top \mathbf{p}_i - 2\mathbf{p}_i^\top \mathbf{p}_j + \mathbf{p}_j^\top \mathbf{p}_j.$$

Thus all entries of Gram Matrix are recovered.

From Gram Matrix, we can almost recover  $\mathbf{p}_i$ 's.

Let

$$\mathbf{P} = \begin{bmatrix} \mathbf{p}_1^\top \\ \vdots \\ \mathbf{p}_n^\top \end{bmatrix} \in \mathbb{R}^{n \times d}, \quad \mathbf{G} = \mathbf{P} \mathbf{P}^\top.$$

Observe:  $\text{rank}(\mathbf{G}) \leq d$ .

Let  $\mathbf{Q} \mathbf{D} \mathbf{Q}^\top$  be eigendecomposition of  $\mathbf{G}$ . At most  $d$  entries of  $\mathbf{D}$  are nonzero, say entries  $1:d$ .

$$\mathbf{G} = \mathbf{Q}(:, 1:d) \mathbf{D}(1:d, 1:d) \mathbf{Q}(:, 1:d)^\top = \bar{\mathbf{P}} \bar{\mathbf{P}}^\top,$$

where  $\bar{\mathbf{P}} := \mathbf{Q}(:, 1:d) \mathbf{D}(1:d, 1:d)^{1/2}$ .

Fact: if  $\mathbf{G} = \mathbf{P} \mathbf{P}^\top = \bar{\mathbf{P}} \bar{\mathbf{P}}^\top$  for  $\mathbf{P}, \bar{\mathbf{P}} \in \mathbb{R}^{n \times d}$ , then exists orthogonal  $\mathbf{Q} \in \mathbb{R}^{d \times d}$  such that  $\bar{\mathbf{P}} = \mathbf{P} \mathbf{Q}$ .

So given EDM  $\mathbf{E}$ , recover  $\mathbf{P}$  via

$$\mathbf{E} \xrightarrow{\sum \mathbf{p}_i = \mathbf{0}} \mathbf{G} \xrightarrow{\text{Eigendecomposition}} \mathbf{P} \quad (\text{up to unknown rotation}).$$

Given a data set with knowledge of pairwise distances:

- (i) Square distances.
- (ii) Form *Fake EDM*.
- (iii) Form Gram Matrix.
- (iv) Form Eigendecomposition.
- (v) Keep top  $d$  eigenvalues ( $d$  is the desired embedding dimension).

Eigenvectors yield coordinates in  $\mathbb{R}^d$ .

## Step 5 SGD

Say  $p, q$  are probability distributions over same set of events. **Cross-Entropy** of  $q$  with respect to  $p$  is  $-\mathbb{E}_p[\log q]$ .

Example: flipping a coin with true heads probability  $p$  but guessed value  $q$ .

$$\text{CE}(p, q) = -p \log q - (1 - p) \log(1 - q).$$

This is convex in  $q$ , minimized at  $p = q$ .

UMAP idea: assume a hidden undirected graph among data points with edge probabilities from Step 3. For embedded points  $\mathbf{x}_i \in \mathbb{R}^{\bar{d}}$ ,

$$\Pr[(i, j) \text{ edge}] = \frac{1}{1 + a\|\mathbf{x}_i - \mathbf{x}_j\|^{2b}},$$

with hyperparameters  $a, b > 0$ .

UMAP minimizes Cross-Entropy:

$$\min_{\mathbf{x}} - \sum_{e=(i,j)} \mathbb{E}_{p_e}[\log q_e].$$

UMAP uses a variant of *SGD*: on each iteration pick edge  $(i, j)$  with probability proportional to true weight and update  $\mathbf{x}_i, \mathbf{x}_j$  using gradient of

$$-p_{ij} \log \left( \frac{1}{1 + a\|\mathbf{x}_i - \mathbf{x}_j\|^{2b}} \right) - (1 - p_{ij}) \log \left( 1 - \frac{1}{1 + a\|\mathbf{x}_i - \mathbf{x}_j\|^{2b}} \right).$$

Then perform **Negative Sampling**. Choose random  $i, j$  with weight 0 and update similarly.

Negative sampling originates from **Word2Vec** by Mikolov et al.

*Does this really minimize Cross-Entropy?*

*Unsure..(similar but not same)*

A pattern in *SGD* for modern machine learning: start with a principled approach, then insert many heuristics (layer normalizations, etc), then keep trying until it works well.