

# Lecture Notes of Linear Algebra and Convex Optimization (Dr. Jun GJ)

## Linear Algebra

### 1. PSD matrix

A positive semi-definite (PSD) matrix is a symmetric matrix  $A \in \mathbb{R}^{n \times n}$  for which the quadratic form  $x^T A x$  is non-negative for all vectors  $x$ . This means that for any vector  $x$ , the expression  $x^T A x \geq 0$ .

Key properties of PSD:

Symmetry: a PSD  $A$  is always symmetric, i.e.,  $A = A^T$ .

Non-negative Eigenvalues: All eigenvalues of a PSD matrix are non-negative.

Eigenvalue Decomposition:  $A = Q \Lambda Q^T$ . ( $Q$  is an orthogonal matrix of eigenvectors;  $\Lambda$  is a diagonal matrix of eigenvalues).

Square Root Construction: The principal square root  $A^{1/2}$  can be constructed using the eigenvalue decomposition  $A^{1/2} = Q \Lambda^{1/2} Q^T$  ( $\Lambda^{1/2}$  is the diagonal matrix obtained by taking the square roots of the non-negative eigenvalues of  $A$ ).

Lemma 1.1: a PSD matrix  $A$  can be expressed as  $A = A^{1/2} A^{1/2}$  ( $A^{1/2}$  is known as the principal square root of  $A$ .)

Proposition 1.2: a PSD matrix  $A$  has a multiset of eigenvalues  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  and orthonormal eigenvectors  $\{a_1, a_2, \dots, a_n\}$ , then

$$A = \sum_{i=1}^n \lambda_i a_i a_i^T$$
 ( $P_i = a_i a_i^T$  are orthogonal projectors and the eigenvectors  $\{a_i\}$  form an orthonormal basis of the vector space.)

### 2. SVD

Singular Value Decomposition (SVD) decomposes a matrix  $A \in \mathbb{R}^{n \times m}$  into three other matrices, capturing essential geometric and algebraic properties of  $A$ . The SVD of an  $m \times n$  matrix  $A$  is given by:

$$A = U \Sigma V^T$$

$U$ :  $n \times n$  orthogonal matrix. The columns of  $U$  are called the left singular vectors of  $A$ .

$\Sigma$ :  $s \times s$  diagonal matrix with non-negative real numbers on the diagonal. These values are called the singular values of  $A$ , typically ordered in a descending order. If  $A$  is rank  $r$ , then  $\Sigma$  contains  $r$  positive singular values, with remaining entries being zero.

$V^T$ : the transpose of an  $m \times s$  orthogonal matrix  $V$ . The columns of  $V$  (or rows of  $V^T$ ) are called the right singular vectors of  $A$ .

Practical Applications: Dimensionality Reduction (low-rank approximation) and Data Compression.

### 3. Matrix Multiplication

The product of two matrices  $A \in \mathbb{R}^{n \times m}$  and  $B \in \mathbb{R}^{m \times p}$  is the matrix.

$$C = AB \in \mathbb{R}^{n \times p}$$

where  $C_{ij} = \sum_{k=1}^m A_{ik} B_{kj}$ . The number of columns in  $A$  must equal the number of rows in  $B$ .

Vector-Vector Products (inner product / dot product)

$$x^T y \in \mathbb{R} = [x_1, x_2, \dots, x_n] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \sum_{i=1}^n x_i y_i$$

Outer Product of vectors

$$xy^T \in \mathbb{R}^{n \times n} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} [y_1, y_2, \dots, y_n] = \begin{bmatrix} x_1 y_1 & x_1 y_2 & \cdots & x_1 y_n \\ x_2 y_1 & x_2 y_2 & \cdots & x_2 y_n \\ \vdots & \vdots & \ddots & \vdots \\ x_n y_1 & x_n y_2 & \cdots & x_n y_n \end{bmatrix}$$

### Matrix-Vector Products

If we write A in column form, we see that

$$y = Ax = \begin{bmatrix} | & | & \dots & | \\ a_1 & a_2 & \dots & a_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} | \\ a_1 \end{bmatrix} x_1 + \begin{bmatrix} | \\ a_2 \end{bmatrix} x_2 + \dots + \begin{bmatrix} | \\ a_n \end{bmatrix} x_n.$$

In other words,  $y$  is a linear combination of columns of  $A$ , where the coefficients of the linear combination are given by the entries of  $x$ .

### Matrix-Matrix Products

$$\begin{aligned} C = AB &= \begin{bmatrix} | & | & \dots & | \\ a_1 & a_2 & \dots & a_n \end{bmatrix} \begin{bmatrix} -b_1^T \\ -b_2^T \\ \vdots \\ -b_p^T \end{bmatrix} = \sum_{i=1}^n a_i b_i^T \quad (\text{represent } A \text{ by columns and } B \text{ by rows}) \\ &= A \begin{bmatrix} | & | & \dots & | \\ b_1 & b_2 & \dots & b_p \end{bmatrix} = \begin{bmatrix} | & | & \dots & | \\ Ab_1 & Ab_2 & \dots & Ab_p \end{bmatrix} \quad (\text{represent } B \text{ by columns}) \\ &= \begin{bmatrix} -a_1^T \\ -a_2^T \\ \vdots \\ -a_m^T \end{bmatrix} B = \begin{bmatrix} -a_1^T B \\ -a_2^T B \\ \vdots \\ -a_m^T B \end{bmatrix} \quad (\text{represent } A \text{ by rows}) \end{aligned}$$

Exercise: When  $a_i^T B$  represents the sum of elements in the  $i$ -row, what's  $B$ ?

### 4. trace

The trace of a square matrix  $A \in \mathbb{R}^{n \times n}$ .  $\text{Tr}(A) = \sum_{i=1}^n A_{ii}$ .

If  $a_i^T B = 1$ , what's it implied?

key properties:

Ans:  $QB = I \in \mathbb{R}^m$

- For  $A \in \mathbb{R}^{n \times n}$ .  $\text{Tr}(A) = \text{Tr}(A^T)$ .

② It implies the one vector is the eigenvector of  $A$ , and the eigenvalue is 1.

- For  $A, B \in \mathbb{R}^{n \times n}$ ,  $\text{Tr}(A+B) = \text{Tr}(A) + \text{Tr}(B)$ .

- For  $A \in \mathbb{R}^{n \times n}$ ,  $t \in \mathbb{R}$ .  $\text{Tr}(tA) = t\text{Tr}(A)$ .

- For  $A, B$  such that  $AB$  is square.  $\text{Tr}(AB) = \text{Tr}(BA)$ .

- For  $A, B, C$  such that  $ABC$  is square,  $\text{Tr}(ABC) = \text{Tr}(BCA) = \text{Tr}(CAB)$ .

### 5. Norms

(optional)

Definition 5.1: A mapping  $\Phi: \mathbb{H} \rightarrow \mathbb{R}_+$  is said to define a norm on  $\mathbb{H}$  if it verifies the following axioms:

- definiteness:  $\forall x \in \mathbb{H}, \Phi(x) = 0 \Leftrightarrow x = 0$ ;
- homogeneity:  $\forall x \in \mathbb{H}, \forall \lambda \in \mathbb{R}, \Phi(\lambda x) = |\lambda| \Phi(x)$ ;
- triangle inequality:  $\forall x, y \in \mathbb{H}, \Phi(x+y) \leq \Phi(x) + \Phi(y)$ .

A norm is typically denoted by  $\|\cdot\|$ .

More generally, for any  $p \geq 1$ , the  $L_p$ -norm is defined on  $\mathbb{R}^n$  as:

$$\forall x \in \mathbb{R}^n, \|x\|_p = \left( \sum_{j=1}^n |x_j|^p \right)^{\frac{1}{p}}.$$

The  $L_1$ ,  $L_2$ , and  $L_\infty$  norms are some of the most commonly used norms, where  $\|x\|_\infty = \max_{j \in [n]} |x_j|$ . The following general inequalities relating these norms can be proven straightforwardly:

## 6. Dual Norms (optional)

**Definition 6.1:** Let  $\|\cdot\|$  be a norm on  $\mathbb{R}^n$ , then, the dual norm  $\|\cdot\|_*$  associated to  $\|\cdot\|$  is the norm defined by

$$\forall y \in \mathbb{R}^n, \|y\|_* = \sup_{\|x\|=1} |\langle y, x \rangle|.$$

For any  $p, q \geq 1$  that are conjugate, that  $\frac{1}{p} + \frac{1}{q} = 1$ , the  $L_p$  and  $L_q$  norms are dual norms of each other. For example, the dual norm of  $L_2$  is the  $L_2$ -norm, and the dual norm of the  $L_1$ -norm is the  $L_\infty$ -norm.

**Proposition 6.2 (Holder's inequality)** Let  $p, q \geq 1$  be conjugate:  $\frac{1}{p} + \frac{1}{q} = 1$ . Then, for all  $x, y \in \mathbb{R}^n$ ,

$$|\langle x, y \rangle| \leq \|x\|_p \|y\|_q.$$

with equality when  $|y_i| = |x_i|^{\frac{1}{p}}$  for all  $i \in [n]$ .

## 7. Matrix Norms (optional)

the matrix norm induced by the vector norm  $\|\cdot\|_p$  or the operator norm induced by that norm is also denoted by  $\|\cdot\|_p$  and denoted by:  $\|M\|_p = \sup_{\|x\|=1} \|Mx\|_p$ .

the norm induced for  $p=2$  is known as the spectral norm, which equals the largest singular value of  $M$ .

$$\|M\|_2 = \sigma_1(M) = \sqrt{\lambda_{\max}(M^T M)}.$$

The Frobenius norm denoted by  $\|\cdot\|_F$  is the most notable of such norms and is defined by:

$$\|M\|_F = \left( \sum_{i=1}^n \sum_{j=1}^m M_{ij}^2 \right)^{\frac{1}{2}}$$

Frobenius product:  $M, N \in \mathbb{R}^{n \times m}$ ,  $\langle M \cdot N \rangle_F = \text{Tr}(M^T N)$

this relates the Frobenius norm to the singular values of  $M$ :

$$\|M\|_F^2 = \text{Tr}(M^T M) = \sum_{i=1}^r \sigma_i^2(M), \text{ where } r = \text{rank}(M)$$

## 8. Linear Independence

A set of vectors  $\{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^m$  is said to be (linearly) independent if no vector can be represented as a linear combination of the remaining vectors. Conversely, if one vector belonging to the set can be represented as a linear combination of the remaining vectors, then the vectors are said to be (linearly) dependent. That is, if

$$x_n = \sum_{i=1}^{n-1} \alpha_i x_i,$$

for some scalar values  $\alpha_1, \alpha_2, \dots, \alpha_{n-1} \in \mathbb{R}$ , then we say that the vectors  $x_1, x_2, \dots, x_n$  are linearly dependent.

**Exercise:** Given  $Ax_1 = C_1$ ,  $Ax_2 = C_2$ ,  $Ax_3 = C_3$ ,  $C_1, C_2$  and  $C_3$  are linearly independent in  $\mathbb{R}^3$ , then what is  $Ax_4 = ?$

**Ans:**  $C_4 = Ax_4 = \alpha_1 C_1 + \alpha_2 C_2 + \alpha_3 C_3 = A(\alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3)$ , where  $\alpha_1, \alpha_2, \alpha_3 \in \mathbb{R}$ .

**Exercise:** Could a randomized always generate a set of linear independent vectors?

**Ans:** When vectors are initialized randomly, especially over a sufficiently large field of range, the likelihood of picking vectors that are linearly dependent is low because the conditions for linear dependence are very specific.

## 9. The Inverse

The inverse of a square matrix  $A \in \mathbb{R}^{n \times n}$  is denoted  $A^{-1}$ , and is the unique matrix such that

$$A^{-1}A = I = AA^{-1}$$

Note that not all matrices have inverses. Non-square matrices, for example, do not have inverse by definition. For some square matrices  $A$ , it may still be the case that  $A^{-1}$  may not exist. In particular, we say that  $A$  is invertible or non-singular if  $A^{-1}$  exists and non-invertible or singular otherwise.

In order for a square matrix  $A$  to have an inverse  $A^{-1}$ , then  $A$  must be full rank.

In particular, Let's assume  $A^{-1} = [c_1 \ c_2 \ \dots \ c_n]$ , then  $AA^{-1} = [A_{11} \ A_{12} \ \dots \ A_{1n}] = [e_1 \ e_2 \ \dots \ e_n]$ .

## 10. Orthogonal matrices

Two vectors  $x, y \in \mathbb{R}^n$  are orthogonal if  $x^T y = 0$ . A vector  $x \in \mathbb{R}^n$  is normalized if  $\|x\|_2 = 1$ . A square matrix  $U \in \mathbb{R}^{n \times n}$  is orthogonal if all its columns are orthogonal to each other and are normalized. (the columns are referred to as being orthonormal.)

$$U^T U = I = U U^T$$

key property: Operating on a vector with an orthogonal matrix does not change its Euclidean norm, i.e.,

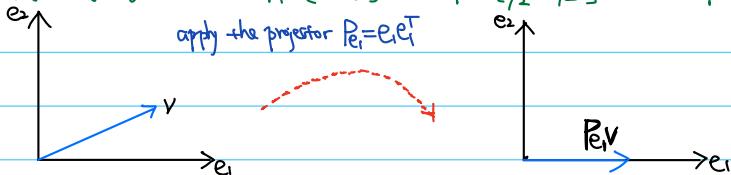
$$\|Ux\|_2 = \|x\|_2, \forall x \in \mathbb{R}^n, U \in \mathbb{R}^{n \times n} \text{ orthogonal.}$$

## 11. Projection Matrix

A matrix  $P$  is a projection matrix if  $P^2 = P$ . Given a set of orthonormal vectors  $e_1, e_2, \dots, e_n$ , a projection matrix onto the subspace  $W = \text{span}(e_1, e_2, \dots, e_n)$  is given by

$$P_W = \sum_{i=1}^n e_i e_i^T$$

Exercise: Show that the matrices  $P_1 = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$  and  $P_2 = \begin{bmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{bmatrix}$  are both projection matrices.



## 12. Eigenvalues and Eigenvectors

Given a square matrix  $A \in \mathbb{R}^{n \times n}$ , we say that  $\lambda \in \mathbb{C}$  is an eigenvalue of  $A$  and  $x \in \mathbb{C}^n$  is the corresponding eigenvector if

$$Ax = \lambda x, x \neq 0.$$

We can rewrite the equation above to state that  $(\lambda, x)$  is an eigenvalue-eigenvector pair of  $A$ , if.

$$(\lambda I - A)x = 0, x \neq 0.$$

But  $(\lambda I - A)x = 0$  has a non-zero solution to  $x$  if and only if  $(\lambda I - A)$  has a non-empty nullspace, which is only the case if  $(\lambda I - A)$  is singular, i.e.,

$$|\lambda I - A| = 0.$$

Key properties:

- The trace of  $A$  is equal to the sum of its eigenvalues.  $\text{Tr}(A) = \sum_{i=1}^n \lambda_i$
- The determinant of  $A$  is equal to the product of its eigenvalues,  $|A| = \prod_{i=1}^n \lambda_i$
- The rank of  $A$  is equal to the number of non-zero eigenvalues of  $A$ .
- If  $A$  is non-singular, then  $\lambda_i$  is an eigenvalue of  $A^{-1}$  with associated eigenvector  $x_i$ , i.e.,  $A^{-1}x_i = (\lambda_i)x_i$ .
- We can write all the eigenvector equations simultaneously as  $Ax = \lambda x$ .

Q: Given  $A$ , can you compute the eigenvalues  $\lambda_i$ ?

## Convex Optimization

### B. The Gradient

Suppose that  $f: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$  is a function that takes as input a matrix  $A$  of size  $n \times m$  and return a real value. Then, the gradient of  $f$  (with respect to  $A \in \mathbb{R}^{n \times m}$ ) is the matrix of partial derivatives, defined as:

$$\nabla_A f(A) \in \mathbb{R}^{n \times m} = \begin{bmatrix} \frac{\partial f(A)}{\partial A_{11}} & \frac{\partial f(A)}{\partial A_{12}} & \cdots & \frac{\partial f(A)}{\partial A_{1m}} \\ \frac{\partial f(A)}{\partial A_{21}} & \frac{\partial f(A)}{\partial A_{22}} & \cdots & \frac{\partial f(A)}{\partial A_{2m}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(A)}{\partial A_{n1}} & \frac{\partial f(A)}{\partial A_{n2}} & \cdots & \frac{\partial f(A)}{\partial A_{nm}} \end{bmatrix}$$

Note that the size of  $\nabla_A f(A)$  is always the same as the size of  $A$ . If  $A$  is a vector  $x \in \mathbb{R}^n$ ,

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix} \quad \text{the 1st-order gradient of } f \text{ is a } n\text{-dimensional vector.}$$

### C. Hessian matrix

The Hessian matrix wrt.  $x$ , written  $\nabla_x^2 f(x)$  or simply as  $H$  is the  $n \times n$  matrix of partial derivatives,

$$\nabla_x^2 f(x) \in \mathbb{R}^{n \times n} = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix}$$

The 2nd-order gradient of  $f$  is an  $n \times n$  matrix

Exercise: Calculate the Hessian matrix of  $f(x_1, x_2, \dots, x_n) = \sum_{i=1}^n e^{x_i}$  and verify the obtained Hessian is positive semi-definite.

Ans:  $\nabla f = [e^{x_1}, e^{x_2}, \dots, e^{x_n}]^T$ .  $\nabla^2 f = \text{diag}(e^{x_1}, e^{x_2}, \dots, e^{x_n}) \in \mathbb{R}^{n \times n}$  (Hessian Matrix)

thus,  $\forall z \in \mathbb{R}^n$ ,  $z^T \nabla^2 f z = \sum_{i=1}^n e^{x_i} z_i^2 \geq 0 \Rightarrow \nabla^2 f$  is PSD  $\Rightarrow f$  is a convex function.

### 15 Convex set

Definition 15.1: A set  $X \subseteq \mathbb{R}^n$  is said to be convex if for any two points  $x, y \in X$ , the segment  $[x, y]$  lies in  $X$ , that is

$$\{ \lambda x + (1-\lambda)y : 0 \leq \lambda \leq 1 \} \subseteq X$$

Definition 15.2 (Convex hull): The convex hull  $\text{conv}(X)$  of a set of points  $X \subseteq \mathbb{R}^n$  is the minimal convex set containing  $X$  and can be equivalently defined as follows:

$$\text{conv}(X) = \left\{ \sum_{i=1}^n \alpha_i x_i : n \geq 1, \forall i \in [n], x_i \in X, \alpha_i \geq 0, \sum_{i=1}^n \alpha_i = 1 \right\}$$

## 16 Convex function

**Definition 16.1.** Let  $X$  be a convex set. A function  $f: X \rightarrow \mathbb{R}$  is said to be convex iff  $\text{Ep}(f)$  is a convex set, or equivalently, if for all  $x, y \in X$  and  $\alpha \in [0, 1]$ ,

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

$f$  is said to be concave when  $-f$  is convex.

**Theorem 16.2.** Let  $f$  be a differentiable function, then  $f$  is convex if and only if  $\text{dom}(f)$  is convex and the following inequalities hold:

$$\forall x, y \in \text{dom}(f), \quad f(y) - f(x) \geq \nabla f(x)^T(y-x).$$

**Theorem 16.3.** Let  $f$  be a twice differentiable function, then  $f$  is convex iff  $\text{dom}(f)$  is convex and its Hessian is positive semi-definite:  $\forall x \in \text{dom}(f), \quad \nabla^2 f(x) \succeq 0$ .

**Example (Quadratic function).** The function  $f(x) = x^2$  defined over  $\mathbb{R}$  is convex since it is twice differentiable and for all  $x \in \mathbb{R}$ ,  $f''(x) = 2 > 0$ .

**Example (norms).** Any norm  $\|\cdot\|$  defined over a convex set  $X$  is convex since by the triangle inequality and the homogeneity property of the norm, for all  $\alpha \in [0, 1]$ ,  $x, y \in X$ , we can write

$$\| \alpha x + (1-\alpha)y \| \leq \| \alpha x \| + \| (1-\alpha)y \| = \alpha \| x \| + (1-\alpha) \| y \|.$$

**Proposition 16.4.** For an unconstrained optimization problem  $\min_{x \in \mathbb{R}^n} f(x)$ , the vector  $x^*$  minimizes  $f$  over  $\mathbb{R}^n$  if and only if  $\nabla f(x^*) = 0$ .

**7 Constrained optimization** **Definition 17.1.** Let  $X \subseteq \mathbb{R}^n$  and  $f, g_i: X \rightarrow \mathbb{R}$ , for all  $i \in [n]$ , then, a constrained optimization problem has the form:

$$\begin{aligned} & \min_{x \in X} f(x) \\ & \text{s.t. } g_i(x) \leq 0, \quad \forall i \in [n]. \end{aligned}$$

**Definition 17.2.** The Lagrangian function or Lagrangian associated to the general constrained optimization problem is the function defined over  $X \times \mathbb{R}^m$  by:

$$\forall x \in X, \quad \forall \lambda \geq 0, \quad L(x, \lambda) = f(x) + \sum_{i=1}^m \lambda_i g_i(x)$$

where the variables  $\lambda_i$  are known as the Lagrangian or dual variables with  $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_m)^T$ .

**Definition 17.3.** The (Lagrangian) dual function associated to the constrained optimization problem is defined by

$$\forall \lambda \geq 0, \quad F(\lambda) = \inf_{x \in X} L(x, \lambda) = \inf_{x \in X} (f(x) + \sum_{i=1}^m \lambda_i g_i(x))$$

Note that  $F$  is always concave (the Lagrangian is linear w.r.t.  $\lambda$  and the infimum preserves concavity).

We further observe that  $\forall \lambda \geq 0, \quad F(\lambda) \leq f^*(\lambda^*)$ ,  $\lambda^* \in X$ . ( $\forall x \in X, \quad f(x) + \sum_{i=1}^m \lambda_i g_i(x) \leq f(x^*)$ ).

Definition 7.4 (Dual Problem): The dual optimization problem associated to the constrained optimization problem is

$$\max_{\lambda} F(\lambda)$$

$$\text{s.t.: } \lambda \geq 0.$$

The dual problem is always a convex optimization problem (as a maximization of a concave problem). Let  $d^*$  denote its optimal value. There is  $d^* \leq p^*$  (weak duality).

The difference  $(p^* - d^*)$  is called the duality gap. The equality case

$$d^* = p^* \quad (\text{strong duality}).$$

does not hold in general. However, strong duality holds when the constrained convex optimization problem is conceived.

8. Conjugate function: Definition 8.1. For a given function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , its conjugate function  $f^*: \mathbb{R}^n \rightarrow \mathbb{R}$  is defined as:

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

key property: the Conjugate function is always a convex function, regardless of whether the original function is convex.

Proof: the supremum of affine functions ( $\langle y, x \rangle - f(x)$ ) is often in  $y$  when  $x$  is fixed) is always convex.