

EECS127 Course Notes

Anmol Parande

Spring 2021 - Professor Laurent El Ghaoui

Disclaimer: These notes reflect 127 when I took the course (Spring 2021). They may not accurately reflect current course content, so use at your own risk. If you find any typos, errors, etc, please raise an issue on the [GitHub repository](#).

Contents

1	Linear Algebra	3
1.1	Norms	3
1.2	Inner Products	4
1.3	Functions	4
1.3.1	Types of Functions	5
1.3.2	Vector Calculus	6
1.4	Matrices	6
1.4.1	Symmetric Matrices	7
1.4.2	QR Factorization	7
1.4.3	Singular Value Decomposition	8
2	Fundamentals of Optimization	9
2.1	Problem Transformations	10
2.2	Robust Optimization	10
3	Linear Algebraic Optimization	11
3.1	Projection	11

3.1.1	Matrix Pseudo-inverses	11
3.2	Explained Variance	12
3.2.1	PCA	12
3.3	Removing Constraints	13
4	Convex Optimization	13
4.1	Convexity	13
4.2	Optimality	15
4.3	Conic Programming	15
4.3.1	Quadratic Programming	16
4.3.2	Linear Programming	17
5	Duality	18
5.1	Strong Duality	19

1 Linear Algebra

Definition 1 An affine set is one of the form $\mathcal{A} = \{\mathbf{x} \in \mathcal{X} : \mathbf{x} = \mathbf{v} + \mathbf{x}_0, \mathbf{v} \in \mathcal{V}\}$ where \mathcal{V} is a subspace of a vector space \mathcal{X} and \mathbf{x}_0 is a given point.

Notice that by definition 1, a subspace is simply an affine set containing the origin. Also notice that the dimension of an affine set \mathcal{A} is the same as the dimension of \mathcal{V} .

1.1 Norms

Definition 2 A norm on the vector space \mathcal{X} is a function $\|\cdot\| : \mathcal{X} \rightarrow \mathbb{R}$ which satisfies:

1. $\|\mathbf{x}\| \geq 0$ with equality if and only if $\mathbf{x} = \mathbf{0}$
2. $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$
3. $\|\alpha\mathbf{x}\| = |\alpha|\|\mathbf{x}\|$ for any scalar α .

Definition 3 The l_p norms are defined by

$$\|\mathbf{x}\|_p = \left(\sum_{k=1}^n |x_k|^p \right)^{\frac{1}{p}}, \quad 1 \leq p \leq \infty$$

In the limit as $p \rightarrow \infty$,

$$\|\mathbf{x}\|_\infty = \max_k |x_k|.$$

Similar to vectors, matrices can also have norms.

Definition 4 A function $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ is a matrix norm if

$$f(A) \geq 0 \quad f(A) = 0 \Leftrightarrow A = 0 \quad f(\alpha A) = |\alpha|f(A) \quad f(A + B) \leq f(A) + f(B)$$

Definition 5 The Frobenius norm is the l_2 norm applied to all elements of the matrix.

$$\|A\|_F = \sqrt{\text{trace}AA^T} = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$

One useful way to characterize matrices is by measuring their “gain” relative to some l_p norm.

Definition 6 *The operator norm is defined as*

$$\|A\|_p = \max_{\mathbf{u} \neq 0} \frac{\|A\mathbf{u}\|_p}{\|\mathbf{u}\|_p}$$

When $p = 2$, the norm is called the spectral norm because it relates to the largest eigenvalue of $A^T A$.

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$$

1.2 Inner Products

Definition 7 *An inner product on real vector space is a function that maps $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ to a non-negative scalar, is distributive, is commutative, and $\langle \mathbf{x}, \mathbf{x} \rangle = 0 \Leftrightarrow \mathbf{x} = 0$.*

Inner products induce a norm $\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$. In \mathbb{R}^n , the standard inner product is $\mathbf{x}^T \mathbf{y}$. The angle between two vectors is given by

$$\cos \theta = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2}.$$

In general, we can bound the absolute value of the standard inner product between two vectors.

Theorem 1 (Holder Inequality)

$$|\mathbf{x}^T \mathbf{y}| \leq \sum_{k=1}^n |x_k y_k| \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q, \quad p, q \geq 1 \text{ s.t. } p^{-1} + q^{-1} = 1.$$

Notice that for $p = q = 2$, theorem 1 turns into the Cauchy-Schwartz Inequality ($|\mathbf{x}^T \mathbf{y}| \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2$).

1.3 Functions

We consider functions to be of the form $f : \mathbb{R}^n \rightarrow \mathbb{R}$. By contrast, a map is of the form $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$. The components of the map f are the scalar valued functions f_i that produce each component of a map.

Definition 8 *The graph of a function f is the set of input-output pairs that f can attain.*

$$\{(x, f(x)) \in \mathbb{R}^{n+1} : x \in \mathbb{R}^n\}$$

Definition 9 The epigraph of a function is the set of input-output pairs that f can achieve and anything above.

$$\{(x, t) \in \mathbb{R}^{n+1} : \mathbf{x} \in \mathbb{R}^n, t \geq f(x)\}$$

Definition 10 The t -level set is the set of points that achieve exactly some value of f .

$$\{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) = t\}$$

Definition 11 The t -sublevel set of f is the set of points achieving at most a value t .

$$\{x \in \mathbb{R}^n : f(x) \leq t\}$$

Definition 12 The half-spaces are the regions of space which a hyper-plane separates.

$$H_- = \{x : \mathbf{a}^T \mathbf{x} \leq b\} \quad H_+ = \{x : \mathbf{a}^T \mathbf{x} > b\}$$

Definition 13 A polyhedron is the intersection of m half-spaces given by $\mathbf{a}_i^T \mathbf{x} \leq b_i$ for $i \in [1, m]$.

When a polyhedron is bounded, it is called a polytope.

1.3.1 Types of Functions

Theorem 2 A function is linear if and only if it can be expressed as $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} + b$ for some unique pair (\mathbf{a}, b) .

An affine function is linear when $b = 0$. A hyperplane is simply a level set of a linear function.

Theorem 3 Any quadratic function can be written as the sum of a quadratic term involving a symmetric matrix and an affine term:

$$q(x) = \frac{1}{2} \mathbf{x}^T H \mathbf{x} + \mathbf{c}^T \mathbf{x} + d.$$

Another special class of functions are polyhedral functions.

Definition 14 A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is polyhedral if its epigraph is a polyhedron.

$$\text{epi } f = \left\{ (x, t) \in \mathbb{R}^{n+1} : C \begin{bmatrix} \mathbf{x} \\ t \end{bmatrix} \leq d \right\}$$

1.3.2 Vector Calculus

We can also do calculus with vector functions.

Definition 15 The gradient of a function at a point x where f is differentiable is a column vector of first derivatives of f with respect to the components of x

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

The gradient is perpendicular to the level sets of f and points from a point x_0 to higher values of the function. In other words, it is the direction of steepest increase. It is akin to the derivative of a 1D function.

Definition 16 The Hessian of a function f at point x is a matrix of second derivatives.

$$H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

The Hessian is akin to the second derivative in a 1D function. Note that the Hessian is a symmetric matrix.

1.4 Matrices

Matrices define a linear map between an input space and an output space. Any linear map $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ can be represented by a matrix.

Theorem 4 (Fundamental Theorem of Linear Algebra) For any matrix $A \in \mathbb{R}^{m \times n}$,

$$\mathcal{N}(A) \oplus \mathcal{R}(A^T) = \mathbb{R}^n \quad \mathcal{R}(A) \oplus \mathcal{N}(A^T) = \mathbb{R}^m.$$

1.4.1 Symmetric Matrices

Recall that a symmetric matrix is one where $A = A^T$.

Theorem 5 (Spectral Theorem) *Any symmetric matrix is orthogonally similar to a real diagonal matrix.*

$$A = A^T \implies A = U\Lambda U^T = \sum_i \lambda_i \mathbf{u}_i \mathbf{u}_i^T, \quad \|\mathbf{u}\| = 1, \quad \mathbf{u}_i^T \mathbf{u}_j = 0 \ (i \neq j)$$

Let $\lambda_{\min}(A)$ be the smallest eigenvalue of symmetric matrix A and $\lambda_{\max}(A)$ be the largest eigenvalue.

Definition 17 *The Rayleigh Quotient for $\mathbf{x} \neq \mathbf{0}$ is $\frac{\mathbf{x}^T A \mathbf{x}}{\|\mathbf{x}\|^2}$.*

Theorem 6 *For any $\mathbf{x} \neq \mathbf{0}$,*

$$\lambda_{\min}(A) \leq \frac{\mathbf{x}^T A \mathbf{x}}{\|\mathbf{x}\|^2} \leq \lambda_{\max}(A).$$

Two special types of symmetric matrices are those with non-negative eigenvalues.

Definition 18 *A symmetric matrix is positive semi-definite if $\mathbf{x}^T A \mathbf{x} \geq 0 \implies \lambda_{\min}(A) \geq 0$.*

Definition 19 *A symmetric matrix is positive definite if $\mathbf{x}^T A \mathbf{x} > 0 \implies \lambda_{\min}(A) > 0$.*

These matrices are important because they often have very clear geometric structures. For example, an ellipsoid in multi-dimensional space can be defined as the set of points

$$\mathcal{E} = \{x \in \mathbb{R}^m : \mathbf{x}^T P^{-1} \mathbf{x} \leq 1\}$$

where P is a positive definite matrix. The eigenvectors of P give the principle axes of this ellipse, and $\sqrt{\lambda}$ are the semi-axis lengths.

1.4.2 QR Factorization

Similar to how spectral theorem allows us to decompose symmetric matrices, QR factorization is another matrix decomposition technique that works for any general matrix.

Definition 20 *The QR factorization matrix are the orthogonal matrix Q and the upper triangular matrix R such that $A = QR$*

An easy way to find the QR factorization of a matrix is to apply Gram Schmidt to the columns of the matrix and express the result in matrix form. Suppose that our matrix A is full rank (i.e its columns \mathbf{a}_i are linearly independent) and we have applied Gram-Schmidt to columns $\mathbf{a}_{i+1} \cdots \mathbf{a}_n$ to get orthogonal vectors $\mathbf{q}_{i+1} \cdots \mathbf{q}_n$. Continuing the procedure, the i th orthogonal vector \mathbf{q}_i is

$$\tilde{\mathbf{q}}_i = \mathbf{a}_i - \sum_{k=i+1}^n (\mathbf{q}_k^T \mathbf{a}_i) \mathbf{q}_k \quad \mathbf{q}_i = \frac{\tilde{\mathbf{q}}_i}{\|\tilde{\mathbf{q}}_i\|_2}$$

If we re-arrange this, to solve for \mathbf{a}_i , we see that

$$\mathbf{a}_i = \|\tilde{\mathbf{q}}_i\|_2 \mathbf{q}_i + \sum_{k=i+1}^n (\mathbf{q}_k^T \mathbf{a}_i) \mathbf{q}_k$$

Putting this in matrix form, we can see that

$$\begin{bmatrix} | & | & \cdots & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \\ | & | & \cdots & | \end{bmatrix} = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{q}_1 & \mathbf{q}_2 & \cdots & \mathbf{q}_n \\ | & | & \cdots & | \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & \cdots & r_{1n} \\ 0 & r_{22} & \cdots & r_{2n} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & r_{nn} \end{bmatrix} \quad r_{ij} = \mathbf{a}_i^T \mathbf{q}_j, r_{ii} = \|\tilde{\mathbf{q}}_i\|_2$$

1.4.3 Singular Value Decomposition

Definition 21 A matrix $A \in \mathbb{R}^{m \times n}$ is a dyad if it can be written as $\mathbf{p}\mathbf{q}^T$.

A dyad is a rank-one matrix. It turns out that all matrices can be decomposed into a sum of dyads.

Definition 22 The Singular Value Decomposition of a matrix A is

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

where σ_i are the singular values of A and \mathbf{u}_i and \mathbf{v}_i are the left and right singular vectors.

The singular values are ordered such that $\sigma_1 \geq \sigma_2 \geq \cdots$. The left singular values are the eigenvectors of AA^T and the right singular values are the eigenvectors of $A^T A$. The singular values are $\sqrt{\lambda_i}$ where λ_i are the eigenvalues of $A^T A$. Since AA^T and $A^T A$ are symmetric, \mathbf{u}_i and \mathbf{v}_i are orthogonal. The number of non-zero singular values is equal to the rank of the matrix. We can write the SVD in matrix form as

$$A = [U_r \quad U_{n-r}] \text{diag}(\sigma_1, \cdots, \sigma_r, 0, \cdots, 0) \begin{bmatrix} V_r^T \\ V_{n-r}^T \end{bmatrix}$$

Writing the SVD tells us that

1. V_{n-r} forms a basis for $\mathcal{N}(A)$
2. U_r form a basis for $\mathcal{R}(A)$

The Frobenius norm and spectral norm are tightly related to the SVD.

$$\|A\|_F = \sum_i \sigma_i^2$$

$$\|A\|_2^2 = \sigma_1^2$$

2 Fundamentals of Optimization

Definition 23 *The standard form of optimization is*

$$p^* = \min_{\mathbf{x}} f_0(\mathbf{x}) \text{ such that } f_i(\mathbf{x}) \leq 0$$

- The vector $\mathbf{x} \in \mathbb{R}^n$ is known as the **decision variable**.
- The function $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$ is the **objective**.
- The functions $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is the **constraints**.
- p^* is the **optimal value**, and the \mathbf{x}^* which achieves the optimal value is called the **optimizer**.

Definition 24 *The feasible set of an optimization problem is*

$$\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^n : f_i(\mathbf{x}) \leq 0\}$$

Definition 25 *A point \mathbf{x} is ϵ -suboptimal if it is feasible and satisfies*

$$p^* \leq f_0(\mathbf{x}) \leq p^* + \epsilon$$

Definition 26 *An optimization problem is strictly feasible if $\exists \mathbf{x}_0$ such that all constraints are strictly satisfied (i.e inequalities are strict inequalities, and equalities are satisfied).*

2.1 Problem Transformations

Sometimes, optimizations in a particular formulation do not admit themselves to be solved easily. In this case, we can sometimes transform the problem into an easier one from which we can easily recover the solution to our original problem. In many cases, we can introduce additional “slack” variable and constraints to massage the problem into a form which is easier to analyze.

Theorem 7 (Epigraphic Constraints) $\min_x f_0(x)$ is equivalent to the problem with epigraphic constraints

$$\min_{x,t} t \quad : \quad f_0(x) \leq t,$$

theorem 7 works because by minimizing t , we are also minimizing how large $f_0(x)$ can get since $f_0(x) \leq t$, so at optimum, $f_0(x) = t$. It can be helpful when $f_0(x) \leq t$ can be massaged further into constraints that are easier to deal with.

Theorem 8 (Monotone Objective Transformation) Let $\Phi : \mathbb{R} \rightarrow \mathbb{R}$ be a continuous and strictly increasing function over a feasible set \mathcal{X} . Then

$$\min_{\mathbf{x} \in \mathcal{X}} f_0(\mathbf{x}) \equiv \min_{\mathbf{x} \in \mathcal{X}} \Phi(f_0(\mathbf{x}))$$

2.2 Robust Optimization

For a “nominal” problem

$$\min_{\mathbf{x}} f_0(\mathbf{x}) \quad : \quad \forall i \in [1, m], f_i(\mathbf{x}) \leq 0,$$

uncertainty can enter in the data used to create the f_0 and f_i . It can also enter during decision time where the \mathbf{x}^* which solves the optimization cannot be implemented exactly. These uncertainties can create unstable solutions or degraded performance. To make our optimization more robust to uncertainty, we add a new variable $\mathbf{u} \in \mathcal{U}$.

Definition 27 For a nominal optimization problem $\min_{\mathbf{x}} f_0(\mathbf{x})$ subject to $f_i(\mathbf{x}) \leq 0$ for $i \in [1, m]$, the robust counterpart is

$$\min_{\mathbf{x}} \max_{\mathbf{u} \in \mathcal{U}} f_0(\mathbf{x}, \mathbf{u}) \quad : \quad \forall i \in [1, m], f_i(\mathbf{x}, \mathbf{u}) \leq 0$$

3 Linear Algebraic Optimization

Many optimization problems can be solved using the machinery of Linear Algebra. These problems do not have inequality constraints or non-euclidean norms in the objective function.

3.1 Projection

The idea behind projection is to find the closest point in a set closest (with respect to particular norm) to a given point.

Definition 28 Given a vector \mathbf{x} in inner product space \mathcal{X} and a subspace $S \subseteq \mathcal{X}$, the projection of \mathbf{x} onto S is given by

$$\Pi_S(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y} \in S} \|\mathbf{y} - \mathbf{x}\|$$

where the norm is the one induced by the inner product.

Theorem 9 There exists a unique vector $\mathbf{x}^* \in S$ which solves

$$\min_{\mathbf{y} \in S} \|\mathbf{y} - \mathbf{x}\|.$$

It is necessary and sufficient for \mathbf{x}^* to be optimal that $(\mathbf{x} - \mathbf{x}^*) \perp S$. The same condition applies when projecting onto an affine set.

3.1.1 Matrix Pseudo-inverses

Definition 29 A pseudoinverse is a matrix A^\dagger that satisfies:

$$AA^\dagger A = A \quad A^\dagger AA^\dagger = A^\dagger \quad (AA^\dagger)^T = AA^\dagger \quad (A^\dagger A)^T = A^\dagger A$$

There are several special cases of pseudoinverses.

1. $A^\dagger = V_r \operatorname{diag} \left(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_r} \right) U_r^T$ is the Moore-Penrose Pseudo-inverse.
2. When A and non-singular, $A^\dagger = A^{-1}$.
3. When A is full column rank, $A^\dagger = (A^T A)^{-1} A^T$.

4. When A is full row rank, $A^\dagger = A^T(AA^T)^{-1}$

The pseudo-inverses are useful because they can easily compute the projection of a vector onto a related subspace of A .

1. $\operatorname{argmin}_{z \in \mathcal{R}(A)} \|z - y\|_2 = AA^\dagger y$
2. $\operatorname{argmin}_{z \in \mathcal{R}(A)^\perp} \|z - y\|_2 = (I - AA^\dagger)y$
3. $\operatorname{argmin}_{z \in \mathcal{N}(A)} \|z - y\|_2 = (I - A^\dagger A)y$
4. $\operatorname{argmin}_{z \in \mathcal{N}(A)^\perp} \|z - y\|_2 = A^\dagger Ay$

3.2 Explained Variance

The Low Rank Approximation problem is to approximate a matrix A with a rank k matrix

$$\min_{A_k} \|A - A_k\|_F^2 \text{ such that } \operatorname{rank}(A_k) = k.$$

The solution to the low rank approximation problem is simply the first k terms of the SVD:

$$A_K^* = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T.$$

This is because the singular values give us a notion of how much of the Frobenius Norm (Total Variance) each dyad explains.

$$\eta = \frac{\|A_k\|_F^2}{\|A\|_F^2} = \frac{\sum_i^k \sigma_i^2}{\sum_i^r \sigma_i^2}$$

3.2.1 PCA

Suppose we had a matrix containing m data points in \mathbb{R}^n (each data point is a column), and without loss of generality, assume this data is centered around 0 (i.e. $\sum_i \mathbf{x}_i = 0$). The variance of this data along a particular direction z is given by $z^T X X^T z$. Principle Component Analysis is finding the directions z such that the variance is maximized.

$$\max_{z \in \mathbb{R}^n} z^T X X^T z \text{ such that } \|z\|_2 = 1$$

The left singular vector corresponding to the largest singular value of the $X X^T$ matrix is the optimizer of this problem, and the variance along this direction is σ_1^2 . If we wanted to find subsequent directions of maximal variance, they are just the left singular vectors corresponding to the largest singular values.

3.3 Removing Constraints

Following from the Fundamental Theorem of Linear Algebra, if $A\mathbf{x} = \mathbf{y}$ has a solution, then the set of solutions can be expressed as

$$S = \{\bar{\mathbf{x}} + N\mathbf{z}\}$$

where $A\bar{\mathbf{x}} = \mathbf{y}$ and N is a basis for $\mathcal{N}(A)$. This means if we have a constrained optimization problem

$$\min_{\mathbf{x}} f_0(\mathbf{x}) : A\mathbf{x} = \mathbf{b},$$

we can write an equivalent unconstrained problem

$$\min_{\mathbf{z}} f_0(\mathbf{x}_0 + N\mathbf{z})$$

where $A\mathbf{x}_0 = \mathbf{b}$

4 Convex Optimization

4.1 Convexity

Definition 30 A subset $C \in \mathbb{R}^n$ is convex if it contains the line segment between any two points in the set.

$$\forall \mathbf{x}_1, \mathbf{x}_2 \in C, \lambda \in [0, 1], \quad \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in C$$

Convexity can be preserved by some operations.

Theorem 10 If C_1, \dots, C_m are convex sets, then their intersection $C = \bigcap_{i=1, \dots, m} C_i$ is also a convex set.

Theorem 11 If a map $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is affine and $C \subset \mathbb{R}^n$ is convex, then $f(C) = \{f(\mathbf{x}) : \mathbf{x} \in C\}$ is convex.

theorems 10 and 11 are important because they allow us to prove sets are convex using sets that we know are convex. For example, theorem 11 tells us that a projection of a convex set onto a subspace must also be convex since projection is a linear operator.

Definition 31 A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if its domain is a convex set and $\forall \mathbf{x}, \mathbf{y}$ in the domain, $\lambda \in [0, 1]$,

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

Loosely, convexity means that the function is bowl shaped since a line connecting any two points on the function is above the function itself. A concave function is simply one where $-f$ is convex, and these appear like a “hill”. Because convex functions are bowl shaped, they must be ∞ outside their domain.

Theorem 12 *A function f is convex if and only if its epigraph is a convex set.*

Just like convex sets, some operations preserve convexity for functions.

Theorem 13 *If $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ are convex functions, then $f(\mathbf{x}) = \sum_{i=1}^m \alpha_i f_i(\mathbf{x})$ where $\alpha_i \geq 0$ is also convex.*

A similar property to theorem 11 exists for convex functions.

Theorem 14 *If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, then $g(\mathbf{x}) = f(A\mathbf{x} + b)$ is also convex.*

We can also look at the first and second order derivatives to determine the convexity of a function.

Theorem 15 *If f is differentiable, then f is convex if and only if*

$$\forall \mathbf{x}, \mathbf{y}, \quad f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla_x^T(\mathbf{y} - \mathbf{x})$$

theorem 15 can be understood geometrically by saying the graph of f is bounded below everywhere by its tangent hyperplanes.

Theorem 16 *If f is twice differentiable, then f is convex if and only if the Hessian ∇^2 is positive semi-definite everywhere.*

Geometrically, the second-order condition says that f looks bowl-shaped.

Theorem 17 *A function f is convex if and only if its restriction to any line $g(t) = f(\mathbf{x}_0 + t\mathbf{v})$ is convex.*

Theorem 18 *If $(f_\alpha)_{\alpha \in A}$ is a family of convex functions, then the pointwise maximum $f(\mathbf{x}) = \max_{\alpha \in A} f_\alpha(\mathbf{x})$ is convex.*

Because of the nice geometry that convexity gives, optimization problems which involve convex functions and sets are reliably solvable.

Definition 32 A convex optimization problem in standard form is

$$p^* = \min_{\mathbf{x}} f_0(\mathbf{x}) : \quad \forall i \in [1, m], f_i(\mathbf{x}) \leq 0, A\mathbf{x} = \mathbf{b}$$

where f_0, f_1, \dots are convex functions and the equality constraints are affine.

Since the constraints form a convex set, definition 32 is equivalent to minimizing a convex function over a convex set \mathcal{X} .

Theorem 19 A locally optimal solution to a convex problem is also globally optimal, and this set \mathcal{X} is convex.

theorem 19 is why convex problems are nice to solve.

4.2 Optimality

When problems are convex, we can define conditions that any optimal solution must satisfy.

Theorem 20 For a convex optimization problem with a differentiable objective function $f_0(\mathbf{x})$ and feasible set \mathcal{X} ,

$$\mathbf{x} \text{ is optimal} \Leftrightarrow \forall \mathbf{y} \in \mathcal{X}, \nabla_{\mathbf{x}} f_0(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \geq 0$$

Since the gradient points in the direction of greatest increase, the dot product of the gradient with the different between any vector and the optimal solution being positive means other solutions will only increase the value of $f_0(\mathbf{x})$. For unconstrained problems, we can make this condition even sharper.

Theorem 21 In a convex unconstrained problem with a differentiable objective function $f_0(\mathbf{x})$, \mathbf{x} is optimal if and only if $\nabla_{\mathbf{x}} f_0(\mathbf{x}) = \mathbf{0}$

4.3 Conic Programming

Conic programming is the set of optimization problems which deal with variables constrained to a second-order cone.

Definition 33 A n -dimensional second-order cone is the set

$$\mathcal{K}_n = \{(\mathbf{x}, t), \mathbf{x} \in \mathbb{R}^n, t \in \mathbb{R} : \|\mathbf{x}\|_2 \leq t\}$$

By Cauchy-Schwartz, $\|\mathbf{x}\|_2 = \max_{\mathbf{u}: \|\mathbf{u}\| \leq 1} \mathbf{u}^T \mathbf{x} \leq t$. This means that second order cones are convex sets since they are the intersection of half-spaces. In spaces 3-dimensions and higher, we can rotate these cones.

Definition 34 A rotated second order cone in \mathbb{R}^{n+2} is the set

$$\mathcal{K}_n^r = \{(\mathbf{x}, y, z), \mathbf{x} \in \mathbb{R}^n, y \in \mathbb{R}, z \in \mathbb{R} : \mathbf{x}^T \mathbf{x} \leq yz, y \geq 0, z \geq 0\}.$$

The rotated second-order cone can be interpreted as a rotation because the hyperbolic constraint $\|\mathbf{x}\|_2^2 \leq yz$ can be expressed equivalently as

$$\left\| \begin{bmatrix} 2\mathbf{x} \\ y - z \end{bmatrix} \right\|_2 \leq y + z.$$

Definition 35 The standard Second Order Cone Constraint is

$$\|A\mathbf{x} + \mathbf{b}\|_2 \leq \mathbf{c}^T \mathbf{x} + d.$$

A SOC constraint will confine \mathbf{x} to a second order cone since if we let $\mathbf{y} = A\mathbf{x} + \mathbf{b} \in \mathbb{R}^m$ and $t = \mathbf{c}^T \mathbf{x} + d$, then $(\mathbf{y}, t) \in \mathcal{K}_m$.

Definition 36 A second-order cone program in standard inequality form is given by

$$\min \mathbf{c}^T \mathbf{x} \text{ such that } \|A_i \mathbf{x} + \mathbf{b}_i\|_2 \leq \mathbf{c}_i^T \mathbf{x} + d_i.$$

An SOC program is a convex problem since its objective is linear, and hence convex, and the SOC constraints are also convex.

4.3.1 Quadratic Programming

A special case of SOCPs are Quadratic Programs. These programs have constraints and an objective function which can be expressed as a quadratic function. In SOCP form, they look like

$$\begin{aligned} \min_{\mathbf{x}, t} \quad & \mathbf{a}_0^T \mathbf{x} + t \\ \text{s.t:} \quad & \left\| \begin{bmatrix} 2Q_0^{\frac{1}{2}} \mathbf{x} \\ t - 1 \end{bmatrix} \right\|_2 \leq t + 1 \\ & \left\| \begin{bmatrix} 2Q_i^{\frac{1}{2}} \mathbf{x} \\ b_i - \mathbf{a}_i^T \mathbf{x} - 1 \end{bmatrix} \right\|_2 \leq b_i - \mathbf{a}_i^T \mathbf{x} + 1 \end{aligned}$$

Since they are a special case of SOCPs, Quadratic Programs are also convex.

Definition 37 *The standard form of a quadratic constrained quadratic program is*

$$\min_{\mathbf{x}} \mathbf{x}^T Q_0 \mathbf{x} + \mathbf{a}_0^T \mathbf{x} \quad : \quad \forall i \in [1, m], \mathbf{x}^T Q_i \mathbf{x} + \mathbf{a}_i^T \mathbf{x} \leq b_i$$

To be a quadratic program, the matrix H must be positive semi-definite. If the $Q_i = 0$ in the constraints, then we get a normal quadratic program.

Definition 38 *The standard form of a quadratic program is given by*

$$\min_{\mathbf{x}} \frac{1}{2} \mathbf{x}^T H \mathbf{x} + \mathbf{c}^T \mathbf{x} \quad : \quad \forall i \in [1, m], \mathbf{a}_i^T \mathbf{x} \leq b_i$$

Its SOCP form looks like

$$\begin{aligned} \min_{\mathbf{x}, y} \quad & \mathbf{c}^T \mathbf{x} + y \\ \text{s.t:} \quad & \left\| \begin{bmatrix} 2H^{\frac{1}{2}} \mathbf{x} \\ y - 1 \end{bmatrix} \right\|_2 \leq y + 1, \\ & \mathbf{a}_i^T \mathbf{x} \leq b_i \end{aligned}$$

In the special case where H is positive definite and we have no constraints, then

$$\frac{1}{2} \mathbf{x}^T H \mathbf{x} + \mathbf{c}^T \mathbf{x} + d = \frac{1}{2} (\mathbf{x} + H^{-1} \mathbf{c})^T H (\mathbf{x} + H^{-1} \mathbf{c}) + d - (H^{-1} \mathbf{c})^T H (H^{-1} \mathbf{c})$$

Thus

$$\operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \mathbf{x}^T H \mathbf{x} + \mathbf{c}^T \mathbf{x} + d = -H^{-1} \mathbf{c}$$

4.3.2 Linear Programming

If the matrix in the objective function of a quadratic program is 0 (and there are no quadratic constraints), then the resulting objective and constraints are affine functions. This is a linear program.

Definition 39 *The inequality form of a linear program is given by*

$$\min_{\mathbf{x}} \mathbf{c}^T \mathbf{x} + d \quad : \quad \forall i \in [1, m], \mathbf{a}_i^T \mathbf{x} \leq b_i$$

Since linear program is a special case of a quadratic program, it can also be expressed as an SOCP.

$$\begin{aligned} \min_{\mathbf{x}} \quad & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & \forall i \in [1, m], \|\mathbf{0}\mathbf{x} + 0\|_2 \leq b_i - \mathbf{a}_i^T \mathbf{x} \end{aligned}$$

Because of the constraints, the feasible set of a linear program is a polyhedron. Thus linear programs are also convex.

5 Duality

Definition 40 A primal optimization problem is given by

$$p^* = \min_{\mathbf{x} \in \mathbb{R}^n} f_0(\mathbf{x}) : \forall i \in [1, m] f_i(\mathbf{x}) \leq 0, \forall k \in [1, n] h_k(\mathbf{x}) = 0$$

The primal problem is essentially the standard form of optimization. There are no assumptions of convexity on any of the functions involved. We would like to express primal problems as a min-max optimization with no constraints.

Definition 41 The Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ using Lagrange Multipliers $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ is given by

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{k=1}^n \mu_k h_k(\mathbf{x})$$

The Lagrangian achieves the goal of removing the constraints in the min-max optimization

$$p^* = \min_{\mathbf{x} \in \mathbb{R}^n} \max_{\boldsymbol{\lambda} \geq 0, \boldsymbol{\mu}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$$

This is true because if any inequality constraints are violated, then $f_i(\mathbf{x}) \geq 0$, and the maximization could set λ_i very large to make the overall problem ∞ , and if any equality constraints are violated, then $h_k(\mathbf{x}) \neq 0$, and the maximization would set μ_k to a very large number of the same sign as $h_k(\mathbf{x})$ to make the overall problem ∞ . Thus the minimax problem is equivalent to the original problem. At this point, it might be easier to solve the problem if the order of min and max were switched.

Theorem 22 (Minimax Inequality) For any sets X, Y and any function $F : X \times Y \rightarrow \mathbb{R}$

$$\min_{\mathbf{x} \in X} \max_{\mathbf{y} \in Y} F(\mathbf{x}, \mathbf{y}) \geq \max_{\mathbf{y} \in Y} \min_{\mathbf{x} \in X} F(\mathbf{x}, \mathbf{y})$$

theorem 22 can be interpreted as a game where there is a minimizing player and a maximizing player. If the maximizer goes first, it will always produce a higher score than if the minimizer goes first (unless they are equal). We can now apply theorem 22 to switch the min and max in our optimization with the Lagrangian.

Theorem 23 (Weak Duality)

$$\min_{\mathbf{x} \in \mathbb{R}^n} \max_{\lambda \geq 0, \mu} \mathcal{L}(\mathbf{x}, \lambda, \mu) \geq \max_{\lambda \geq 0, \mu} \min_{\mathbf{x} \in \mathbb{R}^n} \mathcal{L}(\mathbf{x}, \lambda, \mu)$$

What weak duality does is convert our minimization problem to a maximization problem.

Definition 42 *The dual function of the primal problem is given by*

$$g(\lambda, \mu) = \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda, \mu)$$

Note that g is a concave function because it is the pointwise minimum of functions that are affine in μ and λ . A maximization of a concave function over a convex set is a convex problem, so the dual problem (minimizing g) is convex. Thus duality achieves two primary purposes.

1. It removes constraints, potentially making the problem easier to solve.
2. It can turn a non-convex problems into a convex one.

Even when there are no constraints, we can sometimes introduce constraints to leverage duality by adding slack variables that are equal to expressions in the objective.

5.1 Strong Duality

In some cases, duality gives not just a lower bound, but an exact value. When this happens, we have **Strong Duality**.

Theorem 24 (Sion's MiniMax Theorem) *Let $X \subseteq \mathbb{R}^n$ be convex, and $Y \subseteq \mathbb{R}^m$ be bounded and closed (compact). Let $F : X \times Y \rightarrow \mathbb{R}$ be a function such that $\forall y, F(\cdot, y)$ is convex and continuous, and $\forall x, F(x, \cdot)$ is concave and continuous, then*

$$\min_{\mathbf{x} \in X} \max_{\mathbf{y} \in Y} F(\mathbf{x}, \mathbf{y}) = \max_{\mathbf{y} \in Y} \min_{\mathbf{x} \in X} F(\mathbf{x}, \mathbf{y})$$

If we focus on convex problems, then we can find conditions which indicate strong duality holds.

Theorem 25 (Slater's Condition) *If a convex optimization problem is strictly feasible, then strong duality holds*

Once we find a solution to the dual problem, then the solution to the primal problem is recovered by minimizing $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ where $\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*$ are the optimal dual variables, and if no such feasible point \mathbf{x} exists, then the primal itself is infeasible. When searching for strong duality and an optimal solution $(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$, it can be useful to consider particular conditions.

Theorem 26 *For a convex primal problem which is feasible and has a feasible dual where strong duality holds, a primal dual pair $(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ is optimal if and only if the KKT conditions are satisfied.*

1. **Primal Feasibility:** \mathbf{x} satisfies $\forall i \in [1, m], f_i(\mathbf{x}) \leq 0$ and $\forall k \in [1, n], h_k(\mathbf{x}) = 0$.
2. **Dual Feasibility:** $\boldsymbol{\lambda} \geq \mathbf{0}$.
3. **Complementary Slackness:** $\forall i \in [1, m], \lambda_i f_i(\mathbf{x}) = 0$
4. **Lagrangian Stationarity:** *If the lagrangian is differentiable, then*

$$\nabla_x f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i \nabla_x f_i(\mathbf{x}) + \sum_{k=1}^n \mu_k \nabla_x h_k(\mathbf{x}) = \mathbf{0}$$

The complementary slackness requirement essentially says that if a primal constraint is slack ($f_i(\mathbf{x}) < 0$), then $\lambda_i = 0$, and if $\lambda_i > 0$, then $f_i(\mathbf{x}) = 0$.