

Jonas Hübotter
based on Rasmus Kyng's lectures at ETH Zurich

Graph Algorithms and Optimization

Contents. Convex optimization and duality. Spectral Graph Theory. Combinatorial Graph Algorithms. Electrical Flows.

Contributions are welcome at <https://github.com/jonhue/graph-algorithms-and-optimization>.

ACKNOWLEDGEMENT

The contents of this summary are based on the lecture “Advanced Graph Algorithms and Optimization” given by Rasmus Kyng at ETH Zurich in the spring of 2022. Certain parts are similar/taken from the lecture notes.

Contents

I	Preliminaries	7
1	Electrical Flows	9
1.1	The Laplacian Matrix	9
1.2	An Optimization Problem	12
1.3	Energy & Duality	13
2	Linear Algebra	15
2.1	Square Roots	17
2.2	Matrix Norms	18
2.3	Loewner Order	19
2.4	Kernel and Image	20
2.5	Matrix Functions	20
2.6	Pseudoinverses	22
2.7	Solving Linear Systems	23
3	Probability	25
3.1	Random Walks	25
3.2	Concentration	29
3.3	Martingales	30
4	Analysis	31
4.1	First-order Taylor Approximations	31
4.2	Directional Derivatives	32
4.3	Second-order Taylor Approximations	33

II Convex Optimization 37

5	Convex Geometry	39
5.1	Convex Sets & Functions	39
5.2	First-order Characterization of Convexity	40
5.3	Second-order Characterization of Convexity	41
6	Gradient Descent	43
6.1	Smoothness	44
6.2	Strong Convexity	45
6.3	Acceleration	47
7	Non-Euclidean Geometries	49
7.1	Mirror Descent	49
8	Lagrange Multipliers and Duality	51
8.1	Separating Hyperplanes	51
8.2	Lagrange Multipliers and KKT Conditions	52
8.3	Lagrangian Duality	54
8.4	Slater's Condition	56
8.5	Example: Duality of Maximum Flow and Minimum Cut	56
8.6	Fenchel Conjugates	57
9	Newton's Method	59

III Spectral Graph Theory 61

10	Introduction to Spectral Graph Theory	63
10.1	Eigenvalues of the Laplacian Matrix	63
10.2	Examples	65
11	Conductance and Expanders	67
11.1	Conductance	67
11.2	Cheeger's Inequality	69

11.3	<i>Sparsity</i>	69
12	<i>Effective Resistance</i>	71
12.1	<i>Effective Resistance as a Metric</i>	73
13	<i>Spectral Graph Sparsification</i>	75
14	<i>Solving Laplacian Linear Systems</i>	77
14.1	<i>Dealing with pseudoinverses</i>	77
14.2	<i>Computing the Cholesky Decomposition</i>	78
14.3	<i>Approximate Almost Linear-Time Solvers</i>	79
<i>IV Combinatorial Graph Algorithms</i>		81
15	<i>Algorithms for Maximum Flow</i>	83
15.1	<i>The Ford-Fulkerson Algorithm</i>	85
15.2	<i>Dinitz's Algorithm</i>	86
15.3	<i>The Push-Relabel Algorithm</i>	89
15.4	<i>Outlook</i>	89
16	<i>Link-Cut Trees</i>	91
17	<i>Finding Expanders using Maximum Flow</i>	93
17.1	<i>Graph Embedding</i>	93
18	<i>Distance Oracles</i>	95
<i>V Further Topics</i>		97
19	<i>Interior Point Methods for Maximum Flow</i>	99
A	<i>Solutions</i>	101
A.1	<i>Part I</i>	101

<i>A.2 Part II</i>	103
<i>A.3 Part III</i>	104

<i>Summary of Notation</i>	105
----------------------------	-----

<i>Bibliography</i>	109
---------------------	-----

<i>Index</i>	111
--------------	-----

PART I

Preliminaries

Electrical Flows

A classical graph problem is the flow of electrical currents through a network of resistors. Such a network $G = (V, E, r)$ can be described by a set of vertices V , set of wires (or edges) E , and resistances $r \in \mathbb{R}_{>0}^{|E|}$ of wires. We are interested in finding the *electrical flow* $\tilde{f} \in \mathbb{R}^{|E|}$ through the network, assigning to each wire the current that is transported per unit time. Alternatively, we can think of *electrical voltages* $\tilde{x} \in \mathbb{R}^{|V|}$ at the vertices, which Ohm's law relates to the electrical flow.

By *Ohm's law*, we have that for any wire $e \in E$,

$$\tilde{f}(e) = \frac{\tilde{x}(e)}{r(e)}, \quad \tilde{x}(e) = \tilde{f}(e) \cdot r(e), \quad (1.1)$$

where $\tilde{x}(\{u, v\}) \doteq \tilde{x}(u) - \tilde{x}(v)$ is the voltage difference of vertices u and v . For any flow $f \in \mathbb{R}^{|E|}$, the *net flow* of current at a vertex $u \in V$ is given as,

$$\sum_{v \sim u} f(v, u). \quad (1.2)$$

We say that a flow routes *demand* $d \in \mathbb{R}^{|V|}$ if the net flow at every vertex is $d(v)$. The fact that at vertices with zero demand, the flow is conserved¹ is also known as *Kirchhoff's current law*.

To keep track of the direction of flow on each edge, we assign an arbitrary direction to each edge (we "orient" G) and only consider non-negative flows, $f \in \mathbb{R}_{\geq 0}^{|E|}$. Clearly, for any previously feasible flow, we can assign directions in such a way that the flow remains feasible.

1.1 The Laplacian Matrix

Definition 1.1 (Adjacency matrix). The *adjacency matrix* of a graph G , $\tilde{A} \in \mathbb{R}^{|V| \times |V|}$, is defined as,²

We use $v \sim u$ to denote all v that are adjacent to u .

¹ As much current is flowing into the vertex as is flowing out of it.

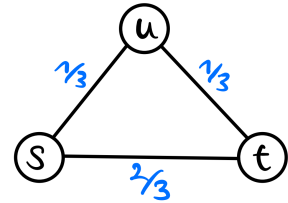


Figure 1.1: Example of an electrical flow (shown in blue) with voltages

$$x(s) = 0, \quad x(u) = 1, \quad x(t) = 2$$

and unit resistances, routing demands

$$d(s) = -1, \quad d(u) = 0, \quad d(t) = 1.$$

² By A we will later denote the weighted adjacency matrix.

$$\tilde{A}(u, v) \doteq \begin{cases} 1 & \text{if } u \sim v \\ 0 & \text{otherwise.} \end{cases} \quad (1.3)$$

Definition 1.2 (Incidence matrix). The *incidence matrix* of an oriented graph G , $B \in \mathbb{R}^{|V| \times |E|}$, is defined as,

$$B(v, e) \doteq \begin{cases} 1 & \text{if } e = (u, v) \text{ for some } u \in V \\ -1 & \text{if } e = (v, u) \text{ for some } u \in V \\ 0 & \text{otherwise.} \end{cases} \quad (1.4)$$

Each column of B only has two non-zero entries and sums to one.

Lemma 1.3. $BB^\top = \text{diag}\{\deg(v)\}_{v \in V} - \tilde{A}$.

Proof. The dot product of the rows, corresponding to the same vertex v , produces exactly $\deg(v)$. All other dot products between rows corresponding to vertices u and v are -1 iff $u \sim v$ and 0 otherwise. \square

We can now also write the net flow constraint,

$$Bf = d. \quad (1.5)$$

We define $R \doteq \text{diag}\{r(e)\}_{e \in E}$ and then have that Ohm's law can be expressed as,

$$B^\top \tilde{x} = R\tilde{f}, \quad \text{or equivalently, } R^{-1}B^\top \tilde{x} = \tilde{f}. \quad (1.6)$$

If the net flow constraint is satisfied, this yields,

$$\underbrace{BR^{-1}B^\top}_{\text{Laplacian}} \tilde{x} = B\tilde{f} = d. \quad (1.7)$$

Definition 1.4 (Laplacian matrix). The *Laplacian matrix* of an oriented graph G is defined as,

$$L \doteq BR^{-1}B^\top = BWB^\top \in \mathbb{R}^{|V| \times |V|}, \quad (1.8)$$

where $W \doteq R^{-1}$ is a diagonal matrix of weights $w(e) \doteq \frac{1}{r(e)}$.

Intuitively, the weight of an edge can be understood as how “connected” the two vertices at its endpoints are. In contrast, the resistance of an edge is smaller when endpoints are well-connected.

We will now learn a little more about Laplacian matrices.

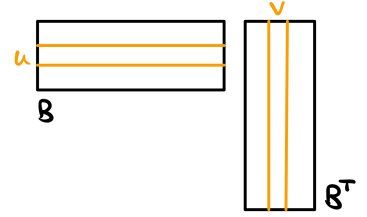


Figure 1.2: Illustration of the matrix product BB^\top .

Definition 1.5 (Weighted adjacency matrix). The *weighted adjacency matrix* of a graph G , $A \in \mathbb{R}^{|V| \times |V|}$, is defined as,

$$A(u, v) \doteq \begin{cases} w(\{u, v\}) & \text{if } u \sim v \\ 0 & \text{otherwise.} \end{cases} \quad (1.9)$$

Lemma 1.6. A is symmetric, that is, $A = A^\top$.

Proof. This follows immediately from the fact that G is undirected. \square

Definition 1.7 (Weighted degree). The *weighted degree* of a vertex $v \in V$ is given as,

$$d(v) \doteq \sum_{\{u, v\} \in E} w(\{u, v\}). \quad (1.10)$$

We write $D \doteq \text{diag}\{d(v)\}_{v \in V}$.

Lemma 1.8. $L = D - A$.

Proof. The proof is identical to the proof of lemma 1.3, only that every entry is now weighted, due to the additional factor W . \square

Corollary 1.9. L is symmetric.

Proof. This directly follows from the fact that D and A are symmetric.³ \square

³ Diagonal matrices like D are trivially symmetric.

Lemma 1.10. For any $x \in \mathbb{R}^{|V|}$, we have,

$$x^\top Lx = \sum_{\{u, v\} \in E} w(\{u, v\})[x(u) - x(v)]^2 \geq 0. \quad (1.11)$$

Proof. We have,

$$\begin{aligned} x^\top Lx &= x^\top Dx - x^\top Ax. \\ x^\top Dx &= \sum_{v \in V} d(v)x(v)^2 = \sum_{\{u, v\} \in E} w(\{u, v\})[x(u)^2 + x(v)^2]. \\ x^\top Ax &= \sum_{v \in V} x(v)(Ax)(v) \\ &= \sum_{v \in V} x(v) \sum_{u \in V} A(v, u)x(u) \\ &= \sum_{v, u \in V} w(\{u, v\})x(u)x(v) \\ &= 2 \sum_{\{u, v\} \in E} w(\{u, v\})x(u)x(v). \end{aligned}$$

Combining the above equalities, we obtain,

$$\begin{aligned} \mathbf{x}^\top \mathbf{L} \mathbf{x} &= \sum_{\{u,v\} \in E} w(\{u,v\}) [x(u)^2 + x(v)^2] - 2 \sum_{\{u,v\} \in E} w(\{u,v\}) x(u)x(v) \\ &= \sum_{\{u,v\} \in E} w(\{u,v\}) [x(u) - x(v)]^2. \end{aligned} \quad \square$$

Corollary 1.11. \mathbf{L} is positive semi-definite.⁴

Exercise 1.12. A matrix \mathbf{M} is a Laplacian matrix iff it satisfies the following conditions:

1. $\mathbf{M}^\top = \mathbf{M}$;
2. the diagonal entries of \mathbf{M} are non-negative, and the off-diagonal entries of \mathbf{M} are non-positive; and
3. $\mathbf{M}\mathbf{1} = \mathbf{0}$.

It is often useful to look at a normalized Laplacian matrix, where weighted vertex degrees are normalized to one and edges are weighted based on the degrees of their endpoints.

Definition 1.13 (Normalized Laplacian matrix). The *normalized Laplacian matrix* of an oriented graph G is defined as,

$$N(i, j) \doteq \begin{cases} 1 & \text{if } i = j \\ -\frac{1}{\sqrt{d(i)d(j)}} & \text{if } i \sim j \\ 0 & \text{otherwise.} \end{cases} \quad (1.12)$$

This characterization is equivalent to,⁵

$$\mathbf{N} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2} = \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}. \quad (1.13)$$

As we will see later, the normalized Laplacian matrix is intimately related to the probability transition matrix of a random walk on G , where transition probabilities are proportional to edge weights.

Lemma 1.14. \mathbf{N} is positive semi-definite.

Proof. We have for any $\mathbf{x} \in \mathbb{R}^n$,

$$\mathbf{x}^\top \mathbf{N} \mathbf{x} = \mathbf{x}^\top \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2} \mathbf{x} = (\mathbf{D}^{-1/2} \mathbf{x})^\top \mathbf{L} (\mathbf{D}^{-1/2} \mathbf{x}) \geq 0,$$

using positive semi-definiteness of \mathbf{L} . \square

1.2 An Optimization Problem

We have seen that finding electrical voltages $\tilde{\mathbf{x}}$ or the electrical flow $\tilde{\mathbf{f}}$ is equivalent, we can go from one to the other and back. So let us first focus on how we can find electrical voltages $\tilde{\mathbf{x}}$.

⁴ By the previous lemma, \mathbf{L} satisfies the definition of positive semi-definiteness, which we will introduce in the following section.

⁵ $\mathbf{A}^{-1/2}$, where \mathbf{A} is a diagonal matrix, is the diagonal matrix $\text{diag}_i \{A(i, i)^{-1/2}\}$.

In eq. (1.7), we saw that electrical voltages satisfy $L\tilde{x} = d$, obeying by Ohm's law and satisfying the net flow constraint. A standard approach to reframe the solution to such a system of linear equations as the result of an optimization problem, is to consider the cost function,

$$c(x) \doteq \frac{1}{2}x^\top Lx - x^\top d. \quad (1.14)$$

Observe that $\nabla c(x) = Lx - d \stackrel{!}{=} 0$ iff $Lx = d$.

Claim 1.15. c is convex, hence, its critical point coincides with its minimizer.⁶

⁶ We will develop tools to show this in the next part.

We have therefore recast the problem of finding electrical voltages to the convex optimization,

$$\tilde{x} = \arg \min_{x \in \mathbb{R}^{|V|}} c(x). \quad (1.15)$$

1.3 Energy & Duality

Let us now look at the same problem through a different lens. Transporting current through a network of resistors requires energy, which is dissipated as heat by the resistor. By *Joule's law*, sending a current \tilde{f} across a resistor with potential drop \tilde{x} , spends $\tilde{f} \cdot \tilde{x}$ units of energy per unit time.⁷ Using Ohm's law, we have,

$$\tilde{f} \cdot \tilde{x} = \frac{\tilde{x}^2}{r} = r \cdot \tilde{f}^2. \quad (1.16)$$

⁷ We will think about everything as if happening in one unit of time.

We can therefore write the *electrical energy* dissipated by routing a current \tilde{f} (or equivalently with electrical voltages \tilde{x}) as,

$$\mathcal{E}(\tilde{f}) \doteq \sum_{e \in E} r(e) \tilde{f}(e)^2 = \tilde{f}^\top R \tilde{f} = \tilde{x}^\top L \tilde{x} \doteq \mathcal{E}(\tilde{x}). \quad (1.17)$$

using Ohm's law, $\tilde{f} = R^{-1}B^\top \tilde{x}$

Remark 1.16. Given electrical voltages $\tilde{x} \perp \mathbf{1}$, we can also write the electrical energy as

$$\mathcal{E}(\tilde{x}) = \tilde{x}^\top L \tilde{x} = d^\top \tilde{x} = d^\top L^+ d \doteq \mathcal{E}(d), \quad (1.18)$$

using $Lx = d$ and $L^+ Lx = x$ as $x \perp \ker L$

using the pseudoinverse L^+ of L .

Let us consider the *electrical energy-minimizing flow*:⁸

$$f^* \doteq \arg \min_{\substack{f \in \mathbb{R}^{|E|} \\ Bf = d}} \mathcal{E}(f). \quad (1.19)$$

⁸ Note that we could instead (and equivalently) characterize the optimization problem using electrical voltages.

Exercise 1.17. f^* is precisely the electrical flow \tilde{f} , that is, f^* satisfies Ohm's law.

This indicates that the two above optimization problems are intimately related: both yield the electrical flow (or equivalently, electrical voltages). In fact, it can be shown that,

Exercise 1.18. $\mathcal{E}(\tilde{f}) = -c(\tilde{x})$,

where we think about maximizing $-c(\tilde{x})$ instead of minimizing $c(\tilde{x})$. More generally,

Exercise 1.19. $\mathcal{E}(f) \geq -c(x)$ for any flow f routing d and any voltages x .

So, for any voltages x , the value of $-c(x)$ is a lower bound on the minimum electrical energy $\mathcal{E}(\tilde{f})$.

This is an example of a much broader phenomenon known as Lagrangian duality, where we have a minimization problem and a related maximization problem that gives lower bounds on the optimal value of the minimization problem.

Linear Algebra

Claim 2.1. *If a square matrix $A \in \mathbb{R}^{n \times n}$ is symmetric¹, then A has n real eigenvalues $\lambda_1, \dots, \lambda_n$ and eigenvectors $v_1, \dots, v_n \in \mathbb{R}^n$ such that $Av_i = \lambda_i v_i$ and the v_i are orthogonal².*

¹ A is symmetric iff $A = A^\top$.

² that is, $v_i^\top v_j = 0$ for $i \neq j$

Definition 2.2 (Positive (semi-)definiteness). Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix. We say A is

1. *positive definite* iff $x^\top A x > 0$ for any $x \in \mathbb{R}^n \setminus \{0\}$;
2. *positive semi-definite* iff $x^\top A x \geq 0$ for any $x \in \mathbb{R}^n$;
3. if neither A nor $-A$ is positive semi-definite, A is *indefinite*.

We denote by S^n the set of symmetric $n \times n$ matrices, by S_+^n the set of such matrices that are positive semi-definite, and by S_{++}^n the set of such matrices that are positive definite.

Theorem 2.3. *Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix. Then,*

1. *A is positive definite iff all eigenvalues are positive; and*
2. *A is positive semi-definite iff all eigenvalues are non-negative.*

This theorem is a corollary of the Courant-Fischer theorem, which we will work towards now.

Fact 2.4 (Spectral theorem for symmetric matrices). *For all symmetric matrices $A \in \mathbb{R}^{n \times n}$ there exist*

$$V = \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad \text{and} \quad \Lambda = \text{diag}\{\lambda_i\}_{i \in [n]} \in \mathbb{R}^{n \times n}, \quad (2.1)$$

where λ_i and v_i are the eigenvalues and corresponding (normalized) eigenvectors of A , such that

1. $A = V \Lambda V^\top = \sum_{i=1}^n \lambda_i v_i v_i^\top$; and
2. $V^\top V = I$, i.e., the columns of V form an orthonormal basis of \mathbb{R}^n .

Theorem 2.5 (Courant-Fischer min-max theorem). *For symmetric*

matrices $A \in \mathbb{R}^{n \times n}$ with eigenvalues $\lambda_1 \leq \dots \leq \lambda_n$,

$$\lambda_i = \min_{\substack{\text{subspace } W \subseteq \mathbb{R}^n \\ \dim(W)=i}} \max_{\substack{x \in W \\ x \neq 0}} \frac{x^\top A x}{x^\top x} \quad (2.2)$$

$$= \max_{\substack{\text{subspace } W \subseteq \mathbb{R}^n \\ \dim(W)=n-i+1}} \min_{\substack{x \in W \\ x \neq 0}} \frac{x^\top A x}{x^\top x}. \quad (2.3)$$

Proof. We show eq. (2.2). The proof of the other equation proceeds analogously.

- “ \geq ”: We choose $W = \text{span}\{v_1, \dots, v_i\}$. We can write x in the basis of eigenvectors,

$$x = \sum_{j=1}^i c(j) v_j$$

for some $c \in \mathbb{R}^i$. We have,

$$x^\top x = \|x\|_2^2 = \sum_{j=1}^i \sum_{k=1}^i c(j) c(k) v_j^\top v_k = \sum_{j=1}^i c(j)^2,$$

using that $v_j^\top v_k = 0$ if $j \neq k$ and $v_j^\top v_k = 1$ otherwise

and,

$$\begin{aligned} x^\top A x &= x^\top V \Lambda V^\top x = (V^\top x)^\top \underbrace{\Lambda}_{\mathbf{c}} (V^\top x) \\ &= \mathbf{c}^\top \Lambda \mathbf{c} = \sum_{j=1}^i \lambda_j c(j)^2 \leq \lambda_i \sum_{j=1}^i c(j)^2. \end{aligned}$$

Altogether,

$$\frac{x^\top A x}{x^\top x} \leq \lambda_i.$$

- “ \leq ”: Consider any subspace $W \subseteq \mathbb{R}^n$ with $\dim(W) = i$ and fix the subspace $T = \text{span}\{v_i, \dots, v_n\}$ with $\dim(T) = n - i + 1$. We have that $\dim(W \cap T) = \dim(W) + \dim(T) - \dim(W \cup T)$ and $\dim(W \cup T) \leq \dim(\mathbb{R}^n) = n$, so, $\dim(W \cap T) \geq 1$. Therefore,

$$\begin{aligned} \max_{\substack{x \in W \\ x \neq 0}} \frac{x^\top A x}{x^\top x} &\geq \max_{\substack{x \in W \cap T \\ x \neq 0}} \frac{x^\top A x}{x^\top x} \\ &\geq \min_{\substack{\text{subspace } V \subseteq T \\ \dim(V)=1}} \max_{\substack{x \in V \\ x \neq 0}} \frac{x^\top A x}{x^\top x}. \end{aligned}$$

For the last inequality note that V can be chosen as $W \cap T$.

We choose $V = \text{span}\{v_i\}$. For some $c \in \mathbb{R}$, we can write $x = c v_i$.

Similarly to the previous part, we obtain,

$$\frac{x^\top A x}{x^\top x} = \frac{\lambda_i c^2}{c^2} = \lambda_i. \quad \square$$

Proof of theorem 2.3. Using Courant-Fischer, we have for the smallest eigenvalue λ_1 of A ,

$$\lambda_1 = \min_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{x^\top A x}{x^\top x}.$$

Thus, if λ_1 is positive, then $x^\top A x > 0$ for all $x \in \mathbb{R}^n \setminus \{0\}$. In contrast, if for any such x , $x^\top A x > 0$, then λ_1 must be positive. The proof of positive semi-definiteness is analogous. \square

Corollary 2.6. For a symmetric matrix $A \in \mathbb{R}^{n \times n}$ and any $x \in \mathbb{R}^n$,

$$\lambda_{\min}(A) \|x\|_2^2 \leq x^\top A x, \quad (2.4)$$

where $\lambda_{\min}(A)$ is the smallest eigenvalue of A .

Proof. By Courant-Fischer, we have for any $x \in \mathbb{R}^n$ such that $x \neq 0$,

$$\lambda_{\min}(A) = \min_{\substack{y \in \mathbb{R}^n \\ y \neq 0}} \frac{y^\top A y}{\|y\|_2^2} \leq \frac{x^\top A x}{\|x\|_2^2}.$$

If $x = 0$, the inequality trivially holds. \square

Claim 2.7. For any matrix M and invertible matrix T , M and TMT^{-1} have the same eigenvalues.

2.1 Square Roots

We now state a useful fact for positive semi-definite matrices.

Lemma 2.8. Any symmetric and positive semi-definite matrix A has a positive semi-definite square root $A^{1/2}$ such that $A^{1/2} A^{1/2} = A$.

Proof. By the spectral theorem, $A = V\Lambda V^\top$, where V is an orthonormal matrix of eigenvectors and Λ is a diagonal matrix of eigenvalues. Let $A^{1/2} \doteq V\Lambda^{1/2}V^\top$, where $\Lambda^{1/2} = \text{diag}_i\{\Lambda(i, i)^{1/2}\}$. Then,

$$\begin{aligned} A^{1/2} A^{1/2} &= V\Lambda^{1/2}V^\top V\Lambda^{1/2}V^\top \\ &= V\Lambda^{1/2}\Lambda^{1/2}V^\top \\ &= V\Lambda V^\top = A. \end{aligned}$$

It is immediately clear from the definition that $A^{1/2}$ is positive semi-definite. \square

2.2 Matrix Norms

Definition 2.9 (Matrix norm). Given a matrix $A \in \mathbb{R}^{n \times n}$ and norms $\|\cdot\|_\alpha$ and $\|\cdot\|_\beta$ on \mathbb{R}^n , the (induced) norm of A is defined as,³

$$\|A\|_{\alpha \rightarrow \beta} \doteq \sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{\|Ax\|_\beta}{\|x\|_\alpha}. \quad (2.5)$$

³ If you think of A as a linear map, you can think of $\|\cdot\|_\alpha$ as a norm of the input space and $\|\cdot\|_\beta$ as a norm of the output space.

We write $\|A\|_\alpha \doteq \|A\|_{\alpha \rightarrow \alpha}$.

Lemma 2.10. For any matrix $A \in \mathbb{R}^{n \times n}$, any $x \in \mathbb{R}^n$, and any norm $\|\cdot\|$ on \mathbb{R}^n ,

$$|x^\top Ax| \leq \|A\| \|x\|^2. \quad (2.6)$$

Proof. We have,

$$\begin{aligned} |x^\top Ax| &\leq \|x\| \|Ax\| \\ &\leq \|A\| \|x\|^2. \end{aligned}$$

using Cauchy-Schwarz

□ using the definition of the induced matrix norm

Lemma 2.11. For a symmetric matrix $A \in \mathbb{R}^{n \times n}$,

$$\|A\|_2 = \max\{|\lambda_{\min}(A)|, |\lambda_{\max}(A)|\}. \quad (2.7)$$

$\|A\|_2$ is called the spectral norm of A .

Proof. We have,

$$\begin{aligned} \|A\|_2^2 &= \sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{x^\top A^\top Ax}{x^\top x} \\ &= \sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{x^\top A^2 x}{x^\top x} \\ &= \sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{x^\top V \Lambda^2 V^\top x}{x^\top x} \\ &= \sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{V^\top x^\top \Lambda^2 (V^\top x)}{V^\top x^\top (V^\top x)} \\ &= \sup_{\substack{y \in \mathbb{R}^n \\ y \neq 0}} \frac{y^\top \Lambda^2 y}{y^\top y} = \|\Lambda\|_2^2. \end{aligned}$$

using that A is symmetric, $A^\top = A$

using that V is orthogonal, $V^\top = V^{-1}$

using that the columns of V form a basis of \mathbb{R}^n , set $y \doteq V^\top x$

Finally,

$$\|\Lambda\|_2^2 = \sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{x^\top \Lambda^2 x}{x^\top x} = \sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{\sum_{i=1}^n \lambda_i^2 x(i)^2}{\sum_{i=1}^n x(i)^2} = \max_{i \in [n]} \lambda_i^2. \quad \square$$

2.3 Loewner Order

The Loewner order (or positive semi-definite order) is a partial ordering on symmetric matrices.

Definition 2.12 (Loewner order). Given symmetric matrices $A, B \in \mathbb{R}^{n \times n}$, $A \preceq B$ iff $\forall x \in \mathbb{R}^n : x^\top A x \leq x^\top B x$.

Remark 2.13. $A \succeq 0$ iff A is positive semi-definite.

Lemma 2.14 (Properties of the Loewner order). We have that for any symmetric $A, B, C \in \mathbb{R}^{n \times n}$,

1. $A \preceq A$ (reflexivity);
2. $A \preceq B, B \preceq A \implies A = B$ (antisymmetry);
3. $A \preceq B, B \preceq C \implies A \preceq C$ (transitivity);
4. $A \preceq B \implies A + C \preceq B + C$;
5. if $A \succeq 0$, then $\frac{1}{\alpha} A \preceq A \preceq \alpha A$ for any $\alpha \geq 1$; and
6. if $A \preceq B$, then $\forall i \in [n] : \lambda_i(A) \leq \lambda_i(B)$, where $\lambda_i(M)$ is the i -th largest eigenvalue of M .⁴

Remark 2.15. Properties (1), (2), and (3) together imply that the Loewner order is a partial ordering of symmetric matrices.

⁴ The converse is false:

$$A \doteq \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}, \quad B \doteq \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$

have equal eigenvalues, but $A \not\preceq B$ and $B \not\preceq A$.

Proof. Properties (1) through (5) follow directly from the definition using only elementary operations. For property (6), we have by Courant-Fischer,

$$\begin{aligned} \lambda_i(A) &= \min_{\substack{\text{subspace } W \subseteq \mathbb{R}^n \\ \dim(W)=i}} \max_{\substack{x \in W \\ x \neq 0}} \frac{x^\top A x}{x^\top x} \\ &\leq \min_{\substack{\text{subspace } W \subseteq \mathbb{R}^n \\ \dim(W)=i}} \max_{\substack{x \in W \\ x \neq 0}} \frac{x^\top B x}{x^\top x} = \lambda_i(B). \quad \square \end{aligned}$$

Loewner Order on Graphs

Definition 2.16. We write $G \preceq H$ iff we have for the Laplacian matrices L_G and L_H that $L_G \preceq L_H$. For any $c > 0$, we write cG in place of cL_G , which corresponds to scaling the weight of every edge of G by c .

Lemma 2.17. For subgraphs $H \subseteq G$, we have $H \preceq G$.

Proof. Dropping edges can only decrease the quadratic form of the Laplacian of G (see eq. (1.11)). □

2.4 Kernel and Image

Definition 2.18 (Kernel). The *kernel* (or *null space*) of a linear map $A \in \mathbb{R}^{n \times m}$ is the linear subspace,

$$\ker A \doteq \{x \in \mathbb{R}^m \mid Ax = 0\} \subseteq \mathbb{R}^m. \quad (2.8)$$

Definition 2.19 (Image). The *image* (or *range*) of a linear map $A \in \mathbb{R}^{n \times m}$ is the linear subspace,

$$\operatorname{im} A \doteq \{Ax \mid x \in \mathbb{R}^m\} = \operatorname{span}\{A(:,1), \dots, A(:,m)\} \subseteq \mathbb{R}^n, \quad (2.9)$$

where $A(:,i)$ denotes the i -th column vector of A .

We have the following useful property relating image and kernel.

Lemma 2.20. For a matrix $A \in \mathbb{R}^{n \times m}$, we have $\operatorname{im}(A^\top)^\perp = \ker A$ and $(\operatorname{im} A)^\perp = \ker(A^\top)$.⁵

⁵ Note that $(W^\perp)^\perp = W$ for any subspace W .

Proof. If $x \in \operatorname{im}(A^\top)^\perp$, we have $A(i,:)^\top x = 0$ for all $i \in [n]$. Hence, $Ax = 0$ and $x \in \ker A$.

Conversely, if $x \in \ker A$, then $Ax = 0$. Therefore, $A(i,:)^\top x = 0$ for all $i \in [n]$. As $\operatorname{im}(A^\top) = \operatorname{span}\{A(1,:), \dots, A(n,:)\}$, we have $z^\top x = 0$ for any $z \in \operatorname{im}(A^\top)$, yielding, $x \in \operatorname{im}(A^\top)^\perp$.

The same argument goes through for the transpose of A on both sides. \square

2.5 Matrix Functions

First, let us remind ourselves of the notion of the trace of a matrix.

Definition 2.21 (Trace). The *trace* of a symmetric matrix $A \in \mathbb{R}^{n \times n}$ is defined as,

$$\operatorname{tr} A \doteq \sum_{i=1}^n A(i,i). \quad (2.10)$$

Lemma 2.22. The trace is invariant under cyclic permutations. That is, for any $A, B \in \mathbb{R}^{n \times n}$,

$$\operatorname{tr} AB = \operatorname{tr} BA. \quad (2.11)$$

Proof. We have,

$$\operatorname{tr} AB = \sum_{i=1}^n (AB)(i,i)$$

$$\begin{aligned}
 &= \sum_{i=1}^n \sum_{j=1}^n A(i, j) B(j, i) \\
 &= \sum_{j=1}^n \sum_{i=1}^n B(j, i) A(i, j) \\
 &= \sum_{j=1}^n (BA)(j, j) \\
 &= \operatorname{tr} BA. \quad \square
 \end{aligned}$$

Lemma 2.23. For $A \in \mathbb{R}^{n \times n}$,

$$\operatorname{tr} A = \sum_{i=1}^n \lambda_i, \quad (2.12)$$

where λ_i are the eigenvalues of A .

Proof. By the spectral theorem for symmetric matrices and using the cycle property of the trace,

$$\operatorname{tr} A = \operatorname{tr} (V \Lambda V^\top) = \operatorname{tr} (\underbrace{\Lambda V^\top V}_I) = \operatorname{tr} \Lambda = \sum_{i=1}^n \lambda_i,$$

where Λ is a diagonal matrix of eigenvalues of A and V is an orthonormal matrix of the corresponding eigenvectors. \square

Definition 2.24 (Matrix function). A *matrix function* $f : S^n \rightarrow S^n$ given the scalar function $f : \mathbb{R} \rightarrow \mathbb{R}$ is defined as,

$$f(A) \doteq V \operatorname{diag}_i \{f(\lambda_i)\} V^\top, \quad (2.13)$$

where $A = V \operatorname{diag}_i \{\lambda_i\} V^\top$ is the spectral decomposition of A . We say,

1. a function $f : \mathcal{S} \rightarrow \mathbb{R}$ for $\mathcal{S} \subseteq S^n$ is *monotonically increasing* iff $A \preceq B$ implies $f(A) \leq f(B)$; and similarly
2. a matrix function $f : \mathcal{S} \rightarrow \mathcal{T}$ for $\mathcal{S}, \mathcal{T} \subseteq S^n$ is *monotonically increasing* iff $A \preceq B$ implies $f(A) \preceq f(B)$.

Lemma 2.25. If $f : T \rightarrow \mathbb{R}$ for $T \subseteq \mathbb{R}$ is monotone, then $X \mapsto \operatorname{tr} f(X)$ is monotone.

Proof. Let $A \preceq B$. Then, $\lambda_i(A) \leq \lambda_i(B)$ for all $i \in [n]$. Thus, $f(\lambda_i(A)) \leq f(\lambda_i(B))$, and hence,

$$\operatorname{tr} f(A) = \sum_{i=1}^n f(\lambda_i(A)) \leq \sum_{i=1}^n f(\lambda_i(B)) = \operatorname{tr} f(B).$$

The analogous argument works when f is monotonically decreasing. \square

Claim 2.26 (Facts about matrix functions).

1. $X \mapsto X^{-1}$ is monotonically decreasing on S_{++}^n ;
2. \log is monotonically increasing on S_{++}^n ;
3. $\exp(A) \preceq I + A + A^2$ for $\|A\|_2 \leq 1$;
4. $\log(I + A) \preceq A$ for $A \succeq -I$; and
5. (Lieb's theorem) if $f(A) = \text{tr}(\exp(H + \log(A)))$ for $A \in S_{++}^n$ and some $H \in S^n$, then $-f$ is convex.

2.6 Pseudoinverses

We often want to solve systems of linear equations such as

$$Ax = b \quad (2.14)$$

where A and b are given and we seek to identify x . If A is invertible, a solution to eq. (2.14) is given by $x \doteq A^{-1}b$. Recall that a square matrix A is invertible iff $\det A \neq 0$. Can we still explicitly write x when A is not invertible?

When A is not invertible, the corresponding system of linear equations does not have a unique solution. However, we may still be able to find *some* solution. The Moore-Penrose inverse is such a generalization of the inverse.

Definition 2.27 (Moore-Penrose inverse). Given a symmetric matrix $A \in \mathbb{R}^{n \times n}$,⁶ its *Moore-Penrose inverse* (or simply *pseudoinverse*) is a matrix $A^+ \in \mathbb{R}^{n \times n}$ such that

1. A^+ is symmetric;
2. $\ker A^+ = \ker A$;⁷ and
3. for $v \perp \ker A$, $A^+Av = v$ and $AA^+v = v$.

Thus, provided $b \perp \ker A$, $x \doteq A^+b$ is a solution to eq. (2.14).

Lemma 2.28. The pseudoinverse of the symmetric matrix $A \in \mathbb{R}^{n \times n}$ is (uniquely⁸) characterized as,

$$A^+ = V\Lambda^+V^\top, \quad (2.15)$$

where V is the matrix of orthogonal eigenvectors and

$$\Lambda^+ = \text{diag}_{i \in [n]} \begin{cases} \lambda_i^{-1} & \lambda_i \neq 0 \\ 0 & \text{otherwise} \end{cases} \quad (2.16)$$

for eigenvalues λ_i .

Proof. The properties of a pseudoinverse follow from the definition immediately. \square

Corollary 2.29. $A \succeq B$ implies $A^+ \preceq B^+$ when $\ker A = \ker B$.

⁶ The Moore-Penrose inverse can also be defined for non-symmetric and even non-square matrices, but this will not be important for us.

⁷ that is, for any $v \in \mathbb{R}^n$,

$$A^+v = 0 \iff Av = 0$$

⁸ this we will not prove

2.7 Solving Linear Systems

Linear systems such as eq. (2.14) can often be solved much faster than by computing the pseudoinverse of A .

We first need to remind ourselves of the notion of a projection.

Definition 2.30 (Projection matrix). A *projection matrix* is a matrix $\Pi \in \mathbb{R}^{n \times n}$ such that $\Pi^2 = \Pi$.⁹ We say that Π is *orthogonal* iff $\text{im } \Pi \perp \ker \Pi$.

⁹ This property is called *idempotency*.

Remark 2.31. Note that Π is the identity operator on $\text{im } \Pi$, i.e., for all $x \in \text{im } \Pi$ we have that $\Pi x = x$.

Definition 2.32 (Orthogonal projection to the complement of the kernel). Given a matrix A , Π_A is the *orthogonal projection to $(\ker A)^\perp$* , that is,

1. $\Pi_A v = 0$ if $v \in \ker A$; and
2. $\Pi_A v = v$ if $v \in (\ker A)^\perp$.¹⁰

¹⁰ Recall that $(\ker A)^\perp = \text{im } A^\top$.

Claim 2.33. For $A \in S^n$ with spectral decomposition $A = V\Lambda V^\top = \sum_i \lambda_i v_i v_i^\top$ where V is orthonormal, we have,

$$\Pi_A = \sum_{\substack{i \\ \lambda_i \neq 0}} v_i v_i^\top = A^{+/2} A A^{+/2} = A A^+ = A^+ A. \quad (2.17)$$

Corollary 2.34. In particular, $\Pi_A^+ = \Pi_A$.

Claim 2.35. Consider a real symmetric matrix $A = XYX^\top$, where X is real and invertible and Y is real and symmetric. Then,

$$A^+ = \Pi_A (X^\top)^{-1} Y^+ X^{-1} \Pi_A. \quad (2.18)$$

Now, let us return to solving linear systems.

Lemma 2.36. Given an invertible square lower triangular matrix \mathcal{L} or an invertible square upper triangular matrix \mathcal{U} , we can solve the linear systems $\mathcal{L}y = b$ and $\mathcal{U}z = b$ in time $\mathcal{O}(\text{nnz } \mathcal{L})$ and $\mathcal{O}(\text{nnz } \mathcal{U})$, respectively, where $\text{nnz } A$ denotes the number of non-zero entries of A .¹¹

¹¹ We need the matrices to be stored as an adjacency list to have fast access to their non-zero entries.

Proof sketch. We can iteratively solve the linear equations of $\mathcal{L}y = b$. As \mathcal{L} has full rank and is lower triangular, there is always one linear equation in a single variable.

These algorithms are known as forward and back substitution for lower and upper triangular matrices, respectively. \square

Thus, if we can decompose $A = \mathcal{L}\mathcal{L}^\top$ such that \mathcal{L} is invertible and lower triangular, we can solve the linear system of eq. (2.14) in time $\mathcal{O}(\text{nnz } \mathcal{L})$ by solving the two linear systems,

$$\mathcal{L}\mathbf{y} = \mathbf{b} \quad \text{and} \quad (2.19)$$

$$\mathcal{L}^\top \mathbf{x} = \mathbf{y}. \quad (2.20)$$

Fact 2.37 (Cholesky decomposition). *For any positive semi-definite matrix $A \in \mathbb{R}^{n \times n}$, there is a decomposition of the form $A = \mathcal{L}\mathcal{L}^\top$ where $\mathcal{L} \in \mathbb{R}^{n \times n}$ is lower triangular and positive semi-definite.*

We will see that the Cholesky decomposition can be computed efficiently when A is a graph Laplacian.

3

Probability

3.1 Random Walks

In this section, we study random walks on undirected weighted graphs $G = (V, E, w)$ with self-loops. A *random walk* visits a random sequence of vertices X_1, X_2, \dots , where

$$\mathbb{P}[X_{t+1} = v \mid X_t = u] = \frac{w(\{v, u\})}{d(u)} \quad (3.1)$$

and $d(u)$ is the weighted degree of vertex u , as we have defined previously.

Remark 3.1. The random walks considered here satisfy the *Markov property*, that is,

$$X_{t+1} \perp X_1, \dots, X_{t-1} \mid X_t.^1 \quad (3.2)$$

¹ So you can think of these random walks as Markov chains.

Moreover, we restrict our attention to time-homogeneous random walks, that is, the transition probabilities remain constant over time.

We can therefore model the update of a single round using the linear map $W \in \mathbb{R}^{|V| \times |V|}$ (called *transition matrix*),

$$W \doteq \begin{bmatrix} \frac{w(\{1,1\})}{d(1)} & \dots & \frac{w(\{n,1\})}{d(n)} \\ \vdots & \ddots & \vdots \\ \frac{w(\{1,n\})}{d(1)} & \dots & \frac{w(\{n,n\})}{d(n)} \end{bmatrix} = AD^{-1}. \quad (3.3)$$

A probability distribution over vertices is a vector $p \in \mathbb{R}^{|V|}$ such that $\mathbf{1}^\top p = 1$ and $p \geq \mathbf{0}$. If our initial distribution is p_0 , we have,

$$p_t = W^t p_0. \quad (3.4)$$

Observe that $W^t(u, v)$ denotes the probability to reach v from u in exactly t steps.

Definition 3.2 (Mixing). We say that a random walk W is *mixing* at step t iff for each $u, v \in V$, $W^t(u, v) \geq \frac{1}{2n}$.²

Definition 3.3 (Stationary distribution). A distribution $\pi \in \mathbb{R}^{|V|}$ is *stationary* iff $\pi = W\pi$.

Lemma 3.4. Every graph has the stationary distribution $\pi \doteq \frac{d}{1^\top d}$.

Proof. First, π is a distribution as,

$$\mathbf{1}^\top \pi = \frac{\mathbf{1}^\top d}{\mathbf{1}^\top d}$$

and clearly $\pi \geq 0$. We have,

$$W\pi = \frac{1}{\mathbf{1}^\top d} AD^{-1}d = \frac{1}{\mathbf{1}^\top d} A\mathbf{1} = \frac{d}{\mathbf{1}^\top d} = \pi. \quad \square$$

Remark 3.5. When the graph is connected, this is the unique stationary distribution.³

Lazy Random Walk

We would also like to have that we converge to this stationary distribution regardless of the initial distribution p_0 , but this is not true for general graphs, as is shown in fig. 3.1. A sufficient condition for convergence to the stationary distribution is, however, that all vertices have self-loops.⁴

Given the random walk W , the associated *lazy random walk* is given by,

$$\tilde{W} \doteq \frac{1}{2}I + \frac{1}{2}W,$$

that is, we add self-loops to each vertex with weight $1/2$ and halve all other weights. Observe that this does not change the stationary distribution of the random walk. This ensures that the following holds.

Theorem 3.6 (Convergence of lazy random walk). *For a connected graph, the lazy random walk converges to its unique stationary distribution irrespectively of the initial distribution p_0 ,*

$$\lim_{t \rightarrow \infty} \tilde{W}^t p_0 = \tilde{\pi} = \pi. \quad (3.5)$$

To prove this theorem, let us first understand the transition matrix in terms of the graph Laplacian.

² That is, the random walk is “half way” to being completely mixed.

³ In the context of Markov chains, irreducibility is sufficient for a unique stationary condition and equivalent to the transition graph being connected.



Figure 3.1: Consider the initial distribution $p_0(1) = 1, p_0(2) = 0$. Clearly, the random walk will forever oscillate between the two states.

⁴ It is easy to check that this ensures that the Markov chain is aperiodic, which together with irreducibility implies convergence to the unique stationary distribution.

Lemma 3.7. When v_1, \dots, v_n are the eigenvalues and ψ_1, \dots, ψ_n the corresponding eigenvectors of the normalized Laplacian matrix N , then \tilde{W} has eigenvalues $1 - v_i/2$ and (not necessarily orthogonal) eigenvectors $D^{1/2}\psi_i$.

Proof. Let us first express the transition matrix of the original random walk in terms of the normalized graph Laplacian,

$$\begin{aligned} W &= AD^{-1} = D^{1/2}(D^{-1/2}AD^{-1/2})D^{-1/2} \\ &= I + D^{1/2} \underbrace{(D^{-1/2}AD^{-1/2} - I)}_{-N} D^{-1/2} \\ &= I - D^{1/2}ND^{-1/2} \\ &= D^{1/2}(I - N)D^{-1/2}. \end{aligned} \quad (3.6)$$

By claim 2.7, $D^{1/2}(I - N)D^{-1/2}$ and $I - N$ have the same eigenvalues, namely $1 - v_i$. We also have,

$$\tilde{W} = \frac{1}{2}I + \frac{1}{2}(I - D^{1/2}ND^{-1/2}) = I - \frac{1}{2}D^{1/2}ND^{-1/2}, \quad (3.7)$$

implying that the eigenvalues of \tilde{W} are $1 - v_i/2$. Finally, we have,

$$\begin{aligned} \tilde{W}D^{1/2}\psi_i &= (I - \frac{1}{2}D^{1/2}ND^{-1/2})D^{1/2}\psi_i \\ &= D^{1/2}\psi_i - \frac{1}{2}D^{1/2}N\psi_i \\ &= D^{1/2}\psi_i - \frac{v_i}{2}D^{1/2}\psi_i \\ &= \left(1 - \frac{v_i}{2}\right)D^{1/2}\psi_i. \end{aligned} \quad \square$$

using that ψ_i is an eigenvector of N with corresponding eigenvalue v_i

Proof of theorem 3.6. TBD \square

Theorem 3.8 (Convergence rate of lazy random walk). For any unweighted connected graph G , we have that at time step t ,⁵

$$\|p_t - \pi\|_\infty \leq e^{-v_2 t/2} \sqrt{n}. \quad (3.8)$$

⁵ We will later see that v_2 is an indicator of the “connectedness” of G .

Proof. TBD \square

Hitting Time

Definition 3.9 (Hitting time). The *hitting time*,

$$H_{a,s} \doteq \min\{t \geq 1 \mid X_t = s, X_0 = a\}, \quad (3.9)$$

is the number of steps to reach s starting from a . We have,

$$h_s(a) \doteq \mathbb{E}[H_{a,s}] = 1 + \sum_{b \sim a} \frac{w(\{a, b\})}{d(a)} h_s(b). \quad (3.10)$$

Lemma 3.10. *If \tilde{x} is a solution to $L\tilde{x} = d - \|d\|_1 \mathbf{1}_s$,⁶ then*

$$h_s = \tilde{x} - \tilde{x}(s)\mathbf{1}. \quad (3.11)$$

⁶ We use $\mathbf{1}_s$ as a shorthand notation for $\mathbf{1}_{\{s\}}$.

Proof. For any $a \neq s$, we can equivalently write eq. (3.10) as,

$$\mathbf{1}_a^\top h_s = 1 + (W\mathbf{1}_a)^\top h_s \iff \mathbf{1}_a^\top (I - W^\top)h_s = 1.$$

This yields a linear system of $n - 1$ equations,

$$\mathbf{1} - \alpha \mathbf{1}_s = (I - \underbrace{W^\top}_{D^{-1}A})h_s, \quad (3.12)$$

where α is due to the remaining degree of freedom, as the entry corresponding to s is not fixed. Multiplying from the left with D , we obtain,

$$d - \alpha d(s)\mathbf{1}_s = (D - A)h_s = Lh_s.$$

Recall that $\ker L = \text{span}\{\mathbf{1}\}$, and hence, for h_s to exist, we must choose α such that $d - \alpha d(s)\mathbf{1}_s \perp \mathbf{1}$. We have,

$$\mathbf{1}^\top (d - \alpha d(s)\mathbf{1}_s) = \|d\|_1 - \alpha d(s) \stackrel{!}{=} 0 \iff \alpha = \frac{\|d\|_1}{d(s)}.$$

Finally, note that the solution \tilde{x} to $L\tilde{x} = d - \|d\|_1 \mathbf{1}_s$ is not unique.

Given that h_s is one solution, we have that any $\tilde{x} = h_s + c\mathbf{1}$ for $c \in \mathbb{R}$ is also a solution.⁷ Yet, we know that $h_s(s) = 0$, implying that $h_s = \tilde{x} - \tilde{x}(s)\mathbf{1}$. \square

⁷ This follows directly from the fact that $\ker L = \text{span}\{\mathbf{1}\}$.

Commutate Time

An issue with hitting times is that they do not need to be symmetric. This motivates the consideration of commute times, which correspond to the number of steps it takes to reach b from a and return to a .

Definition 3.11 (Commutate time). The *commute time* between a and b is defined as,

$$C_{a,b} \doteq H_{a,b} + H_{b,a}. \quad (3.13)$$

Remark 3.12. By definition, commute times are symmetric.

Lemma 3.13. *If \tilde{x} is a solution to $L\tilde{x} = \|d\|_1 (\mathbf{1}_a - \mathbf{1}_b)$, then*

$$\mathbb{E}[C_{a,b}] = (\mathbf{1}_a - \mathbf{1}_b)^\top \tilde{x} = \tilde{x}(a) - \tilde{x}(b). \quad (3.14)$$

The \tilde{x} can be interpreted as electrical voltages inducing flow that routes $\|d\|_1$ units from a to b .

TBD

Figure 3.2: Example where hitting times are not symmetric.

Proof. We write $\mathbf{b}_v \doteq \mathbf{d} - \|\mathbf{d}\|_1 \mathbf{1}_v$ and let $L\tilde{\mathbf{y}} = \mathbf{b}_b$ and $L\tilde{\mathbf{z}} = \mathbf{b}_a$. Then,

$$\begin{aligned} \mathbb{E}[C_{a,b}] &= \mathbf{h}_b(a) + \mathbf{h}_a(b) \\ &= \tilde{\mathbf{y}}(a) - \tilde{\mathbf{y}}(b) + \tilde{\mathbf{z}}(b) - \tilde{\mathbf{z}}(a) \\ &= (\mathbf{1}_a - \mathbf{1}_b)^\top (\tilde{\mathbf{y}} - \tilde{\mathbf{z}}). \end{aligned}$$

Observe that $\tilde{\mathbf{x}} \doteq \tilde{\mathbf{y}} - \tilde{\mathbf{z}}$ solves $L\tilde{\mathbf{x}} = \mathbf{b}_b - \mathbf{b}_a = \|\mathbf{d}\|_1 (\mathbf{1}_a - \mathbf{1}_b)$. \square

We will see in chapter 12 that the expected commute time is intimately related to the electrical energy required to route flow from a to b , also called the effective resistance between a and b .

3.2 Concentration

Theorem 3.14 (Markov's inequality). *For any random variable $X \geq 0$ and $t > 0$,*

$$\mathbb{P}[X \geq t] \leq \frac{\mathbb{E}[X]}{t}. \quad (3.15)$$

Proof. We have,

$$\mathbb{E}[X] = \int_0^\infty xf(x) dx \geq \int_t^\infty xf(x) dx \geq t \int_t^\infty f(x) dx = t\mathbb{P}[X \geq t],$$

where f is the probability density function of X . \square

Fact 3.15 (Jensen's inequality). *For a random variable X , if f is convex, then $\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$.⁸*

Theorem 3.16 (Bernstein concentration bound). *Given independent real-valued random variables $X_1, \dots, X_k \in \mathbb{R}$ such that $\mathbb{E}[X_i] = 0$ and $|X_i| \leq R$. Let $X \doteq \sum_i X_i$ and $\sigma^2 \doteq \text{Var}[X] = \sum_i \mathbb{E}[X_i^2]$. Then, for $t > 0$,*

$$\mathbb{P}[|X| \geq t] \leq 2 \exp\left(\frac{-t^2}{2Rt + 4\sigma^2}\right). \quad (3.16)$$

Proof. TBD \square

Corollary 3.17 (Chernoff bound). TBD

Proof. TBD \square

Theorem 3.18 (Bernstein matrix concentration bound). *Suppose $\mathbf{X}_1, \dots, \mathbf{X}_k \in \mathbb{R}^{n \times n}$ are independent symmetric matrix-valued random variables satisfying $\mathbb{E}[\mathbf{X}_i] = \mathbf{0}$ and $\|\mathbf{X}_i\|_2 \leq R$. Let $\mathbf{X} \doteq \sum_i \mathbf{X}_i$ and $\sigma^2 \doteq \text{Var}[\mathbf{X}] = \sum_i \mathbb{E}[\mathbf{X}_i^2]$. Then, for $t > 0$,*

$$\mathbb{P}[\|\mathbf{X}\|_2 \geq t] \leq 2n \exp\left(\frac{-t^2}{2Rt + 4\sigma^2}\right). \quad (3.17)$$

⁸ We prove the finite form in theorem A.1.

TBD

Figure 3.3: Jensen's inequality.

Proof. TBD □

3.3 Martingales

Definition 3.19 (Martingale). A *martingale* is a sequence of random variables Z_0, \dots, Z_k such that

$$\mathbb{E}[Z_i \mid Z_0, \dots, Z_{i-1}] = Z_{i-1}. \quad (3.18)$$

That is, conditional on the outcome of all the previous random variables, the expectation of Z_i equals Z_{i-1} .

Typically, we use martingales to show a statement such as “ Z_k is concentrated around $\mathbb{E}[Z_k]$ ”.

We can alternatively think of a martingale as the sequence of changes in $\{Z_i\}_i$. Let $X_i \doteq Z_i - Z_{i-1}$. The sequence of $\{X_i\}_i$ is called *martingale difference sequence*. The martingale condition is equivalent to,

$$\mathbb{E}[X_i \mid Z_0, \dots, Z_{i-1}] = \mathbb{E}[X_i \mid Z_0, X_1, \dots, X_{i-1}] = 0. \quad (3.19)$$

We can write,

$$Z_k = Z_0 + \sum_{i=1}^k Z_i - Z_{i-1} = Z_0 + \sum_{i=1}^k X_i. \quad (3.20)$$

Theorem 3.20. Suppose that $\{X_i\}_i$ form a scalar martingale difference sequence and $|X_i| \leq R$. Define the “pseudo-variance”,

$$W_i \doteq \sum_{j=1}^i \mathbb{E}[X_j^2 \mid X_1, \dots, X_{j-1}], \quad (3.21)$$

Then,

$$\mathbb{P}\left[\left|\sum_{i=1}^k X_i\right| \geq k \text{ and } W_k \leq \sigma^2\right] \leq C_2 \exp\left(-C_1 \frac{t^2}{Rt + \sigma^2}\right), \quad (3.22)$$

for constants C_1, C_2 .

By a simple union bound,

$$\mathbb{P}\left[\left|\sum_{i=1}^k X_i\right| \geq k\right] \leq \mathbb{P}\left[\left|\sum_{i=1}^k X_i\right| \geq k \text{ and } W_k \leq \sigma^2\right] + \mathbb{P}[W_k > \sigma^2]. \quad (3.23)$$

Proof. TBD □

4

Analysis

In this chapter, we study functions $f : S \rightarrow \mathbb{R}$ where $S \subseteq \mathbb{R}^n$.

4.1 First-order Taylor Approximations

Definition 4.1 (Gradient). The *gradient* of a function $f : S \rightarrow \mathbb{R}$ at point $x \in S$ is,

$$\nabla f(x) \doteq \begin{bmatrix} \frac{\partial f(x)}{\partial x(1)} & \cdots & \frac{\partial f(x)}{\partial x(n)} \end{bmatrix}^\top. \quad (4.1)$$

For a single-variable function $f : \mathbb{R} \rightarrow \mathbb{R}$ that is differentiable, we have for any $x, \delta \in \mathbb{R}$,

$$f(x + \delta) = f(x) + \frac{df(x)}{dx} \delta + o(|\delta|), \quad \text{where } \lim_{\delta \rightarrow 0} \frac{o(|\delta|)}{|\delta|} = 0,$$

using a first-order Taylor approximation around x . We can use a similar approximation when f is a multi-variable function.

Definition 4.2 (Fréchet differentiable). A function $f : S \rightarrow \mathbb{R}$ is (Fréchet) differentiable at $x \in S$ if there exists $g \in \mathbb{R}^n$ such that,¹

$$\lim_{\substack{\delta \in \mathbb{R}^n \\ \delta \rightarrow 0}} \frac{|f(x + \delta) - [f(x) + g^\top \delta]|}{\|\delta\|_2} = 0. \quad (4.2)$$

¹ Here, g can be understood as a candidate for $\nabla f(x)$ and $f(x) + g^\top \delta$ is a linear approximation of f around x .

This is equivalent to,

$$f(x + \delta) = f(x) + g^\top \delta + o(\|\delta\|_2), \quad (4.3)$$

for any $x \in S$ and $\delta \in \mathbb{R}^n$ where $\lim_{\delta \rightarrow 0} \frac{o(\|\delta\|_2)}{\|\delta\|_2} = 0$. This is also called a *first-order expansion* of f around x .

Remark 4.3. If this holds, $g = \nabla f(x)$.

Definition 4.4 (Continuously differentiable). We say that $f : S \rightarrow \mathbb{R}$ is *continuously differentiable* on S if it is differentiable and its gradient is continuous on S .

Fact 4.5 (Taylor's theorem, first-order form). If $f : S \rightarrow \mathbb{R}$ is continuously differentiable, then for all $x, y \in S$,

$$f(y) = f(x) + \nabla f(z)^\top (y - x), \quad (4.4)$$

for some $z \in [x, y] \doteq \{\theta x + (1 - \theta)y \mid \theta \in [0, 1]\}$.

Taylor's theorem implies that f can be approximated by the affine function,

$$y \rightarrow f(x) + \nabla f(x)^\top (y - x),$$

when y is "close to" x .

First-order Optimality Conditions

Definition 4.6 (Stationary point). Given a function $f : S \rightarrow \mathbb{R}$, a point $x \in S$ where $\nabla f(x) = 0$ is called a *stationary point* of f .²

Theorem 4.7 (First-order optimality condition). If $x \in S$ is a local extremum of a differentiable function $f : S \rightarrow \mathbb{R}$, then $\nabla f(x) = 0$.³

Proof. Assume x is a local minimum of f . Then, for all $d \in \mathbb{R}^n$ and for all small enough $\lambda \in \mathbb{R}$, we have $f(x) \leq f(x + \lambda d)$, so,

$$\begin{aligned} 0 &\leq f(x + \lambda d) - f(x) \\ &= \lambda \nabla f(x)^\top d + o(\lambda \|d\|_2). \end{aligned}$$

Dividing by λ and taking the limit $\lambda \rightarrow 0$, we obtain,

$$0 \leq \nabla f(x)^\top d + \lim_{\lambda \rightarrow 0} \frac{o(\lambda \|d\|_2)}{\lambda} = \nabla f(x)^\top d.$$

Take $d \doteq -\nabla f(x)$.⁴ Then, $0 \leq -\|\nabla f(x)\|_2^2$, so $\nabla f(x) = 0$. □

4.2 Directional Derivatives

Definition 4.8 (Jacobian). Given a vector-valued function,

$$g : \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad x \mapsto \begin{bmatrix} g_1(x) \\ \vdots \\ g_m(x) \end{bmatrix},$$

where $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$, the *Jacobian* of g at $x \in \mathbb{R}^n$ is,

$$Dg(x) \doteq \begin{bmatrix} Dg_1(x) \\ \vdots \\ Dg_m(x) \end{bmatrix} \doteq \begin{bmatrix} \frac{\partial g(x)}{\partial x(1)} & \cdots & \frac{\partial g(x)}{\partial x(n)} \end{bmatrix}. \quad (4.5)$$

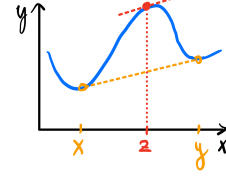


Figure 4.1: Illustration of Taylor's theorem. The affine approximation is shown in orange.

² Being a stationary point is not sufficient for optimality. Take for example the point $x \doteq 0$ of $f(x) \doteq x^3$.

³ Here it is important that we have chosen $S \subseteq \mathbb{R}^n$ to be open. When S is not open, an extremum could be on the boundary of the domain, where the gradient is non-zero.

using a first-order expansion of f around x

⁴ We can only take this step because we assumed that S is open.

Remark 4.9. For $f : S \rightarrow \mathbb{R}$ and any $\mathbf{x} \in S$, $Df(\mathbf{x}) = \nabla f(\mathbf{x})^\top$.

Definition 4.10. For a vector-valued function $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, we define,

$$\nabla \mathbf{g}(\mathbf{x}) \doteq D\mathbf{g}(\mathbf{x})^\top = \begin{bmatrix} \nabla g_1(\mathbf{x}) & \cdots & \nabla g_m(\mathbf{x}) \end{bmatrix}. \quad (4.6)$$

Definition 4.11 (Directional derivative). Let $f : S \rightarrow \mathbb{R}$ be differentiable at $\mathbf{x} \in \mathbb{R}^n$. Given $\mathbf{d} \in \mathbb{R}^n$, the *directional derivative* of f at \mathbf{x} in the direction \mathbf{d} is,

$$Df(\mathbf{x})[\mathbf{d}] \doteq \lim_{\lambda \rightarrow 0} \frac{f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x})}{\lambda}. \quad (4.7)$$

Lemma 4.12. $Df(\mathbf{x})[\mathbf{d}] = \nabla f(\mathbf{x})^\top \mathbf{d} = Df(\mathbf{x})\mathbf{d}$.

Proof. Using a first-order expansion, we have,

$$f(\mathbf{x} + \lambda \mathbf{d}) = f(\mathbf{x}) + \lambda \nabla f(\mathbf{x})^\top \mathbf{d} + o(\lambda \|\mathbf{d}\|_2).$$

Dividing by λ yields,

$$\frac{f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x})}{\lambda} = \nabla f(\mathbf{x})^\top \mathbf{d} + \underbrace{\frac{o(\lambda \|\mathbf{d}\|_2)}{\lambda}}_{\rightarrow 0}.$$

Taking the limit $\lambda \rightarrow 0$ gives the desired result. \square

4.3 Second-order Taylor Approximations

Example 4.13. Given $f : S \rightarrow \mathbb{R}$ and any $\mathbf{x} \in S$, consider the vector-valued function ∇f . We have,

$$D\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial \nabla f(\mathbf{x})}{\partial x(1)} & \cdots & \frac{\partial \nabla f(\mathbf{x})}{\partial x(n)} \end{bmatrix} = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x(1) \partial x(1)} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x(n) \partial x(1)} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x(1) \partial x(n)} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x(n) \partial x(n)} \end{bmatrix}. \quad (4.8)$$

Definition 4.14 (Hessian). Given $f : S \rightarrow \mathbb{R}$ that is twice differentiable coordinatewise, we define the *Hessian* $\mathbf{H}_f(\mathbf{x}) \in \mathbb{R}^{n \times n}$ of f at a point $\mathbf{x} \in S$ as,

$$\mathbf{H}_f(\mathbf{x})(i, j) \doteq \frac{\partial^2 f(\mathbf{x})}{\partial x(i) \partial x(j)} \quad (4.9)$$

Remark 4.15. By eq. (4.8), $\mathbf{H}_f(\mathbf{x}) = (D\nabla f(\mathbf{x}))^\top$.⁵

⁵ Sometimes, $\mathbf{H}_f(\mathbf{x}) = \nabla^2 f(\mathbf{x})$ is written (informally!).

Definition 4.16 (Twice Fréchet Differentiable). We say a function $f : S \rightarrow \mathbb{R}$ is *twice (Fréchet) differentiable* at $\mathbf{x} \in S$ if there exists $\mathbf{g} \in \mathbb{R}^n$ and a matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$ such that,⁶

$$\lim_{\substack{\delta \in \mathbb{R}^n \\ \delta \rightarrow 0}} \frac{|f(\mathbf{x} + \delta) - [f(\mathbf{x}) + \mathbf{g}^\top \delta + \frac{1}{2} \delta^\top \mathbf{M} \delta]|}{\|\delta\|_2^2} = 0. \quad (4.10)$$

This is equivalent to,

$$f(\mathbf{x} + \delta) = f(\mathbf{x}) + \mathbf{g}^\top \delta + \frac{1}{2} \delta^\top \mathbf{M} \delta + o(\|\delta\|_2^2), \quad (4.11)$$

for any $\mathbf{x} \in S$ and $\delta \in \mathbb{R}^n$ where $\lim_{\delta \rightarrow 0} \frac{o(\|\delta\|_2^2)}{\|\delta\|_2^2} = 0$. This is also called a *second-order expansion* of f around \mathbf{x} .

Remark 4.17. If this holds,

1. $\mathbf{g} = \nabla f(\mathbf{x})$,
2. $\mathbf{M} = \mathbf{H}_f(\mathbf{x})$, and
3. $\mathbf{H}_f(\mathbf{x}) = \mathbf{H}_f(\mathbf{x})^\top$.

Definition 4.18 (Twice continuously differentiable). We say that $f : S \rightarrow \mathbb{R}$ is *twice continuously differentiable* on S if it is twice differentiable and the gradient and Hessian are continuous on S .

Fact 4.19 (Taylor's theorem, second-order form). If $f : S \rightarrow \mathbb{R}$ is twice continuously differentiable, then for all $\mathbf{x}, \mathbf{y} \in S$,

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2} (\mathbf{y} - \mathbf{x})^\top \mathbf{H}_f(\mathbf{z}) (\mathbf{y} - \mathbf{x}), \quad (4.12)$$

for some $\mathbf{z} \in [\mathbf{x}, \mathbf{y}]$.

Second-order Optimality Conditions

Theorem 4.20 (Necessary second-order optimality condition). Let $f : S \rightarrow \mathbb{R}$ be twice continuously differentiable at $\mathbf{x} \in S$ and S be open. Then, if \mathbf{x} is a local minimum, $\mathbf{H}_f(\mathbf{x})$ is positive semi-definite.⁷

Proof. Suppose $\mathbf{x} \in S$ is a local minimum. By the first-order optimality condition, $\nabla f(\mathbf{x}) = 0$. For any direction $\mathbf{d} \in \mathbb{R}^n$ and small enough $\lambda \in [-\epsilon, \epsilon] \setminus \{0\}$,

$$\begin{aligned} 0 &\leq f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x}) \\ &= \frac{1}{2} \lambda^2 \mathbf{d}^\top \mathbf{H}_f(\mathbf{x}) \mathbf{d} + o(\lambda^2 \|\mathbf{d}\|_2^2). \end{aligned}$$

Multiplying both sides by $2/\lambda^2$ and taking the limit $\lambda \rightarrow 0$, we obtain,

$$0 \leq \mathbf{d}^\top \mathbf{H}_f(\mathbf{x}) \mathbf{d} + \lim_{\lambda \rightarrow 0} \frac{o(\lambda^2 \|\mathbf{d}\|_2^2)}{\lambda^2} = \mathbf{d}^\top \mathbf{H}_f(\mathbf{x}) \mathbf{d}. \quad \square$$

⁶ You can think of \mathbf{g} as a candidate for $\nabla f(\mathbf{x})$, \mathbf{M} as a candidate for $\mathbf{H}_f(\mathbf{x})$, and $f(\mathbf{x}) + \mathbf{g}^\top \delta + \frac{1}{2} \delta^\top \mathbf{M} \delta$ is a quadratic approximation of f around \mathbf{x} .

⁷ In other words, if $\mathbf{H}_f(\mathbf{x})$ has a negative eigenvalue, \mathbf{x} cannot be a minimum. Intuitively, you can think of $\mathbf{H}_f(\mathbf{x})$ as the curvature of f around \mathbf{x} , and therefore, a negative eigenvalue indicates that the function value can be decreased by moving in some direction.

using that f is locally minimized at \mathbf{x}

using a second-order expansion

Theorem 4.21 (Sufficient second-order optimality condition). *Let $f : S \rightarrow \mathbb{R}$ be twice continuously differentiable at $\mathbf{x} \in S$ and S be open. Then, if \mathbf{x} is a stationary point and $\mathbf{H}_f(\mathbf{x})$ is positive definite, \mathbf{x} is a local minimum.*

Proof. Suppose that $\mathbf{x} \in S$ is stationary, i.e., $\nabla f(\mathbf{x}) = \mathbf{0}$, and $\mathbf{H}_f(\mathbf{x})$ is positive definite. For any direction $\mathbf{d} \in \mathbb{R}^n$ and small enough $\lambda \in [-\epsilon, \epsilon] \setminus \{0\}$,

$$\begin{aligned}
 f(\mathbf{x} + \lambda \mathbf{d}) &= f(\mathbf{x}) + \underbrace{\nabla f(\mathbf{x})^\top \mathbf{d}}_0 + \frac{1}{2} \lambda^2 \mathbf{d}^\top \mathbf{H}_f(\mathbf{x}) \mathbf{d} + o(\lambda^2 \|\mathbf{d}\|_2^2) && \text{using a second-order expansion} \\
 &\geq f(\mathbf{x}) + \frac{1}{2} \lambda^2 \lambda_{\min}(\mathbf{H}_f(\mathbf{x})) \|\mathbf{d}\|_2^2 + o(\lambda^2 \|\mathbf{d}\|_2^2) && \text{by corollary 2.6} \\
 &\geq f(\mathbf{x}) + \frac{1}{4} \lambda_{\min}(\mathbf{H}_f(\mathbf{x})) \|\mathbf{d}\|_2^2 && \text{for small enough } \lambda \\
 &> f(\mathbf{x}) && \square \quad \text{using positive definiteness of } \mathbf{H}_f(\mathbf{x})
 \end{aligned}$$

PART II

Convex Optimization

5

Convex Geometry

We want to develop a better understanding of optimization problems. The general form of an *optimization problem* is,

$$\min_{\substack{\mathbf{y} \in \mathbb{R}^n \\ \mathbf{g}(\mathbf{y}) \leq \mathbf{b}}} f(\mathbf{y}), \quad (5.1)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is the function to be minimized, $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a vector-valued function of m constraints with thresholds $\mathbf{b} \in \mathbb{R}^m$.¹

Definition 5.1 (Feasible set). We call

$$\mathcal{F} \doteq \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{g}(\mathbf{x}) \leq \mathbf{b}\} \quad (5.2)$$

the *feasible set*. We call

- $\mathbf{x} \in \mathcal{F}$ a *feasible point*; and
- $\mathbf{x} \notin \mathcal{F}$ an *infeasible point*.

Definition 5.2 (Optimal solution). We say that $\mathbf{x}^* \in \mathbb{R}^n$ is *optimal* if $\mathbf{x}^* \in \mathcal{F}$ and $\forall \mathbf{x} \in \mathcal{F} : f(\mathbf{x}^*) \leq f(\mathbf{x})$.

Let us look at a sufficient condition for optimal solutions.

Theorem 5.3 (Extreme value theorem). Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be continuous, and let $\mathcal{F} \subseteq \mathbb{R}^n$ be non-empty, bounded, and closed. Then, f is bounded on \mathcal{F} and has an optimal solution.

5.1 Convex Sets & Functions

Definition 5.4 (Convex set). A set $S \subseteq \mathbb{R}^n$ is *convex* iff

$$\forall \mathbf{x}, \mathbf{y} \in S : \forall \theta \in [0, 1] : \theta \mathbf{x} + (1 - \theta) \mathbf{y} \in S. \quad (5.3)$$

¹ It suffices to consider minimization problems. If we want to maximize a function f , this is equivalent to minimizing the function $-f$.

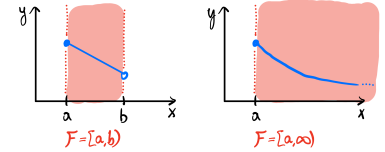


Figure 5.1: Examples of optimization problems without an optimal solution.

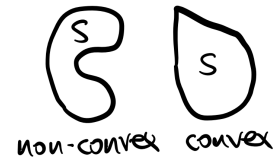


Figure 5.2: Example of a non-convex and a convex set.

Definition 5.5 (Convex function). For a convex set $S \subseteq \mathbb{R}^n$, a function $f : S \rightarrow \mathbb{R}$ is *convex* on S iff

$$\forall \mathbf{x}, \mathbf{y} \in S : \forall \theta \in [0, 1] : f(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta) f(\mathbf{y}). \quad (5.4)$$

Similarly, we call f *strictly convex* on S iff

$$\forall \mathbf{x}, \mathbf{y} \in S : \forall \theta \in [0, 1] : f(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) < \theta f(\mathbf{x}) + (1 - \theta) f(\mathbf{y}). \quad (5.5)$$

Remark 5.6. If the function f is convex on S , we say that the function $-f$ is *concave* on S .

Remark 5.7. We call the optimization problem from eq. (5.1) *convex* iff f and g_i are convex functions. Sometimes it is useful to equivalently write

$$\min_{\mathbf{y} \in S} f(\mathbf{y}) \quad \text{for } S \doteq \{\mathbf{y} \in \mathbb{R}^n \mid g(\mathbf{y}) \leq \mathbf{b}\}. \quad (5.6)$$

Observe that these characterizations are equivalent, as any convex set S can be characterized in terms of convex constraints.

Fact 5.8. A differentiable and convex function f , whose domain $S \subseteq \mathbb{R}^n$ is open and convex, is always continuously differentiable.

In the following, we will assume that $S \subseteq \mathbb{R}^n$ is open.

Definition 5.9 (Epigraph). The *epigraph* of a function $f : S \rightarrow \mathbb{R}$ is

$$\text{epi}(f) \doteq \{(\mathbf{x}, y) \mid f(\mathbf{x}) \leq y\} \subseteq \mathbb{R}^{n+1}. \quad (5.7)$$

Exercise 5.10. The function $f : S \rightarrow \mathbb{R}$ is convex iff $\text{epi}(f)$ is convex.

Definition 5.11 ((Sub-)level set). Given a function $f : S \rightarrow \mathbb{R}$, we call,

$$S_\alpha(f) \doteq \{\mathbf{x} \in S \mid f(\mathbf{x}) \leq \alpha\}, \quad (5.8)$$

$$L_\alpha(f) \doteq \{\mathbf{x} \in S \mid f(\mathbf{x}) = \alpha\}, \quad (5.9)$$

its α -sub-level set and α -level set, respectively.

Exercise 5.12. Any α -sub-level set of a convex function is convex.²

5.2 First-order Characterization of Convexity

Theorem 5.13 (First-order characterization of convexity). Consider a differentiable function $f : S \rightarrow \mathbb{R}$. Then, f is convex iff

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \quad (5.10)$$

for all $\mathbf{x}, \mathbf{y} \in S$. Moreover, f is strictly convex iff

$$f(\mathbf{y}) > f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \quad (5.11)$$

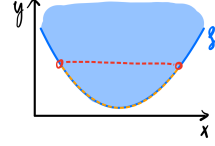


Figure 5.3: Example of a convex function. Any line between two points on f , lies “above” f . The epigraph of f is shown in blue.

² Note that the other direction does not hold! Take $f(x) \doteq x^3$ as an example. Functions whose sub-level sets are convex are called *quasiconvex*.

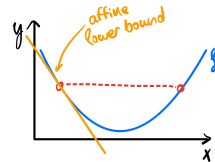


Figure 5.4: The first-order characterization characterizes convexity in terms of affine lower bounds.

Proof. We first prove the statement about convexity.

- “ \Rightarrow ”: Fix any $\mathbf{x}, \mathbf{y} \in S$. As f is convex,

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

for all $\theta \in [0, 1]$. We can rearrange to,

$$f(\underbrace{(1 - \theta)\mathbf{x} + \theta\mathbf{y}}_{\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})}) - f(\mathbf{y}) \leq \theta(f(\mathbf{x}) - f(\mathbf{y})).$$

Dividing by θ yields,

$$\frac{f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\theta} \leq f(\mathbf{y}) - f(\mathbf{x}).$$

Taking the limit $\theta \rightarrow 0$ on both sides gives the directional derivative at \mathbf{x} in direction $\mathbf{y} - \mathbf{x}$,

$$\nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) = Df(\mathbf{x})[\mathbf{y} - \mathbf{x}] \leq f(\mathbf{y}) - f(\mathbf{x}).$$

- “ \Leftarrow ”: Fix any $\mathbf{x}, \mathbf{y} \in S$ and let $\mathbf{z} \doteq \theta\mathbf{y} + (1 - \theta)\mathbf{x}$. We have,

$$\begin{aligned} f(\mathbf{y}) &\geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{z}), \quad \text{and} \\ f(\mathbf{x}) &\geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{z}). \end{aligned}$$

We also have $(1 - \theta)(\mathbf{y} - \mathbf{x}) = \mathbf{y} - \mathbf{z}$ and $\theta(\mathbf{y} - \mathbf{x}) = \mathbf{x} - \mathbf{z}$. Hence,

$$\begin{aligned} \theta f(\mathbf{y}) + (1 - \theta)f(\mathbf{x}) &\geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top \underbrace{(\theta(\mathbf{y} - \mathbf{z}) + (1 - \theta)(\mathbf{x} - \mathbf{z}))}_0 \\ &= f(\theta\mathbf{y} + (1 - \theta)\mathbf{x}). \end{aligned}$$

Finally, observe that the statement about strict convexity can be proven analogously by making the inequalities strict. \square

Theorem 5.14. Let $f : S \rightarrow \mathbb{R}$ be a convex and differentiable function. Then, if $\mathbf{x} \in S$ is a stationary point of f , then \mathbf{x} is a global minimum of f .

Proof. By the first-order characterization of convexity, we have for any $\mathbf{y} \in S$,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \underbrace{\nabla f(\mathbf{x})^\top}_{0} (\mathbf{y} - \mathbf{x}) = f(\mathbf{x}). \quad \square$$

5.3 Second-order Characterization of Convexity

Theorem 5.15 (Second-order characterization of convexity). Consider a twice continuously differentiable function $f : S \rightarrow \mathbb{R}$.³

1. f is convex iff $\mathbf{H}_f(\mathbf{x})$ is positive semi-definite for all $\mathbf{x} \in S$.
2. f is strictly convex iff $\mathbf{H}_f(\mathbf{x})$ is positive definite for all $\mathbf{x} \in S$.

³ Here we need our assumption that S is open.

Proof. We first prove the statement about convexity.

- “ \Leftarrow ”: Fix any $\mathbf{x}, \mathbf{y} \in S$. By the second-order form of Taylor’s theorem,

$$\begin{aligned} f(\mathbf{y}) &= f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2} \underbrace{(\mathbf{y} - \mathbf{x})^\top \mathbf{H}_f(\mathbf{z})(\mathbf{y} - \mathbf{x})}_{\geq 0} \\ &\geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}), \end{aligned}$$

for some $\mathbf{z} \in [\mathbf{x}, \mathbf{y}]$. This coincides with the first-order characterization of convexity.

- “ \Rightarrow ”: Fix any $\mathbf{x} \in S$ and $\mathbf{d} \in \mathbb{R}^n$. Note that, as S is open, for small enough $\lambda \in [-\epsilon, \epsilon] \setminus \{0\}$, $\mathbf{x} + \lambda \mathbf{d} \in S$. We have,

$$\begin{aligned} 0 &\leq f(\mathbf{x} + \lambda \mathbf{d}) - [f(\mathbf{x}) + \lambda \nabla f(\mathbf{x})^\top \mathbf{d}] \\ &= \frac{1}{2} \lambda^2 \mathbf{d}^\top \mathbf{H}_f(\mathbf{x}) \mathbf{d} + o(\lambda^2 \|\mathbf{d}\|_2^2). \end{aligned}$$

using the first-order characterization of convexity

using a second-order expansion

Multiplying both sides by $2/\lambda^2$ and taking the limit $\lambda \rightarrow 0$, we obtain,

$$0 \leq \mathbf{d}^\top \mathbf{H}_f(\mathbf{x}) \mathbf{d} + \lim_{\lambda \rightarrow 0} \frac{o(\lambda^2 \|\mathbf{d}\|_2^2)}{\lambda^2} = \mathbf{d}^\top \mathbf{H}_f(\mathbf{x}) \mathbf{d}.$$

The statement about strict convexity follows by using the first-order characterization of strict convexity instead and replacing inequalities with strict inequalities. \square

Gradient Descent

Gradient descent is a method for solving minimization problems such as eq. (5.1).

Definition 6.1 (Approximate solution). We say that a solution \mathbf{x}_k to the optimization problem $\min_{\mathbf{x} \in S} f(\mathbf{x})$ is ϵ -approximate iff

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \epsilon \quad (6.1)$$

for some $\epsilon > 0$, where $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in S} f(\mathbf{x})$.

For this chapter, we assume that the optimization problem is unconstrained, i.e., $S = \mathbb{R}^n$. In chapter 8, we explore how we can solve constrained optimization problems using Lagrange multipliers.

The idea of gradient descent is to iteratively take a step in the opposite direction of the gradient starting from some initial point $\mathbf{x}_0 \in \mathbb{R}^n$,

$$\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i - \alpha \nabla f(\mathbf{x}_i), \quad (6.2)$$

where $\alpha > 0$ is some learning rate.

When does gradient descent work? Clearly, we need that f is convex, otherwise gradient descent might not converge to the global minimum at all. But this is not enough! We also need to ensure that the gradient of f does not change arbitrarily when making very small steps, else the gradient direction would not be useful for us. This property is often called *smoothness*.

Finally, it is intuitively clear that we can do much better when we can ensure that the gradient of f is only close to zero around its minimum. If not, that is, we (almost) have “saddle points”, the step size of gradient descent will slow down and depending on the stopping criterion we might even return a point that is not the minimizer. To

TBD

Figure 6.1: Non-convex function.

TBD

Figure 6.2: Function whose gradient is close to zero at a non-optimal point.

exclude such functions from our analysis, we often assume that f satisfies the *PL condition*, or else is *strongly convex*.

You can think of smoothness as providing a quadratic upper bound to our function, whereas strong convexity provides a quadratic lower bound.

6.1 Smoothness

Definition 6.2 (Smoothness). Let $f : S \rightarrow \mathbb{R}$ be continuously differentiable. We say, f is β -smooth for some $\beta > 0$ iff for any $\mathbf{x}, \mathbf{y} \in S$

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq \beta \|\mathbf{x} - \mathbf{y}\|_2. \quad (6.3)$$

In other words, the gradient of f is β -Lipschitz.

Lemma 6.3. A twice continuously differentiable function $f : S \rightarrow \mathbb{R}$ is β -smooth iff for any $\mathbf{x} \in S$, $\lambda_{\max}(\mathbf{H}_f(\mathbf{x})) \leq \beta$.

Proof. TBD □

Lemma 6.4. A continuously differentiable function $f : S \rightarrow \mathbb{R}$ is β -smooth iff for any $\mathbf{x}, \mathbf{y} \in S$,

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2. \quad (6.4)$$

In words, $f(\mathbf{y})$ is upper bounded by a quadratic approximation based at $f(\mathbf{x})$.

Proof. TBD □

Analysis of Gradient Descent

A natural approach is to choose the gradient step of each iteration such that we minimize the upper bound (due to smoothness) based at the current solution,

$$\nabla_\delta \left(f(\mathbf{x}_i) + \nabla f(\mathbf{x}_i)^\top \delta + \frac{\beta}{2} \|\delta\|_2^2 \right) = \nabla f(\mathbf{x}_i) + \beta \delta \stackrel{!}{=} 0, \quad (6.5)$$

which is achieved for $\delta = -\frac{1}{\beta} \nabla f(\mathbf{x}_i)$. Thus,

$$f(\mathbf{x}_{i+1}) - f(\mathbf{x}_i) \leq \underbrace{\nabla f(\mathbf{x}_i)^\top \delta}_{-\frac{1}{\beta} \|\nabla f(\mathbf{x}_i)\|_2^2} + \underbrace{\frac{\beta}{2} \|\delta\|_2^2}_{\frac{1}{2\beta} \|\nabla f(\mathbf{x}_i)\|_2^2} = -\frac{1}{2\beta} \|\nabla f(\mathbf{x}_i)\|_2^2. \quad (6.6)$$

Moreover, due to the first-order characterization of convexity,

$$f(\mathbf{x}_i) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}_i)^\top (\mathbf{x}_i - \mathbf{x}^*) \leq \|\nabla f(\mathbf{x}_i)\|_2 \|\mathbf{x}_i - \mathbf{x}^*\|_2, \quad (6.7)$$

where the second inequality follows from Cauchy-Schwarz. Combining the previous two inequalities,

$$f(x_{i+1}) - f(x_i) \leq -\frac{1}{2\beta} \left(\frac{f(x_i) - f(x^*)}{\|x_i - x^*\|_2} \right)^2 \leq -\frac{1}{2\beta} \left(\frac{f(x_i) - f(x^*)}{\|x_0 - x^*\|_2} \right)^2. \quad (6.8)$$

using that $\|x_i - x^*\|_2 \leq \|x_0 - x^*\|_2$,
which follows from f decreasing in
every iteration and convexity

Theorem 6.5 (Convergence of gradient descent). *Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and β -smooth. The gradient descent scheme,*

$$x_{i+1} \doteq x_i - \frac{1}{\beta} \nabla f(x_i), \quad (6.9)$$

yields an ϵ -approximate solution x_k for any

$$k \geq \frac{2\beta \|x_0 - x^*\|_2^2}{\epsilon}.$$

Proof. We prove $f(x_k) - f(x^*) \leq \frac{2\beta \|x_0 - x^*\|_2^2}{k+1}$ by induction on the length of the computation k . Suppose $k = 0$, then by the smoothness of f ,

$$f(x_0) \leq f(x^*) - \underbrace{\nabla f(x^*)^\top}_{0} (x_0 - x^*) + \frac{\beta}{2} \|x_0 - x^*\|_2^2.$$

For the induction step, suppose that the statement holds for the k -th iterate. We write $\text{gap}_i \doteq f(x_i) - f(x^*)$. Using eq. (6.8),

$$\text{gap}_{k+1} - \text{gap}_k \leq -\frac{\text{gap}_k^2}{2\beta \|x_0 - x^*\|_2^2}.$$

Dividing by $\text{gap}_k \cdot \text{gap}_{k+1}$ and using that $\text{gap}_{k+1} > 0$ and $\text{gap}_k \geq \text{gap}_{k+1}$, we have,

$$\frac{1}{\text{gap}_k} - \frac{1}{\text{gap}_{k+1}} \leq -\frac{\text{gap}_k^2}{2\beta \|x_0 - x^*\|_2^2 \text{gap}_k \text{gap}_{k+1}} \leq -\frac{1}{2\beta \|x_0 - x^*\|_2^2}.$$

Thus,

$$\frac{1}{\text{gap}_{k+1}} \geq \frac{1}{2\beta \|x_0 - x^*\|_2^2} + \frac{1}{\text{gap}_k} \geq \frac{(k+1) + 1}{2\beta \|x_0 - x^*\|_2^2} \quad \square \quad \text{using the induction hypothesis}$$

6.2 Strong Convexity

We can improve our analysis, when we assume that f is strongly convex, that is lower bounded by a quadratic. Intuitively, this ensures that our steps are large when we are far away from the optimum.

Definition 6.6 (Strong convexity). Let $f : S \rightarrow \mathbb{R}$ be continuously differentiable. We say, f is μ -strongly convex for some $\mu > 0$ iff for any $\mathbf{x}, \mathbf{y} \in S$,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2. \quad (6.10)$$

Lemma 6.7. A twice continuously differentiable function $f : S \rightarrow \mathbb{R}$ is μ -strongly convex iff for any $\mathbf{x} \in S$, $\lambda_{\min}(\mathbf{H}_f(\mathbf{x})) \geq \mu$.

Proof. TBD □

Corollary 6.8. If f is β -smooth and μ -strongly convex, then $\mu \leq \beta$.

Definition 6.9 (Condition number). We call $\kappa \doteq \frac{\beta}{\mu}$ the *condition number* of a function f that is β -smooth and μ -strongly convex.

Often, a weaker condition known as *PL condition* is sufficient to design fast algorithms.

Definition 6.10 (Polyak-Łojasiewicz inequality). A continuously differentiable function $f : S \rightarrow \mathbb{R}$ satisfies the PL inequality with parameter $\mu > 0$ iff

$$\frac{1}{2} \|\nabla f(\mathbf{x})\|_2^2 \geq \mu(f(\mathbf{x}) - f(\mathbf{x}^*)), \quad (6.11)$$

where $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in S} f(\mathbf{x})$.

Intuitively, the norm of the gradient is tied to the suboptimality of the current solution.

Lemma 6.11. Let $f : S \rightarrow \mathbb{R}$ be continuously differentiable and μ -strongly convex. Then, f satisfies the PL condition.

Proof. As f is μ -strongly convex,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2,$$

for any $\mathbf{x}, \mathbf{y} \in S$. Taking the minimum with respect to \mathbf{y} on both sides yields,

$$f(\mathbf{x}^*) \geq f(\mathbf{x}) - \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|_2^2,$$

as the right-hand side is minimized for $\mathbf{y} = \mathbf{x} - \frac{1}{\mu} \nabla f(\mathbf{x})$. The PL inequality follows from rearranging the terms. □

Analysis of Gradient Descent

Theorem 6.12 (Convergence of gradient descent with a strongly convex objective). *Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is μ -strongly convex and β -smooth. The gradient descent scheme,*

$$\mathbf{x}_{i+1} \doteq \mathbf{x}_i - \frac{1}{\beta} \nabla f(\mathbf{x}_i), \quad (6.12)$$

yields an ϵ -approximate solution \mathbf{x}_k for any

$$k \geq \kappa \log \left(\frac{\beta \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2\epsilon} \right).$$

Proof. See the first graded homework. □

Remark 6.13. It turns out that the PL condition is sufficient to establish this convergence rate.

6.3 Acceleration

We can get an algorithm that converges substantially faster than vanilla gradient descent, using a method known as *accelerated gradient descent*. The key idea is to — instead of only tracking the upper bound that is due to smoothness — also use track lower bounds. In one iteration we might not make much progress in terms of reducing the upper bound (that is, improving our current solution), but instead increase the upper bound, which still reduces the error.

Theorem 6.14 (Convergence of accelerated gradient descent). *Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and β -smooth. The accelerated gradient descent scheme,*

$$\begin{aligned} a_i &\doteq \frac{i+1}{2}, \quad A_i \doteq \frac{(i+1)(i+2)}{4} \\ \mathbf{v}_0 &\doteq \mathbf{x}_0 - \frac{1}{2\beta} \nabla f(\mathbf{x}_0) \\ \mathbf{y}_i &\doteq \mathbf{x}_i - \frac{1}{\beta} \nabla f(\mathbf{x}_i) \\ \mathbf{x}_{i+1} &\doteq \frac{A_i \mathbf{y}_i + a_{i+1} \mathbf{v}_i}{A_{i+1}} \\ \mathbf{v}_{i+1} &\doteq \mathbf{v}_i - \frac{a_{i+1}}{\beta} \nabla f(\mathbf{x}_{i+1}), \end{aligned} \quad (6.13)$$

yields an ϵ -approximate solution \mathbf{x}_k for any

$$k \geq \sqrt{\frac{2\beta \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{\epsilon}}.$$

Here, y_i is the current solution (i.e., an upper bound), v_i is a lower bound, and x_i a point that trades improving the lower/upper bounds.

Proof. TBD □

Acceleration with Strongly Convex Objectives

Theorem 6.15 (Convergence of accelerated gradient descent with a strongly convex objective). *Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is μ -strongly convex and β -smooth. The accelerated gradient descent scheme,*

$$\begin{aligned} y_0 &\doteq x_0 \\ y_{i+1} &\doteq x_i - \frac{1}{\beta} \nabla f(x_i) \\ x_{i+1} &\doteq (1 + \theta) y_{i+1} + \theta y_i \end{aligned} \tag{6.14}$$

for $\theta \doteq \frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}$ yields an ϵ -approximate solution x_k for any

$$k \geq \sqrt{\kappa} \log \left(\frac{\beta \|x_0 - x^*\|_2^2}{\epsilon} \right).$$

Proof. See the first graded homework. □

7

Non-Euclidean Geometries

7.1 Mirror Descent

Lagrange Multipliers and Duality

8.1 Separating Hyperplanes

Definition 8.1 (Hyperplane). A *hyperplane* of dimension n is the subset,

$$H(\mathbf{n}, \mu) \doteq \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{n}^\top \mathbf{x} = \mu\}, \quad (8.1)$$

for some *normal* $\mathbf{n} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ and *threshold* $\mu \in \mathbb{R}$.

Every hyperplane divides \mathbb{R}^n into two half-spaces $\{\mathbf{x} \mid \mathbf{n}^\top \mathbf{x} \leq \mu\}$ and $\{\mathbf{x} \mid \mathbf{n}^\top \mathbf{x} \geq \mu\}$. It separates two sets, if they lie in different half-spaces.

Definition 8.2 (Separating hyperplane). We say a hyperplane H *separates* two sets A, B iff

$$\begin{aligned} \forall \mathbf{a} \in A : \mathbf{n}^\top \mathbf{a} &\leq \mu \quad \text{and} \\ \forall \mathbf{b} \in B : \mathbf{n}^\top \mathbf{b} &\geq \mu. \end{aligned} \quad (8.2)$$

If the inequalities are strict, we say that H *strictly* separates A and B .

If A, B are non-convex, we are not guaranteed that a separating hyperplane exists (e.g., a point cannot be separated from a ring around it). However, if we assume that A and B are convex, a separating hyperplane always exists.

Fact 8.3 (Separating hyperplane theorem). *Given two disjoint and non-empty convex subsets $A, B \subseteq \mathbb{R}^n$, there exists a separating hyperplane.*

However, it is not true that there always exists a strictly separating hyperplane. Consider $A \doteq \{(x, y) \mid x \leq 0\}$ and $B \doteq \{(x, y) \mid x > 0 \text{ and } y \geq \frac{1}{x}\}$. Clearly they are disjoint and convex; however, the only separating hyperplane is $H = \{(x, y) \mid x = 0\}$, which intersects A .

When we also assume that A and B are closed and bounded, a strictly separating hyperplane always exists.

TBD

Figure 8.1: Example where no strictly separating hyperplane exists.

Theorem 8.4 (Separating hyperplane theorem; closed, bounded sets). *Given two disjoint, closed, bounded, and non-empty convex subsets $A, B \subseteq \mathbb{R}^n$, there exists a strictly separating hyperplane.*

If $\mathbf{c} \in A, \mathbf{d} \in B$ are minimizers of $\min_{\mathbf{a} \in A, \mathbf{b} \in B} \|\mathbf{a} - \mathbf{b}\|_2$, then one such hyperplane is given by,

$$\mathbf{n} \doteq \mathbf{d} - \mathbf{c} \quad \text{and} \quad \mu \doteq \frac{1}{2} (\|\mathbf{d}\|_2^2 - \|\mathbf{c}\|_2^2). \quad (8.3)$$

TBD

Figure 8.2: Illustration of strictly separating hyperplane.

Proof. We want to show that $\mathbf{n}^\top \mathbf{b} > \mu$ for all $\mathbf{b} \in B$. Then, $\mathbf{n}^\top \mathbf{a} < \mu$ for all $\mathbf{a} \in A$ follows by symmetry. We have,

$$\begin{aligned} \mathbf{n}^\top \mathbf{d} - \mu &= (\mathbf{d} - \mathbf{c})^\top \mathbf{d} - \frac{1}{2} (\|\mathbf{d}\|_2^2 - \|\mathbf{c}\|_2^2) \\ &= \|\mathbf{d}\|_2^2 - \mathbf{d}^\top \mathbf{c} - \frac{1}{2} \|\mathbf{d}\|_2^2 + \frac{1}{2} \|\mathbf{c}\|_2^2 \\ &= \frac{1}{2} \|\mathbf{d} - \mathbf{c}\|_2^2 > 0. \end{aligned}$$

using the assumption that A, B are disjoint, close, and bounded, their distance is positive

Suppose for a contradiction that there exists a $\mathbf{u} \in B$ such that $\mathbf{n}^\top \mathbf{u} - \mu \leq 0$.

Consider the line defined by the distance minimizer \mathbf{d} and the point on the “wrong side” \mathbf{u} , $\mathbf{b}(\lambda) \doteq \mathbf{d} + \lambda(\mathbf{u} - \mathbf{d})$. Taking the derivative of the distance between $\mathbf{b}(\lambda)$ and \mathbf{c} and evaluating it at $\lambda = 0$ (which is when $\mathbf{b}(\lambda) = \mathbf{d}$), we obtain,

$$\begin{aligned} \left. \frac{d}{d\lambda} \|\mathbf{b}(\lambda) - \mathbf{c}\|_2^2 \right|_{\lambda=0} &= 2(\mathbf{d} - \lambda\mathbf{d} + \lambda\mathbf{u} - \mathbf{c})^\top (\mathbf{u} - \mathbf{d}) \Big|_{\lambda=0} \\ &= 2(\mathbf{d} - \mathbf{c})^\top (\mathbf{u} - \mathbf{d}). \end{aligned}$$

However,

$$\mathbf{n}^\top \mathbf{u} - \mu = (\mathbf{d} - \mathbf{c})^\top (\mathbf{u} - \mathbf{d}) + \underbrace{\mathbf{n}^\top \mathbf{d} - \mu}_{>0} \leq 0,$$

implies that $(\mathbf{d} - \mathbf{c})^\top (\mathbf{u} - \mathbf{d})$, and hence, the gradient are negative, which contradicts the minimality of \mathbf{d} . \square

8.2 Lagrange Multipliers and KKT Conditions

We will now discuss how we can treat constraints in a convex optimization problem,

$$\alpha^* \doteq \min_{\substack{\mathbf{y} \in \mathbb{R}^n \\ \mathbf{A}\mathbf{y} = \mathbf{b} \\ g(\mathbf{y}) \leq 0}} f(\mathbf{y}), \quad (8.4)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$, and we have k convex constraints $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$.

Remark 8.5. The linear constraints $A\mathbf{y} = \mathbf{b}$ are not necessary, as they can also be modeled using the (convex) constraints $A\mathbf{y} - \mathbf{b} \leq \mathbf{0}$ and $\mathbf{b} - A\mathbf{y} \leq \mathbf{0}$. We include them here to see later that programs with only linear constraints can be handled in a slightly different way.

We call this optimization problem the *primal program* and we will later see that it has an associated dual problem. We say that $\mathbf{y} \in \mathbb{R}^n$ is *primal feasible* iff $A\mathbf{y} = \mathbf{b}$ and $\mathbf{g}(\mathbf{y}) \leq \mathbf{0}$.

An Intuition

In the following, we want to answer the question: “When is a feasible point optimal?” To simplify things a bit, let us consider the optimization problem of minimizing f under the single constraint g . Suppose we know the feasible optimum \mathbf{y}^* . Then for infinitesimal δ , if $\delta \perp \nabla g(\mathbf{y}^*)$, then $\mathbf{y}^* + \delta$ and $\mathbf{y}^* - \delta$ are feasible. But then we must have that $\delta \perp \nabla f(\mathbf{y}^*)$, or else one direction would improve the objective. This tells us that,

$$\delta^\top \nabla f(\mathbf{y}^*) = 0 = \delta^\top \nabla g(\mathbf{y}^*),$$

and hence, there exists some $\lambda \in \mathbb{R}$ such that $\nabla f(\mathbf{y}^*) = \lambda \nabla g(\mathbf{y}^*)$. This is the fundamental intuition behind a *Lagrange multiplier*: the gradient of the objective at an optimal point is a linear combination of the gradients of the (tight) constraints.

We say that the constraint g_i is *tight* at \mathbf{y} iff $g_i(\mathbf{y}) = 0$. We know that if \mathbf{y} is feasible, then $\mathbf{y} + \delta$ is feasible for some infinitesimal δ if for all constraints $i \in [k]$ we have that either the constraint is not tight, $g_i(\mathbf{y}) < 0$, or $\mathbf{y} + \delta$ is “more” feasible, $\delta^\top \nabla g_i(\mathbf{y}) \leq 0$. The last condition tells us that for each tight constraint g_i , we get a half-space of feasible directions.

On the other hand, if $\mathbf{y}^* + \delta$ is feasible, then the objective f must “worsen”, $\delta^\top \nabla f(\mathbf{y}^*) \geq 0$.

We can therefore write the gradient of f at \mathbf{y}^* as a linear combination of the gradients of tight constraints.¹ Therefore, for some $\mathbf{x} \in \mathbb{R}^m$ and $\mathbf{s} \in \mathbb{R}^k$ with $s(i) \geq 0$ if $g_i(\mathbf{y}^*) = 0$ and $s(i) = 0$ otherwise,

$$-\nabla f(\mathbf{y}^*) = \sum_{i=1}^k s(i) \nabla g_i(\mathbf{y}^*) + \sum_{j=1}^m x(j) A(j, :). \quad (8.5)$$

The variables \mathbf{s} and \mathbf{x} are called the *dual variables* and are said to be *dual feasible* iff $\mathbf{s} \geq \mathbf{0}$. If \mathbf{y} is also primal feasible, we say that $(\mathbf{y}, \mathbf{x}, \mathbf{s})$ are *primal-dual feasible*. Equation (8.5) is equivalent to

$$\nabla f(\mathbf{y}^*) + \nabla g(\mathbf{y}^*)^\top \mathbf{s} + A^\top \mathbf{x} = \mathbf{0} \quad (8.6)$$

TBD

Figure 8.3: Illustration of simple constrained optimization.

TBD

Figure 8.4: Each tight constraint yields a half-space of feasible directions.

¹ Note that linear constraints are always tight.

along with the condition $\mathbf{g}(\mathbf{y}^*)^\top \mathbf{s} = \mathbf{0}$, ensuring that $s(i) = 0$ when constraint g_i is not tight.

Definition 8.6 (Karush-Kuhn-Tucker (KKT) conditions and Lagrangian). Points $\mathbf{y} \in \mathbb{R}^n, \mathbf{x} \in \mathbb{R}^m, \mathbf{s} \in \mathbb{R}^k$ satisfy the *Karush-Kuhn-Tucker conditions* iff,

$$\nabla_{\mathbf{y}} L(\mathbf{y}, \mathbf{x}, \mathbf{s}) = \mathbf{0} \quad \text{gradient condition} \quad (8.7)$$

$$\mathbf{g}(\mathbf{y})^\top \mathbf{s} = \mathbf{0} \quad \text{complementary slackness} \quad (8.8)$$

$$\mathbf{g}(\mathbf{y}) \leq \mathbf{0} \quad \text{and} \quad \mathbf{A}\mathbf{y} = \mathbf{b} \quad \text{primal feasibility} \quad (8.9)$$

$$\mathbf{s} \geq \mathbf{0} \quad \text{dual feasibility,} \quad (8.10)$$

ensures that $s(i)$ is forced to 0 when constraint g_i is not tight

where

$$L(\mathbf{y}, \mathbf{x}, \mathbf{s}) \doteq f(\mathbf{y}) + \mathbf{s}^\top \mathbf{g}(\mathbf{y}) + \mathbf{x}^\top (\mathbf{b} - \mathbf{A}\mathbf{y}) \quad (8.11)$$

is the *Lagrangian* of the optimization problem.

Remark 8.7. Note that,

$$\nabla_{\mathbf{y}} L(\mathbf{y}, \mathbf{x}, \mathbf{s}) = \nabla f(\mathbf{y}) + \nabla \mathbf{g}(\mathbf{y})^\top \mathbf{s} + \mathbf{A}^\top \mathbf{x},$$

which coincides with our intuition from eq. (8.6).

Intuitively, if $(\mathbf{y}, \mathbf{x}, \mathbf{s})$ satisfy the KKT conditions, then \mathbf{y} is optimal. It turns out that this intuition is correct, and we will prove this in the following section.

8.3 Lagrangian Duality

We have that if $(\mathbf{y}, \mathbf{x}, \mathbf{s})$ are primal-dual feasible,

$$L(\mathbf{y}, \mathbf{x}, \mathbf{s}) = f(\mathbf{y}) + \underbrace{\mathbf{s}^\top}_{\geq 0} \underbrace{\mathbf{g}(\mathbf{y})}_{\leq 0} + \mathbf{x}^\top \underbrace{(\mathbf{b} - \mathbf{A}\mathbf{y})}_{\mathbf{0}} \leq f(\mathbf{y}). \quad (8.12)$$

We can also write the primal problem as a two-player game in terms of the Lagrangian,

$$\alpha^* = \min_{\substack{\mathbf{y} \in \mathbb{R}^n \\ \mathbf{A}\mathbf{y} = \mathbf{b} \\ \mathbf{g}(\mathbf{y}) \leq \mathbf{0}}} f(\mathbf{y}) = \min_{\mathbf{y} \in \mathbb{R}^n} \max_{\substack{\mathbf{x} \in \mathbb{R}^m \\ \mathbf{s} \in \mathbb{R}^k \\ \mathbf{s} \geq \mathbf{0}}} L(\mathbf{y}, \mathbf{x}, \mathbf{s}). \quad (8.13)$$

This is because for a minimizing \mathbf{y} all constraints have to be satisfied and the Lagrangian simplifies to $L(\mathbf{y}, \mathbf{x}, \mathbf{s}) = f(\mathbf{y})$. If $\mathbf{b} - \mathbf{A}\mathbf{y} = \mathbf{0}$ was violated, making \mathbf{x} large sends $L(\mathbf{y}, \mathbf{x}, \mathbf{s}) \rightarrow \infty$. If $\mathbf{g}(\mathbf{y}) \leq \mathbf{0}$ was violated, making \mathbf{s} large sends $L(\mathbf{y}, \mathbf{x}, \mathbf{s}) \rightarrow \infty$.

We therefore have for any dual feasible (\mathbf{x}, \mathbf{s}) ,

$$\alpha^* = f(\mathbf{y}^*) \geq L(\mathbf{y}^*, \mathbf{x}, \mathbf{s}) \geq \min_{\mathbf{y} \in \mathbb{R}^n} L(\mathbf{y}, \mathbf{x}, \mathbf{s}) \doteq L(\mathbf{x}, \mathbf{s}). \quad (8.14)$$

Definition 8.8 (Dual program). The *dual program* is given as,

$$\beta^* \doteq \max_{\substack{\mathbf{x} \in \mathbb{R}^m \\ \mathbf{s} \in \mathbb{R}^k \\ \mathbf{s} \geq \mathbf{0}}} L(\mathbf{x}, \mathbf{s}) = \max_{\substack{\mathbf{x} \in \mathbb{R}^m \\ \mathbf{s} \in \mathbb{R}^k \\ \mathbf{s} \geq \mathbf{0}}} \min_{\mathbf{y} \in \mathbb{R}^n} L(\mathbf{y}, \mathbf{x}, \mathbf{s}). \quad (8.15)$$

Theorem 8.9 (Weak duality). $\beta^* \leq \alpha^*$.

Proof. This follows immediately from eq. (8.14). \square

Remark 8.10. Observe that the dual program is a convex optimization problem in disguise. We can equivalently consider the optimization problem,

$$-\beta^* = \min_{\substack{\mathbf{x} \in \mathbb{R}^m \\ \mathbf{s} \in \mathbb{R}^k \\ \mathbf{s} \geq \mathbf{0}}} -L(\mathbf{x}, \mathbf{s}) = \min_{\substack{\mathbf{x} \in \mathbb{R}^m \\ \mathbf{s} \in \mathbb{R}^k \\ \mathbf{s} \geq \mathbf{0}}} \max_{\mathbf{y} \in \mathbb{R}^n} -L(\mathbf{y}, \mathbf{x}, \mathbf{s}). \quad (8.16)$$

Note that $-L(\mathbf{y}, \mathbf{x}, \mathbf{s})$ is a linear function in the dual variables, hence convex, and $-L(\mathbf{x}, \mathbf{s})$ is a maximum of these functions, so also convex.

Definition 8.11 (Strong duality). If $\alpha^* = \beta^*$, we say that *strong duality* holds.

Remark 8.12. It is immediately clear that if we consider linear programs, that is, we have only linear constraints, strong duality always holds. This is because at any primal-dual feasible point (\mathbf{y}, \mathbf{x}) , we have that by definition $L(\mathbf{y}, \mathbf{x}) = f(\mathbf{y})$.

Before we analyze when strong duality holds, let us return to the KKT conditions and confirm our intuition from the previous section.

Theorem 8.13 (KKT conditions are necessary). *If strong duality holds, then the KKT conditions hold for primal-dual optimal $(\mathbf{y}^*, \mathbf{x}^*, \mathbf{s}^*)$.*

Proof. By strong duality,

$$L(\mathbf{y}^*, \mathbf{x}^*, \mathbf{s}^*) = \alpha^* = \beta^*.$$

As $L(\mathbf{y}, \mathbf{x}^*, \mathbf{s}^*)$ is a convex function in \mathbf{y} , we have that

$$\nabla_{\mathbf{y}} L(\mathbf{y}, \mathbf{x}^*, \mathbf{s}^*)|_{\mathbf{y}=\mathbf{y}^*} = \mathbf{0},$$

so the gradient condition is satisfied. Moreover, we have,

$$f(\mathbf{y}^*) = \alpha^* = L(\mathbf{y}^*, \mathbf{x}^*, \mathbf{s}^*) = f(\mathbf{y}^*) + \mathbf{s}^{\top} \mathbf{g}(\mathbf{y}^*) + \mathbf{x}^{\top} \underbrace{(\mathbf{b} - \mathbf{A}\mathbf{y}^*)}_{\mathbf{0}},$$

so $\mathbf{s}^{\top} \mathbf{g}(\mathbf{y}^*) = \mathbf{0}$ (complementary slack) holds. By assumption, $(\mathbf{y}^*, \mathbf{x}^*, \mathbf{s}^*)$ are primal-dual feasible. \square

Theorem 8.14 (KKT conditions are sufficient). *If the KKT conditions hold at $(\hat{\mathbf{y}}, \hat{\mathbf{x}}, \hat{\mathbf{s}})$, then they are primal-dual optimal and strong duality holds.*

Proof. Because $L(\mathbf{y}, \hat{\mathbf{x}}, \hat{\mathbf{s}})$ is a convex function, by the gradient condition, $\hat{\mathbf{y}}$ is its global minimizer. Therefore,

$$L(\hat{\mathbf{y}}, \hat{\mathbf{x}}, \hat{\mathbf{s}}) = \min_{\mathbf{y} \in \mathbb{R}^n} L(\mathbf{y}, \hat{\mathbf{x}}, \hat{\mathbf{s}}) = L(\hat{\mathbf{x}}, \hat{\mathbf{s}}) \leq \beta^*.$$

At the same time, using primal-dual feasibility and complementary slack,

$$L(\hat{\mathbf{y}}, \hat{\mathbf{x}}, \hat{\mathbf{s}}) = f(\hat{\mathbf{y}}) + \underbrace{\hat{\mathbf{s}}^\top \mathbf{g}(\hat{\mathbf{y}})}_0 + \hat{\mathbf{x}}^\top \underbrace{(\mathbf{b} - \mathbf{A}\hat{\mathbf{y}})}_0 = f(\hat{\mathbf{y}}) \geq \alpha^*.$$

Therefore, $\alpha^* \leq \beta^*$. By weak duality, we get the opposite inequality, and hence strong duality holds. \square

8.4 Slater's Condition

We will now discuss under which circumstances strong duality holds. In general, we can have that strong duality does not hold.

Example 8.15. TBD

Definition 8.16 (Slater's condition). *Slater's condition is satisfied iff there exists some $\mathbf{y} \in \mathbb{R}^n$ such that $\mathbf{A}\mathbf{y} = \mathbf{b}$ and $\mathbf{g}(\mathbf{y}) < \mathbf{0}$.² We say that \mathbf{y} is strictly feasible.*

² This is satisfied by most useful optimization problems.

Theorem 8.17. *If Slater's condition holds, then strong duality holds.*

Proof. TBD \square

8.5 Example: Duality of Maximum Flow and Minimum Cut

We can write the (directed) maximum flow problem in a graph G as,

$$\begin{array}{ll} \max_{\substack{F \in \mathbb{R} \\ f \in \mathbb{R}^{|E|} \\ Bf = F(\mathbf{1}_t - \mathbf{1}_s) \\ 0 \leq f \leq c}} F = - \min_{\substack{F \in \mathbb{R} \\ f \in \mathbb{R}^{|E|} \\ Bf = F(\mathbf{1}_t - \mathbf{1}_s) \\ 0 \leq f \leq c}} -F. \end{array} \quad (8.17)$$

Here, $\mathbf{0} \leq \mathbf{f}$ ensures that directions are respected, and $\mathbf{f} \leq \mathbf{c}$ ensures that edge capacities are respected. $F(\mathbf{1}_t - \mathbf{1}_s)$ is the demand of a flow routing F units from t to s . Observe that when there exists an s - t path, then there are flows \mathbf{f} strictly satisfying the constraint $\mathbf{0} \leq \mathbf{f} \leq \mathbf{c}$, and hence, Slater's criterion is satisfied. By strong duality, the above program is equivalent to its dual program,

$$- \max_{\substack{\mathbf{x} \in \mathbb{R}^{|V|} \\ \mathbf{s} \in \mathbb{R}^{|E|} \\ \mathbf{s} \geq \mathbf{0}}} \min_{\substack{F \in \mathbb{R} \\ f \in \mathbb{R}^{|E|} \\ f \geq \mathbf{0}}} -F + \mathbf{s}^\top (\mathbf{f} - \mathbf{c}) + \mathbf{x}^\top (F(\mathbf{1}_t - \mathbf{1}_s) - B\mathbf{f}) \quad (8.18)$$

$$= \min_{\substack{\mathbf{x} \in \mathbb{R}^{|V|} \\ \mathbf{s} \in \mathbb{R}^{|E|} \\ \mathbf{s} \geq \mathbf{0} \\ (\mathbf{1}_t - \mathbf{1}_s)^\top \mathbf{x} = 1 \\ \mathbf{s} \geq \mathbf{B}^\top \mathbf{x}}} \mathbf{s}^\top \mathbf{c} \quad (8.19)$$

$$= \min_{\substack{\mathbf{x} \in \mathbb{R}^{|V|} \\ (\mathbf{1}_t - \mathbf{1}_s)^\top \mathbf{x} = 1}} \sum_{e \in E} \max\{(\mathbf{B}^\top \mathbf{x})(e), 0\} \cdot c(e). \quad (8.20)$$

Observe that the dual program is a linear program computing the minimum cut.

8.6 Fenchel Conjugates

Definition 8.18 (Fenchel conjugate). Given a function $f : S \rightarrow \mathbb{R}$, its *Fenchel conjugate* is the function $f^* : \mathbb{R}^n \rightarrow \mathbb{R}$,³

$$f^*(\mathbf{z}) \doteq \sup_{\mathbf{y} \in S} \mathbf{z}^\top \mathbf{y} - f(\mathbf{y}). \quad (8.21)$$

³ In principle, f^* is a function defined over the dual space of S , but this will not be very important to us.

Remark 8.19. f^* is convex, as it is the maximum of linear functions.

Example 8.20. Let us consider the function $f(\mathbf{y}) \doteq \|\mathbf{y}\|$. Then we have,

$$\begin{aligned} f^*(\mathbf{z}) &= \sup_{\mathbf{y} \in \mathbb{R}^n} \mathbf{z}^\top \mathbf{y} - \underbrace{\|\mathbf{y}\|}_{\doteq \theta} \\ &= \sup_{\theta \geq 0} \theta (\|\mathbf{z}\|_* - 1) \\ &= \begin{cases} \infty & \|\mathbf{z}\|_* > 1 \\ 0 & \text{otherwise.} \end{cases} \end{aligned} \quad (8.22)$$

Recall that

$$\theta \|\mathbf{z}\|_* = \max_{\substack{\mathbf{y} \in \mathbb{R}^n \\ \|\mathbf{y}\| = \theta}} \mathbf{z}^\top \mathbf{y}.$$

Example 8.21. For the function $f(\mathbf{y}) \doteq \frac{1}{p} \|\mathbf{y}\|_p^p$, we have that $f^*(\mathbf{z}) = \frac{1}{q} \|\mathbf{z}\|_q^q$, where $\frac{1}{p} + \frac{1}{q} = 1$.

Example 8.22. When we have a primal program with only linear constraints,

$$\min_{\substack{\mathbf{y} \in \mathbb{R}^n \\ \mathbf{A}\mathbf{y} = \mathbf{b}}} f(\mathbf{y}),$$

we can write the dual program as,

$$\begin{aligned} \max_{\mathbf{x} \in \mathbb{R}^m} \min_{\mathbf{y} \in \mathbb{R}^n} f(\mathbf{y}) + \mathbf{x}^\top (\mathbf{b} - \mathbf{A}\mathbf{y}) &= \max_{\mathbf{x} \in \mathbb{R}^m} \mathbf{b}^\top \mathbf{x} - \max_{\mathbf{y} \in \mathbb{R}^n} (\mathbf{x}^\top \mathbf{A}\mathbf{y} - f(\mathbf{y})) \\ &= \max_{\mathbf{x} \in \mathbb{R}^m} \mathbf{b}^\top \mathbf{x} - f^*(\mathbf{A}^\top \mathbf{x}). \end{aligned} \quad (8.23)$$

Lemma 8.23 (Properties of the Fenchel conjugate). *If f is strictly convex (i.e., $\mathbf{H}_f \succ \mathbf{0}$) and ∇f is a surjective mapping onto \mathbb{R}^n ,*

1. $\nabla f(\nabla f^*(z)) = z$ and $\nabla f^*(\nabla f(y)) = y$;
2. $(f^*)^* = f$; and
3. $H_{f^*}(\nabla f(y)) = H_f^{-1}(y)$.⁴

Proof. TBD

⁴ Thus, if f is β -smooth, f^* is β -strongly convex. In other words, f^* has the “opposite” curvature of f .

□

9

Newton's Method

PART III

Spectral Graph Theory

Introduction to Spectral Graph Theory

Spectral graph theory studies graphs through linear algebra. The fundamental object that we will work with is the Laplacian matrix that we introduced in definition 1.4.

10.1 Eigenvalues of the Laplacian Matrix

Lemma 10.1. Denote by $X_1 \cup \dots \cup X_k = V$ the connected components of a graph G . Then we have for the Laplacian matrix \mathbf{L} of G ,

$$\ker \mathbf{L} = \text{span}\{\mathbf{1}_{X_1}, \dots, \mathbf{1}_{X_k}\}, \quad (10.1)$$

where $\mathbf{1}_X(v) = \mathbb{1}\{v \in X\}$. In particular, $\ker \mathbf{L} = \text{span}\{\mathbf{1}\}$ if G is connected.

Proof. TBD □

Corollary 10.2. When G has k connected components, then $\lambda_i(\mathbf{L}) = 0$ for all $i \in [k]$.

In the following, we will study connected graphs. If a graph consists of multiple connected components, we may study each individually.

Lemma 10.3. $\mathbf{L} \preceq 2\mathbf{D}$.

Proof. It can be shown that $\mathbf{D} + \mathbf{A} \succeq \mathbf{0}$,¹ so, $\mathbf{D} \succeq -\mathbf{A}$, and hence, $\mathbf{L} \preceq 2\mathbf{D}$. □

¹ The proof is similar to the proof that $\mathbf{L} = \mathbf{D} - \mathbf{A} \succeq \mathbf{0}$. $\mathbf{D} + \mathbf{A}$ is also called the *signless Laplacian matrix* and is an interesting object in its own right: its eigenvalues can be used to count edges and identify bipartite components.

Electrical Energy

Lemma 10.4. For any voltages $\mathbf{x} \in \mathbb{R}^{|V|}$,

$$\lambda_2(\mathbf{L}) \|\mathbf{x}\|_2^2 \leq \mathcal{E}(\mathbf{x}) \leq \lambda_n(\mathbf{L}) \|\mathbf{x}\|_2^2. \quad (10.2)$$

Proof. TBD □

Lemma 10.5. Voltages $\mathbf{x} \in \mathbb{R}^{|V|}$ routing demands $\mathbf{d} \perp \mathbf{1}$ satisfy,

$$\frac{\|\mathbf{d}\|_2^2}{\lambda_n(\mathbf{L})} \leq \mathcal{E}(\mathbf{x}) \leq \frac{\|\mathbf{d}\|_2^2}{\lambda_2(\mathbf{L})}. \quad (10.3)$$

Proof. Recall from eq. (1.18) that $\mathcal{E}(\mathbf{x}) = \mathbf{d}^\top \mathbf{L}^+ \mathbf{d}$. Using Courant-Fischer,

$$\begin{aligned} \lambda_n(\mathbf{L}^+) &= \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \mathbf{L}^+ \mathbf{x}}{\|\mathbf{x}\|_2^2} \geq \frac{\mathbf{d}^\top \mathbf{L}^+ \mathbf{d}}{\|\mathbf{d}\|_2^2} \\ \lambda_2(\mathbf{L}^+) &= \min_{\substack{\mathbf{x} \neq \mathbf{0} \\ \mathbf{x} \perp \mathbf{1}}} \frac{\mathbf{x}^\top \mathbf{L}^+ \mathbf{x}}{\|\mathbf{x}\|_2^2} \leq \frac{\mathbf{d}^\top \mathbf{L}^+ \mathbf{d}}{\|\mathbf{d}\|_2^2}, \end{aligned} \quad \text{using that } \ker \mathbf{L}^+ = \ker \mathbf{L} = \text{span}\{\mathbf{1}\}$$

using that $\mathbf{d} \neq \mathbf{0}$ and $\mathbf{d} \perp \mathbf{1}$. Rearranging the inequalities, we obtain,

$$\lambda_2(\mathbf{L}^+) \|\mathbf{d}\|_2^2 \leq \mathbf{d}^\top \mathbf{L}^+ \mathbf{d} \leq \lambda_n(\mathbf{L}^+) \|\mathbf{d}\|_2^2.$$

Finally, recall from eq. (2.16) that $\lambda_2(\mathbf{L}^+) = \lambda_n(\mathbf{L})^{-1}$ and $\lambda_n(\mathbf{L}^+) = \lambda_2(\mathbf{L})^{-1}$. □

Useful Inequalities

Lemma 10.6 (Path inequality). *We have that*

$$(n-1)P_n \succeq G_{1,n}, \quad (10.4)$$

where P_n is the unit weight path graph on n vertices and $G_{i,j}$ is the unit weight graph on n vertices with the single edge $\{i, j\}$.

Proof. Fix any $\mathbf{x} \in \mathbb{R}^n$ and let $\Delta(i) \doteq \mathbf{x}(i+1) - \mathbf{x}(i)$. We have,

$$\begin{aligned} \mathbf{x}^\top G_{1,n} \mathbf{x} &= (\mathbf{x}(n) - \mathbf{x}(1))^2 \\ &= \left(\sum_{i=1}^{n-1} \Delta(i) \right)^2 \\ &= (\mathbf{1}_{n-1}^\top \Delta)^2 \\ &\leq \|\mathbf{1}_{n-1}\|_2^2 \|\Delta\|_2^2 \\ &= (n-1) \sum_{i=1}^{n-1} \Delta(i)^2 \\ &= (n-1) \sum_{i=1}^{n-1} (\mathbf{x}(i+1) - \mathbf{x}(i))^2 \\ &= (n-1) \mathbf{x}^\top P_n \mathbf{x}. \end{aligned} \quad \text{using Cauchy-Schwarz} \quad \square$$

Lemma 10.7. For any unit weight graph G on n vertices with diameter² D ,

$$\lambda_2(G) \geq \frac{1}{nD}. \quad (10.5)$$

² The *diameter* of a graph is the maximum shortest distance between any two vertices.

Proof. We denote by $G^{i,j}$ the subgraph of G consisting of the shortest $i - j$ path. We have,

$$\begin{aligned} K_n &= \sum_{i < j} G_{i,j} \preceq \sum_{i < j} \underbrace{(j-i)}_{\leq D} \underbrace{G^{i,j}}_{\subseteq G} \\ &\preceq n^2 DG. \end{aligned}$$

analogously to the path inequality

Thus, $n^2 D \lambda_2(G) \geq \lambda_2(K_n) = n$, and hence, $\lambda_2(G) \geq \frac{1}{nD}$. \square

10.2 Examples

Lemma 10.8 (Spectrum of the complete graph).

$$\lambda_2(K_n) = \dots = \lambda_n(K_n) = n. \quad (10.6)$$

Proof. We have, $A = \mathbf{1}\mathbf{1}^\top - I$ and $D = (n-1)I$, so $L = nI - \mathbf{1}\mathbf{1}^\top$. For any $x \perp \mathbf{1}$, $Lx = nx$. \square

We now want to better understand λ_2 and λ_n for some common graphs. The tools we will use are the following:

- to lower bound $\lambda_2(G)$: Relate the eigenvalues of K_n and G , yielding $K_n \preceq f(n)G$. Knowing that $\lambda_2(K_n) = n$, we have, $f(n)\lambda_2(G) \geq n$, and hence, $\lambda_2(G) \geq \frac{n}{f(n)}$.
- to upper bound $\lambda_2(G)$: Due to Courant-Fischer,

$$\lambda_2(G) = \min_{\substack{x \perp \mathbf{1} \\ x \neq 0}} \frac{x^\top Lx}{x^\top x} \leq \frac{y^\top Ly}{y^\top y},$$

for any $y \perp \mathbf{1}, y \neq 0$. We can therefore find a so-called *test vector* y with these properties.

- to lower bound $\lambda_n(G)$: Similarly, due to Courant-Fischer,

$$\lambda_n(G) = \max_{x \neq 0} \frac{x^\top Lx}{x^\top x} \geq \frac{y^\top Ly}{y^\top y},$$

for any $y \neq 0$.

- to upper bound $\lambda_n(G)$: Using that $L \preceq 2D$, we have, $\lambda_n(G) \leq 2 \max_{v \in V} d(v)$, where $d(v)$ is the weighted degree of v .

Path Graph

Exercise 10.9. $\lambda_2(P_n) = \Theta\left(\frac{1}{n^2}\right)$.

Exercise 10.10. $\lambda_n(P_n) \in [1, 4]$.

Complete Binary Tree

Exercise 10.11. $\lambda_2(T_d) = \Theta\left(\frac{1}{n}\right)$.

Exercise 10.12. $\lambda_n(T_d) \in [1, 6]$.

Conductance and Expanders

11.1 Conductance

Definition 11.1 (Volume). The *volume* of a set of vertices $S \subseteq V$ is the sum of weighted degrees,

$$\text{vol}(S) \doteq \sum_{v \in S} d(v) = \mathbf{1}_S^\top \mathbf{d} = \mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S. \quad (11.1)$$

A *cut* $(S, V \setminus S)$ is a proper subset of vertices, $\emptyset \subset S \subset V$ partitioning vertices into two sets S and $V \setminus S$.

Definition 11.2 (Value of a cut). The *value* of a cut $(S, V \setminus S)$ is the sum of weights of crossing edges,

$$\begin{aligned} c(S) &\doteq \sum_{\substack{\{a,b\} \in E \\ a \in S, b \in V \setminus S}} w(\{a,b\}) \\ &= \sum_{\{a,b\} \in E} w(\{a,b\}) [\mathbf{1}_S(a) - \mathbf{1}_S(b)]^2 = \mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S. \end{aligned} \quad (11.2)$$

Definition 11.3 (Conductance of a cut). The *conductance* of a cut $(S, V \setminus S)$ is,

$$\phi(S) \doteq \frac{c(S)}{\min\{\text{vol}(S), \text{vol}(V \setminus S)\}} = \phi(V \setminus S) \in [0, 1], \quad (11.3)$$

where, if the graph has unit weights, $c(S) = |E(S, V \setminus S)|$ counts the number of crossing edges.

Remark 11.4. If $\text{vol}(S) \leq \text{vol}(V \setminus S)$, then we can write

$$\phi(S) = \frac{\mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S}{\mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S}. \quad (11.4)$$

Definition 11.5 (Conductance of a graph). The *conductance* of a graph G is the smallest conductance of all cuts,

$$\phi(G) \doteq \min_{\emptyset \subset S \subset V} \phi(S). \quad (11.5)$$

Thus, $\phi(G)$ is small if there is a “good” cut (with few crossing edges relative to the volume of the parts). In contrast, if $\phi(G)$ is large, then G is well-connected, i.e., there is no “good” cut.

Definition 11.6 (Expander and expander decomposition). For any $\phi \in (0, 1]$, if $\phi(G) \geq \phi$, then G is called a ϕ -*expander*.

A ϕ -*expander decomposition* of *quality* q is a partition of the vertex set $V = X_1 \cup \dots \cup X_k$ such that

1. $G[X_i]$ is a ϕ -expander; and
2. the number of edges not contained in any $G[X_i]$ is at most $q\phi m$, i.e. only “few” edges cross the parts.

In chapter 17, we discuss how we can efficiently find an expander decomposition. Let us consider a few examples.

Exercise 11.7. $\phi(K_n) = \frac{n}{2(n-1)}$. So K_n is a $\frac{1}{2}$ -expander.

Exercise 11.8. $\phi(P_n) = \frac{1}{n-1}$. So P_n is a $\frac{1}{n}$ -expander.

Lemma 11.9. If G is a connected ϕ -expander with unit weights, then we have for the diameter D of G ,

$$D = \mathcal{O}\left(\frac{\log m}{\phi}\right). \quad (11.6)$$

Proof. Fix any pair of vertices $s, t \in V$. Let $B(s, d)$ be the closed ball around s of radius d . Let $E(B(s, d))$ be the internal edge set of $B(s, d)$.

Observe that since G is connected, we have $|E(B(s, 0))| \geq 1$. Moreover, for each $d \geq 0$ where $|E(B(s, d))| \leq \frac{m}{2}$, we have by the definition of a ϕ -expander that

$$|E(B(s, d), V \setminus B(s, d))| \geq \phi \cdot |E(B(s, d))|.$$

Thus, $|E(B(s, d+1))| \geq (1 + \phi)|E(B(s, d))|$. Let d_{\max} be the largest integer such that $|E(B(s, d_{\max}))| \leq m/2$. We have $d_{\max} \leq 2 \log m / \phi$ as otherwise,

$$|E(B(s, d_{\max}))| > (1 + \phi)^{2 \log m / \phi} \geq (1 + \phi/2 + (\phi/2)^2)^{2 \log m / \phi} \geq e^{\log m} = m, \quad \text{using } e^x < 1 + x + x^2 \text{ for } x < 1.79$$

which gives a contradiction to the fact that the number of edges in G is m . Thus, for a radius of $2 \log m / \phi + 1$, the ball centered at s has more than $m/2$ edges.

Finally, follow the same argument from t . As both balls contain more than $m/2$ edges, they must intersect in at least one edge. But this implies that there is a s - t path of length $\mathcal{O}(\log m/\phi)$. \square

11.2 Cheeger's Inequality

Theorem 11.10 (Cheeger's inequality). *We have for a graph G and its normalized Laplacian matrix N ,*

$$\frac{\lambda_2(N)}{2} \leq \phi(G) \leq \sqrt{2\lambda_2(N)}. \quad (11.7)$$

In words, $\lambda_2(N)$ approximates the conductance $\phi(G)$ up to a square root. This is why we say that $\lambda_2(N)$ is a measure of connectivity of a graph.

Proof. TBD \square

11.3 Sparsity

A concept related to conductance is the notion of *sparsity*.

Definition 11.11 (Sparsity). The *sparsity* of a cut $(S, V \setminus S)$ is,

$$\psi(S) \doteq \frac{c(S)}{\min\{|S|, |V \setminus S|\}} = \psi(V \setminus S) \in [0, \max_{v \in V} d(v)]. \quad (11.8)$$

The sparsity of a graph is again defined as the smallest sparsity of all cuts.

Remark 11.12. If $|S| \leq |V \setminus S|$, then we can write

$$\psi(S) = \frac{\mathbf{1}_S^\top L \mathbf{1}_S}{\mathbf{1}_S^\top \mathbf{1}_S}. \quad (11.9)$$

Intuitively, sparsity corresponds to a non-normalized variant of conductance. Sometimes it is easier to reason about sparsity than it is to reason about conductance.

Lemma 11.13. *We have for any cut $(S, V \setminus S)$ in a connected unit weight graph that $\psi(S) \geq \phi(S)$.*

Proof. As the graph is connected and has unit weights, $\text{vol}(S) = \sum_{v \in S} d(v) \geq |S|$. \square

An alternative version of Cheeger's inequality relates the second eigenvalue of L (not N !) to the sparsity of the graph.

Fact 11.14 (Cheeger's inequality for sparsity). *We have for a graph G and its Laplacian matrix L ,*

$$\frac{\lambda_2(L)}{2} \leq \psi(G) \leq \sqrt{2\lambda_2(L) \max_{v \in V} d(v)}. \quad (11.10)$$

Effective Resistance

Definition 12.1 (Effective resistance). The *effective resistance*,

$$R_{\text{eff}}(a, b) \doteq \min_{\substack{f \in \mathbb{R}^{|E|} \\ Bf = \mathbf{1}_b - \mathbf{1}_a}} \mathcal{E}(f) = \min_{\substack{f \in \mathbb{R}^{|E|} \\ Bf = \mathbf{1}_b - \mathbf{1}_a}} f^\top R f, \quad (12.1)$$

is the minimum electrical energy required to route one unit of flow from a to b .

Remark 12.2. Per definition of electrical energy, routing F units of flow from a to b costs $F^2 R_{\text{eff}}(a, b)$.

Let us first consider a few examples.

Example 12.3. For the graph of fig. 12.1, $R_{\text{eff}}(1, k+1) = \sum_{i=1}^k r(i)$.

Proof sketch. For the flow to be 1, by Ohm's law, the voltage difference across edge i must be $r(i)$. □

Example 12.4. For the graph of fig. 12.2, $R_{\text{eff}}(1, 2) = \frac{1}{\sum_{i=1}^k 1/r(i)}$.

Proof sketch. For the flow to be 1, by Ohm's law, we must have,

$$1 = \sum_{i=1}^k \frac{\tilde{x}(\{1, 2\})}{r(i)},$$

where $\tilde{x}(\{1, 2\})$ is the voltage difference between vertices 1 and 2. Note that $R_{\text{eff}}(1, 2) = \tilde{x}(\{1, 2\})$. □

Lemma 12.5. $R_{\text{eff}}(a, b) = \left\| L^{+1/2}(\mathbf{1}_b - \mathbf{1}_a) \right\|_2^2$.

Proof. As the electrical flow \tilde{f} is energy-minimizing, we have that $R_{\text{eff}}(a, b) = \tilde{f}^\top R \tilde{f}$. Recall that by Ohm's law this flow corresponds to

TBD

Figure 12.1: Sequential resistors.

TBD

Figure 12.2: Parallel resistors.

voltages \tilde{x} solving $L\tilde{x} = \mathbf{1}_b - \mathbf{1}_a$, that is, $\tilde{x} = L^+(\mathbf{1}_b - \mathbf{1}_a)$. We obtain,

$$\begin{aligned} R_{\text{eff}}(a, b) &= \tilde{f}^\top R \tilde{f} = \tilde{x}^\top L \tilde{x} = (\mathbf{1}_b - \mathbf{1}_a)^\top L^+ L L^+ (\mathbf{1}_b - \mathbf{1}_a) \\ &= (\mathbf{1}_b - \mathbf{1}_a)^\top L^+ (\mathbf{1}_b - \mathbf{1}_a) && \text{using that } \mathbf{1}_b - \mathbf{1}_a \perp \mathbf{1} \\ &= \left\| L^{+1/2} (\mathbf{1}_b - \mathbf{1}_a) \right\|_2^2. \quad \square \end{aligned}$$

Lemma 12.6. *If G is a ϕ -expander, then*

$$R_{\text{eff}}(a, b) \leq 2\phi^{-2} \left(\frac{1}{d(b)} + \frac{1}{d(a)} \right). \quad (12.2)$$

Proof. By Cheeger's inequality,

$$\phi \leq \phi(G) \leq \sqrt{2\lambda_2(N)} \implies \frac{\phi^2}{2} \leq \lambda_2(N).$$

By Courant-Fischer, we have that for any $y \perp \ker N$,

$$\frac{\phi^2}{2} \leq \lambda_2(N) \leq \frac{y^\top N y}{y^\top y} \implies \frac{\phi^2}{2} y^\top y \leq y^\top N y.$$

Equivalently,

$$\frac{\phi^2}{2} \Pi_N \preceq N,$$

using that $\Pi_N v = v$ for $v \perp \ker N$, and $\Pi_N v = 0$ if $v \in \ker N$

where Π_N is the projection orthogonal to the kernel of N . From this we conclude that

$$2\phi^{-2} \Pi_N = 2\phi^{-2} \Pi_N^+ \succeq N^+,$$

using $\Pi_N^+ = \Pi_N$

as $A \succeq B$ implies $A^+ \preceq B^+$ when $\ker A = \ker B$. By eq. (2.18),

$$N^+ = (D^{-1/2} L D^{-1/2})^+ = \Pi_N D^{1/2} L^+ D^{1/2} \Pi_N \quad (12.3)$$

Therefore, for any $y \perp \ker N$,

$$2\phi^{-2} y^\top y \geq y^\top N^+ y = y^\top D^{1/2} L^+ D^{1/2} y.$$

Substituting $z \doteq D^{-1/2} y$, we obtain,

$$2\phi^{-2} z^\top D^{-1} z \geq z^\top L^+ z.$$

Finally, observe that for $z \doteq \mathbf{1}_b - \mathbf{1}_a$, we have that $y = D^{1/2}(\mathbf{1}_b - \mathbf{1}_a) \perp \ker N$ as $\mathbf{1}_b - \mathbf{1}_a \perp \ker L$ and therefore,¹

$$R_{\text{eff}}(a, b) = z^\top L^+ z \leq 2\phi^{-2} z^\top D^{-1} z = 2\phi^{-2} \left(\frac{1}{d(b)} + \frac{1}{d(a)} \right). \quad \square$$

¹ We have for the kernel of the normalized Laplacian matrix, $N = D^{-1/2} L D^{-1/2}$, that $\ker N = D^{1/2} \ker L = \text{span}\{D^{1/2} \mathbf{1}\}$.

Lemma 12.7. $\mathbb{E}[C_{a,b}] = \|d\|_1 R_{\text{eff}}(a, b).$

Proof. Recall that $\mathbb{E}[C_{a,b}] = (\mathbf{1}_a - \mathbf{1}_b)^\top \tilde{\mathbf{x}}$ for a solution $\tilde{\mathbf{x}}$ to $L\tilde{\mathbf{x}} = \|\mathbf{d}\|_1 (\mathbf{1}_a - \mathbf{1}_b)$, that is, $\tilde{\mathbf{x}} = \|\mathbf{d}\|_1 L^+(\mathbf{1}_a - \mathbf{1}_b)$. Now, observe that,

$$R_{\text{eff}}(b, a) = (\mathbf{1}_a - \mathbf{1}_b)^\top L^+(\mathbf{1}_a - \mathbf{1}_b) = \frac{1}{\|\mathbf{d}\|_1} (\mathbf{1}_a - \mathbf{1}_b)^\top \tilde{\mathbf{x}}.$$

Thus, $\mathbb{E}[C_{a,b}] = \|\mathbf{d}\|_1 R_{\text{eff}}(b, a)$. Using symmetry of the commute time, $\mathbb{E}[C_{a,b}] = \mathbb{E}[C_{b,a}] = \|\mathbf{d}\|_1 R_{\text{eff}}(a, b)$. \square

Corollary 12.8. *Effective resistance is symmetric.*

Remark 12.9. For an edge $e = \{a, b\} \in E$, we write,

$$R_{\text{eff}}(e) \doteq R_{\text{eff}}(a, b) = R_{\text{eff}}(b, a). \quad (12.4)$$

12.1 Effective Resistance as a Metric

Before showing that effective resistance is a metric on the set of vertices, we consider the following lemma. We will write,

$$\tilde{\mathbf{x}}_{a,b} \doteq L^+(\mathbf{1}_b - \mathbf{1}_a), \quad (12.5)$$

for the electrical voltages required to route one unit of current from a to b .

Lemma 12.10. *If $\tilde{\mathbf{x}}_{a,b}$ is a solution to $L\tilde{\mathbf{x}}_{a,b} = \mathbf{1}_b - \mathbf{1}_a$, then we have for all $c \in V$ that $\tilde{\mathbf{x}}_{a,b}(b) \geq \tilde{\mathbf{x}}_{a,b}(c) \geq \tilde{\mathbf{x}}_{a,b}(a)$.*

Proof sketch. Consider any $c \in V \setminus \{a, b\}$. Then, $(L\tilde{\mathbf{x}}_{a,b})(c) = 0$. Thus,

$$\left(\sum_{v \sim c} w(\{v, c\}) \right) \tilde{\mathbf{x}}_{a,b}(c) - \left(\sum_{v \sim c} w(\{v, c\}) \tilde{\mathbf{x}}_{a,b}(v) \right) = 0.$$

So, we have,

$$\tilde{\mathbf{x}}_{a,b}(c) = \frac{\sum_{v \sim c} w(\{v, c\}) \tilde{\mathbf{x}}_{a,b}(v)}{\sum_{v \sim c} w(\{v, c\})}.$$

In words, the electrical voltage of c is a weighted average of the voltages of its neighbors. It follows that the voltages of a and b take the largest absolute values. \square

Definition 12.11 (Metric). A *metric* on a set S is a function $d : S \times S \rightarrow \mathbb{R}$ such that for any $a, b, c \in S$,

1. $d(a, b) = 0 \iff a = b$;
2. $d(a, b) \geq 0$;
3. $d(a, b) = d(b, a)$; and
4. $d(a, b) \leq d(a, c) + d(c, b)$.

Lemma 12.12. *Effective resistance is a metric on V .*

Proof. It is easy to check that properties (1) and (2) are satisfied. We have that property (3) is satisfied by corollary 12.8.

Let us therefore consider property (4), the triangle inequality. We have,

$$\tilde{x}_{a,b} = L^+(\mathbf{1}_b - \mathbf{1}_a) = L^+(\mathbf{1}_c - \mathbf{1}_a + \mathbf{1}_b - \mathbf{1}_c) = \tilde{x}_{a,c} + \tilde{x}_{c,b},$$

This we can use to rephrase the effective resistance,

$$\begin{aligned} R_{\text{eff}}(a, b) &= (\mathbf{1}_b - \mathbf{1}_a)^\top \tilde{x}_{a,b} = (\mathbf{1}_b - \mathbf{1}_a)^\top (\tilde{x}_{a,c} + \tilde{x}_{c,b}) \\ &= \tilde{x}_{a,c}(b) - \tilde{x}_{a,c}(a) + \tilde{x}_{c,b}(b) - \tilde{x}_{c,b}(a) \\ &\leq \tilde{x}_{a,c}(c) - \tilde{x}_{a,c}(a) + \tilde{x}_{c,b}(b) - \tilde{x}_{c,b}(c) && \text{using lemma 12.10} \\ &= (\mathbf{1}_c - \mathbf{1}_a)^\top \tilde{x}_{a,c} + (\mathbf{1}_b - \mathbf{1}_c)^\top \tilde{x}_{c,b} \\ &= R_{\text{eff}}(a, c) + R_{\text{eff}}(c, b). \quad \square \end{aligned}$$

Spectral Graph Sparsification

Many combinatorial graph algorithms perform better on sparse graphs. In this chapter, we will see that for any dense graph, we can find a sparse graph with approximately the same Laplacian matrix as measured by quadratic forms.

Definition 13.1 (Matrix approximation). Given $A, B \in S_+^n$ and $\epsilon > 0$, we say,

$$A \approx_\epsilon B \quad \text{iff} \quad \frac{1}{1+\epsilon}A \preceq B(1+\epsilon)A. \quad (13.1)$$

Given some graph $G = (V, E, w)$, our goal is to find a graph $\tilde{G} = (V, \tilde{E}, \tilde{w})$ such that $|\tilde{E}| \ll |E|$ and $L_G \approx_\epsilon L_{\tilde{G}}$. We will write $L \doteq L_G$ and $\tilde{L} \doteq L_{\tilde{G}}$.

Lemma 13.2. *If $L \approx_\epsilon \tilde{L}$, then for any cut $(S, V \setminus S)$,*

$$\frac{1}{1+\epsilon}c_G(S) \leq c_{\tilde{G}}(S) \leq (1+\epsilon)c_G(S). \quad (13.2)$$

Proof. Recall that $c_G(S) = \mathbf{1}_S^\top L \mathbf{1}_S$. The statement follows immediately by comparing the quadratic forms. \square

Theorem 13.3 (Spectral graph approximation by sampling). *For any $\epsilon, \delta \in (0, 1)$ and sampling probabilities,*

$$p_e \doteq \min\{\alpha w(e) R_{\text{eff}}(e), 1\}, \quad (13.3)$$

of each edge $e \in E$ for some scaling parameter $\alpha > 0$ such that if $e \in \tilde{E}$ with probability p_e and $\tilde{w}(e) \doteq w(e)/p_e$, then with probability at least $1 - \delta$ the graph $\tilde{G} = (V, \tilde{E}, \tilde{w})$ satisfies,¹

$$L \approx_\epsilon \tilde{L} \quad \text{and} \quad |\tilde{E}| = \mathcal{O}\left(\frac{n}{\epsilon} \log\left(\frac{n}{\delta}\right)\right).$$

¹ Daniel A Spielman and Nikhil Srivastava. Graph sparsification by effective resistances. *SIAM Journal on Computing*, 40(6):1913–1926, 2011

Proof. TBD

□

Remark 13.4. It can be shown that to compute the sampling probabilities p_e , it is sufficient to solve Laplacian linear systems to find sufficiently good approximations to the effective resistances.

Solving Laplacian Linear Systems

We have seen by now that many problems can be reduced to solving a Laplacian linear system $Lx = d$ where $d \in \text{im } L = (\ker L)^\perp$. Recall that we can obtain a Cholesky decomposition $L = \mathcal{L}\mathcal{L}^\top$ where \mathcal{L} is lower-triangular and positive semi-definite. If L (and therefore \mathcal{L}) were invertible, then we have seen that the linear system can be solved in time $\mathcal{O}(\text{nnz } \mathcal{L})$. However, we know that L is not invertible as $\ker L \neq \{0\}$. It turns out that we can still use of a Cholesky decomposition in solving Laplacian linear systems because we have a simple characterization of the kernel of L . This we will discuss first, then we discuss how to efficiently compute the Cholesky decomposition.

Theorem 14.1. *We can solve the linear system $Lx = d$ where $d \perp \mathbf{1}$ in time $\mathcal{O}(n^3)$ by first computing a Cholesky decomposition of L using Gaussian elimination and then applying the pseudoinverse L^+ to d .*

Finally, we will see that we can find approximate solutions in almost linear time.

14.1 Dealing with pseudoinverses

A natural approach is to characterize the pseudoinverse L^+ in terms of the lower triangular matrix \mathcal{L} .

Lemma 14.2. *Given a factorization $L = \mathcal{L}\mathcal{L}^\top$ where \mathcal{L} is lower triangular and all diagonal entries are strictly non-zero except that $\mathcal{L}(n, n) = 0$, we can apply L^+ in time $\mathcal{O}(n)$.*

Proof. Consider the matrix $\hat{\mathcal{L}}$, which is identical to \mathcal{L} except that $\hat{\mathcal{L}}(n, n) = 1$. Let \mathcal{D} be the diagonal matrix with $\mathcal{D}(i, i) = 1$ for $i < n$ and $\mathcal{D}(n, n) = 0$. Then, $\mathcal{L}\mathcal{L}^\top = \hat{\mathcal{L}}\mathcal{D}\hat{\mathcal{L}}^\top$ and $\hat{\mathcal{L}}$ is invertible. By

claim 2.35, we have

$$\mathbf{L}^+ = \Pi_L (\hat{\mathcal{L}}^\top)^{-1} \mathcal{D}^+ \hat{\mathcal{L}}^{-1} \Pi_L,$$

where Π_L is the orthogonal projection to the kernel of L . Note that $\Pi_L = I - \frac{1}{n} \mathbf{1}\mathbf{1}^\top$, as this satisfies $\Pi_L v = v$ for $v \perp \mathbf{1}$ and $\Pi_L v = \mathbf{0}$ for $v \in \text{span}\{\mathbf{1}\}$. We also have $\mathcal{D}^+ = \mathcal{D}$.

Finally, observe that Π_L and \mathcal{D} can be applied in time $\mathcal{O}(n)$; and by lemma 2.36, we can apply $\hat{\mathcal{L}}^{-1}$ and $(\hat{\mathcal{L}}^\top)^{-1}$ in time $\mathcal{O}(\text{nnz } \mathcal{L})$. \square

14.2 Computing the Cholesky Decomposition

Theorem 14.3 (Cholesky decomposition on graph Laplacians). *Using Gaussian elimination, we can compute in $\mathcal{O}(n^3)$ time a factorization $L = \mathcal{L}\mathcal{L}^\top$, where \mathcal{L} is lower triangular and has positive diagonal entries except that $\mathcal{L}(n, n) = 0$.*

Proof. Let $L^{(0)} \doteq L$. For $i \in [n-1]$, we define,

$$l_i \doteq \frac{1}{\sqrt{L^{(i-1)}(i, i)}} L^{(i-1)}(:, i) \quad \text{and} \\ L^{(i)} \doteq L^{(i-1)} - l_i l_i^\top.$$

Claim 14.4. *Fix some $i < n$. Let $U \doteq \{i+1, \dots, n\}$. Then, $L^{(i)}(i, j) = 0$ if $i \notin U$ or $j \notin U$; and $L^{(i)}$ is a graph Laplacian matrix of a connected graph on the vertex set U .*

It follows that $L^{(n-1)} = \mathbf{0}_{n \times n}$ because the only graph on a single vertex is the empty graph. From this, we have that $L = \sum_{i=1}^{n-1} l_i l_i^\top$, so $\mathcal{L} = [l_1 \cdots l_{n-1} \mathbf{0}]$, provided that l_i is well-defined, i.e., $L^{(i-1)}(i, i) \neq 0$ for all $i < n$. But this follows immediately from the claim, as the diagonal of the Laplacian matrix of a connected graph on more than one vertex must be strictly positive (as the degrees must be non-zero).

In each iteration, we compute $L^{(i)}$ in time $\mathcal{O}(n^2)$ and we proceed for $\mathcal{O}(n)$ iterations. \square

Proof sketch of claim 14.4. We focus on the first elimination, the remaining are similar. We write,

$$L^{(0)} = L \doteq \begin{bmatrix} w & -\mathbf{a}^\top \\ -\mathbf{a} & \text{diag}(\mathbf{a}) + L_{-1} \end{bmatrix},$$

where L_{-1} is defined to make the equality work. We have that,

$$l_1 = \frac{1}{\sqrt{w}} \begin{bmatrix} w \\ -\mathbf{a} \end{bmatrix} \quad \text{and} \quad l_1 l_1^\top = \begin{bmatrix} w & -\mathbf{a}^\top \\ -\mathbf{a} & \frac{1}{w} \mathbf{a} \mathbf{a}^\top \end{bmatrix}$$

Therefore,

$$\mathbf{L}^{(1)} = \mathbf{L}^{(0)} - \mathbf{l}_1 \mathbf{l}_1^\top = \begin{bmatrix} 0 & \mathbf{0}^\top \\ \mathbf{0} & \mathbf{S}^{(1)} \end{bmatrix},$$

where $\mathbf{S}^{(1)} \doteq \mathbf{L}_{+1} + \mathbf{L}_{-1}$ and $\mathbf{L}_{+1} \doteq \text{diag}(\mathbf{a}) - \frac{1}{w} \mathbf{a} \mathbf{a}^\top$. $\mathbf{S}^{(1)}$ is also called the *Schur complement*. We can also phrase it differently as,

$$\mathbf{l}_1 \mathbf{l}_1^\top = \begin{bmatrix} w & -\mathbf{a}^\top \\ -\mathbf{a} & \text{diag}(\mathbf{a}) \end{bmatrix} - \begin{bmatrix} 0 & \mathbf{0}^\top \\ \mathbf{0} & \text{diag}(\mathbf{a}) - \frac{1}{w} \mathbf{a} \mathbf{a}^\top \end{bmatrix} \quad (14.1)$$

$$\doteq \text{STAR}(1, \mathbf{L}) - \text{CLIQUE}(1, \mathbf{L}). \quad (14.2)$$

It remains to show that $\mathbf{S}^{(1)}$ is the Laplacian matrix of a connected graph on the vertex set $\{2, \dots, n\}$.

Observe that the sum of two Laplacian matrices is again a graph Laplacian. Therefore, \mathbf{L}_{-1} is a graph Laplacian by definition. It is easy to check that by the characterization of exercise 1.12, \mathbf{L}_{+1} also is a graph Laplacian, which represents a clique formed by the neighbors of vertex 1. Hence, $\mathbf{S}^{(1)} = \mathbf{L}_{+1} + \mathbf{L}_{-1}$ is a graph Laplacian.

It remains to show that the underlying graph is connected. Consider any $i, j \in V \setminus \{1\}$. There exists an i - j path P in the graph of \mathbf{L} . If P does not use vertex 1, then P is a path in the graph of \mathbf{L}_{-1} . If P does use vertex 1, it does so using some edges $(u, 1)$ and $(1, v)$. Replacing the two edges with the edge (u, v) , which appears in \mathbf{L}_{+1} as $\mathbf{L}_{+1}(u, v) < 0$, yields a path in the graph of $\mathbf{S}^{(1)}$. \square

14.3 Approximate Almost Linear-Time Solvers

Definition 14.5 (Approximate solution to linear system). We say that $\tilde{\mathbf{x}}$ is an ϵ -approximate solution to the linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ iff

$$\|\tilde{\mathbf{x}} - \mathbf{x}^*\|_A^2 \leq \epsilon \|\mathbf{x}^*\|_A^2, \quad (14.3)$$

where \mathbf{x}^* is a solution and $\|\mathbf{y}\|_A \doteq \sqrt{\mathbf{y}^\top \mathbf{A} \mathbf{y}}$ denotes the *Mahalanobis norm* with respect to \mathbf{A} .

The fast algorithm has two main steps.

Theorem 14.6 (Approximate Cholesky decomposition on graph Laplacians). *We can find $\mathcal{L} \mathcal{L}^\top \approx_{1/2} \mathbf{L}$ such that \mathcal{L} is lower triangular and $\text{nnz } \mathcal{L} = \mathcal{O}(m \log^3 n)$, with probability at least $1 - 3/n^5$ in time $\mathcal{O}(m \log^3 n)$.¹*

Proof. TBD \square

TBD

Figure 14.1: $\mathbf{S}^{(1)}$ is the Laplacian matrix of the graph, where the first vertex was removed and all of its neighbors are made to be in a clique.

¹ Rasmus Kyng and Sushant Sachdeva. Approximate gaussian elimination for laplacians-fast, sparse, and simple. In 2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS), pages 573–582. IEEE, 2016

Here, the main idea is to replace the clique of the neighbors of the eliminated vertex in each round of Gaussian elimination by a sparse approximation.

Theorem 14.7. *We can find an ϵ -approximate solution \tilde{x} to $Lx = d$, using an algorithm that takes $\mathcal{O}\left(m \log^3 n \log(1/\epsilon)\right)$ time and succeeds with probability at least $1 - 1/n^{10}$.*

Proof. TBD

□

PART IV

Combinatorial Graph Algorithms

Algorithms for Maximum Flow

In this chapter, we study classical combinatorial algorithms for the maximum flow problem. Many of the algorithms can be adapted to solve the more general minimum-cost flow problem. In recent years, significant progress was made by using convex optimization and interior point methods. To begin with, we will derive a combinatorial proof for the strong duality of maximum flow and minimum cut.¹

¹ We have already seen a proof of this in section 8.5.

We consider directed graphs $G = (V, E, c)$ with edge capacities $c \in \mathbb{R}^{|E|}$. Recall that a *flow* is a vector $f \in \mathbb{R}^{|E|}$ and f routes demands d iff $Bf = d$. We say that f is *feasible* iff $\mathbf{0} \leq f \leq c$, where $\mathbf{0} \leq f$ ensures that the flow respects edge directions and $f \leq c$ ensures that the flow respects edge capacities.

We call a flow f that routes demands $F(\mathbf{1}_t - \mathbf{1}_s)$ for some $F \in \mathbb{R}$ and vertices $s, t \in V$, so F units of flow from s to t , an *s-t flow* with *value* $\text{val}(f) = F$. We say that f is *optimal* iff there is no feasible s - t flow f' such that $\text{val}(f') > \text{val}(f)$.

We can decompose any s - t flow into two kinds of flows: path flows and cycle flows.

Definition 15.1 (Path flow). An *s-t path flow* f is an s - t flow with $\text{val}(f) = \alpha$ for some $\alpha > 0$ that can be expressed as,

$$f = \alpha \sum_{e \in P} \mathbf{1}_e, \quad (15.1)$$

for some s - t path P .

Definition 15.2 (Cycle flow). A *cycle flow* f is a flow routing demands $\mathbf{0}$ that can be expressed as,

$$f = \alpha \sum_{e \in C} \mathbf{1}_e, \quad (15.2)$$

for some $\alpha > 0$, and cycle C .

TBD
Figure 15.1: Example of an s - t path flow.

TBD

Figure 15.2: Example of a cycle flow.

Lemma 15.3 (Path-cycle decomposition). *Any s - t flow $f \geq 0$ can be decomposed into $k \leq m$ s - t path flows and l cycle flows.*

Proof. TBD □

Lemma 15.4. *There exists an optimal flow with a path-cycle decomposition that has only paths and no cycles.*

Proof. TBD □

Lemma 15.5. *There exists an s - t flow $f \geq 0$ iff there exists a directed s - t path.*

Proof. TBD □

Recall that a *cut* is a proper subset of vertices $\emptyset \subset S \subset V$. An s - t cut is a cut $(S, V \setminus S)$ separating s and t , i.e., $s \in S, t \in V \setminus S$. We say that the capacity of a cut is the sum of capacities of crossing edges,

$$\text{cap}(S) \doteq \sum_{\substack{\{a,b\} \in E \\ a \in S, b \in V \setminus S}} c(\{a,b\}). \quad (15.3)$$

Theorem 15.6 (Weak duality of maximum flow/minimum cut). *For any feasible s - t flow f and any s - t cut $(S, V \setminus S)$,*

$$\text{val}(f) \leq \text{cap}(S). \quad (15.4)$$

Proof. Let $(S, V \setminus S)$ be any s - t cut. Suppose for a contradiction that $\text{val}(f) > \text{cap}(S)$ for some flow f . But this contradicts feasibility of f because the crossing edges of the cut form a bottleneck. □

An important concept in the analysis of flow algorithms is the so-called residual graph.

Definition 15.7 (Residual graph). The *residual graph* G_f of some s - t flow $f \geq 0$ is the graph G with edge capacities $[-f, c - f]$. That is, we say that a flow \hat{f} is *feasible* in the residual graph iff $-f \leq \hat{f} \leq c - f$.

Intuitively, sending positive flow $c - f$ along an edge in G_f corresponds to sending the maximum additional flow without violating the capacity constraint within G , whereas sending negative flow $-f$ along an edge in G_f corresponds to “undoing” the flow that was sent along this edge by f . This simple argument shows that if \hat{f} is feasible in G_f , then $f + \hat{f}$ is feasible in G .

We call an s - t flow \hat{f} in G_f an *augmenting flow* of f .

Lemma 15.8 (Flow optimality condition). *A feasible s - t flow f in G is optimal iff there is no s - t path in G_f , or equivalently, iff there is no f -augmenting flow.*

Proof. TBD □

Theorem 15.9 (Strong duality of maximum flow/minimum cut). *We have that,*

$$\max_{\substack{F \in \mathbb{R}, f \in \mathbb{R}^{|E|} \\ Bf = F(\mathbf{1}_t - \mathbf{1}_s)}} F = \min_{s-t \text{ cut } (S, V \setminus S)} \text{cap}(S). \quad (15.5)$$

Proof. By weak duality of maximum flow and minimum cut, we have the direction \leq . For the direction \geq , let f^* be an optimal flow and consider the cut,

$$S \doteq \{v \in V \mid \text{there exists an } s\text{-}v \text{ path in } G_{f^*}\}.$$

We make two observations.

1. By definition, there are no edges from S to $V \setminus S$ in G_{f^*} , that is, f^* saturates² all crossing edges of the cut.
2. By definition, f^* routes no flow from $V \setminus S$ to S .

² That is, f^* sends flow equal to the capacity of the edge.

This implies that $\text{val}(f^*) \geq \text{cap}(S)$ for this cut $(S, V \setminus S)$. □

15.1 The Ford-Fulkerson Algorithm

Our prior discussion gives rise to a very natural algorithm.

Algorithm 15.10: FORDFULKERSON(G)

```

1  $f \leftarrow \mathbf{0}$ 
2 while there exists any  $s$ - $t$  path flow  $\hat{f}$  in  $G_f$  do
3    $f \leftarrow f + \hat{f}$ 
4 return  $f$ 
```

Given a feasible flow f , we can find an f -augmenting flow, or determine that none exists, in time $\mathcal{O}(m)$ using breadth-first search or depth-first search.

Theorem 15.11 (Ford-Fulkerson). *If capacities are integral, FORDFULKERSON converges to an optimal flow f^* in $\text{val}(f^*)$ iterations.*

Proof. Observe that the initial flow $\mathbf{0}$ is trivially feasible. In each iteration, we add the augmenting flow \hat{f} with $\text{val}(\hat{f}) > 0$, and due to

the integral capacities, $\text{val}(\hat{f}) \geq 1$. Therefore, the flow value $\text{val}(f)$ increases in each iteration by at least one. \square

Improving Ford-Fulkerson

It turns out that if we always choose the shortest augmenting path, we converge in time $\mathcal{O}(nm^2)$. This is known as the *Edmonds-Karp algorithm*.

We can do still better, by choosing that path with the maximum bottleneck capacity. That is, we choose,

$$P^* = \arg \max_{\text{augmenting paths } P} \min_{e \in P} c(e). \quad (15.6)$$

Within the framework of Ford-Fulkerson this corresponds to the augmenting path that allows us to route the most additional flow.

Theorem 15.12. *FORDFULKERSON, where in each iteration we choose the augmenting path with maximum bottleneck capacity, converges in time $\mathcal{O}(m^2 \log mU)$.*

Proof. We can find P^* using a binary search on $[1, U]$, where $U \doteq \max_e c(e)$, by removing all edges with absolute capacity in G_f below the current threshold and testing if an s - t path in G_f exists: if it does, we increase the threshold; if it does not, we decrease the threshold. This procedure takes $\mathcal{O}(m \log U)$ time. If we only consider the occurring capacities, the runtime improves to $\mathcal{O}(m \log m)$.

Suppose \hat{F} is the flow left in G_f . By the path decomposition lemma, this flow can be decomposed into at most m path flows (the “best” of which is P^*) and P^* must route at least the average amount of flow. Hence, P^* routes at least \hat{F}/m units of flow. Thus, the algorithm converged if,

$$\left(1 - \frac{1}{m}\right)^T F^* < 1,$$

where T is the number of augmentations and F^* is the value of an optimal flow. So, some $T = \mathcal{O}(m \log F^*)$ is sufficient.

Overall, we get,

$$\mathcal{O}(m \log m \cdot T) = \mathcal{O}(m^2 \log(m + F^*)) = \mathcal{O}(m^2 \log mU) \quad \square \quad \text{using } F^* \leq mU$$

15.2 Dinitz's Algorithm

We will now look at Dinitz's algorithm, which also belongs to the family of Ford-Fulkerson algorithms. That is, we still start with the

initial feasible flow $f \leftarrow \mathbf{0}$ and iteratively improve this flow by finding an augmenting flow in G_f .

In Dinitz's algorithm, our strategy to finding an augmenting path is to find so-called *blocking flows* in G_f , which in each iteration “block” one shortest s - t path in G_f . Dinitz's algorithm runs in $\mathcal{O}(n^2 + nm \log^2 n)$ time in general graphs and in time $\mathcal{O}(\min\{m^{3/2}, mn^{2/3}\})$ in unit capacity graphs.

Remark 15.13. It can also be shown that Dinitz's algorithm converges in time $\mathcal{O}(m\sqrt{n})$ in bipartite matching graphs, but we will not show this here.

Definition 15.14 (Level). The *level* $\ell(v)$ of a vertex $v \in V$ in G_f is the length (i.e., number of edges) of the shortest s - v path in G_f .

We call an edge $e = (u, v)$ in G_f *admissible* iff $\ell(v) = \ell(u) + 1$. Intuitively, e is admissible iff it is part of a shortest s - t path.

The *level graph* L_f of G_f is the subgraph of only admissible edges.

Definition 15.15 (Blocking flow). A *blocking flow* in G_f is a flow \hat{f} such that

1. \hat{f} is feasible;
2. \hat{f} uses only admissible edges; and
3. for every s - t path in the level graph, \hat{f} saturates at least one edge.

Remark 15.16. By definition, a blocking flow in G_f is a f -augmenting flow.

TBD

Figure 15.3: Illustration of admissible edges.

TBD

Figure 15.4: Examples of blocking flows.

Algorithm 15.17: DINITZ(G)

```

1  $f \leftarrow \mathbf{0}$ 
2 repeat
3    $\hat{f} \leftarrow$  blocking flow in  $G_f$ 
4    $f \leftarrow f + \hat{f}$ 
5 until  $G_f$  is disconnected
6 return  $f$ 
```

Lemma 15.18. Let f be a feasible flow, \hat{f} a blocking flow in G_f , and let $f' \doteq f + \hat{f}$. Then $\ell_{G_{f'}}(t) > \ell_{G_f}(t)$.

Proof. TBD

□

Theorem 15.19 (Dinitz). Dinitz's algorithm converges in $\mathcal{O}(n)$ iterations.

Proof. By the previous lemma, $\ell(t)$ increases in every iteration by at least one. As path can contain each vertex at most once, the level of any vertex can never be larger than n . \square

Finding Blocking Flows

A naïve approach is to repeatedly use depth-first search to find an unsaturated s - t path in the level graph L_f . Whenever we find such a path, we route the maximum possible flow along this path, saturating at least one of its edges.

Algorithm 15.20: $\text{FINDBLOCKINGFLOW}(L_f)$

```

1  $\hat{f} \leftarrow 0$ 
2 while there exists an  $s$ - $t$  path  $P$  in  $L_f$  do
3   Let  $\hat{f}'$  be a flow saturating  $P$ 
4    $\hat{f} \leftarrow \hat{f} + \hat{f}'$ 
5   Remove all edges saturated by  $\hat{f}'$  from  $L_f$ 
6 return  $\hat{f}$ 

```

Lemma 15.21. FINDBLOCKINGFLOW returns a blocking flow in $\mathcal{O}(mn)$ time.

Remark 15.22. This can be improved to $\mathcal{O}(m \log^2 n + n)$ by using link-cut trees, which we will discuss in the next chapter.

Proof. TBD \square

Unit Capacity Graphs

Lemma 15.23. In unit capacity graphs, FINDBLOCKINGFLOW returns a blocking flow in $\mathcal{O}(m)$ time.

Proof. TBD \square

Theorem 15.24. In unit capacity graphs, Dinic's algorithm terminates in $\mathcal{O}(\min\{m^{1/2}, n^{2/3}\})$ time.

Proof. TBD \square

15.3 The Push-Relabel Algorithm

15.4 Outlook

We have seen two approaches to solving maximum flow problems: Ford-Fulkerson maintains a feasible flow and augments this flow until it is optimal. In contrast, Push-Relabel maintains that there is no augmenting path and terminates when the flow is feasible.

Currently, the best known algorithm for real-valued capacities is due to Orlin and takes $\mathcal{O}(mn)$ time.³

Geometrically speaking, Ford-Fulkerson is analogous to a simplex method and Push-Relabel is analogous to an exterior-point method. In recent years, interior-point methods were used to find efficient algorithms when capacities are integral, the best known algorithm taking $\mathcal{O}(m^{1+o(1)} \log U)$ time.⁴

³ James B Orlin. Max flows in $\mathcal{O}(nm)$ time, or better. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 765–774, 2013

⁴ Li Chen, Rasmus Kyng, Yang P Liu, Richard Peng, Maximilian Probst Gutenberg, and Sushant Sachdeva. Maximum flow and minimum-cost flow in almost-linear time. *arXiv preprint arXiv:2203.00671*, 2022

16

Link-Cut Trees

17

Finding Expanders using Maximum Flow

17.1 *Graph Embedding*

18

Distance Oracles

PART V

Further Topics

19

Interior Point Methods for Maximum Flow

A

Solutions

A.1 Part I

Electrical Flows

Solution to exercise 1.12. TBD □

Proof of claim 1.15. We have that $H_c(\mathbf{x}) = L(\mathbf{x}) \succeq \mathbf{0}$. □

Solution to exercise 1.17. Consider any $\mathbf{f} \in \mathbb{R}^{|E|}$ satisfying $\mathbf{B}\mathbf{f} = \mathbf{d}$. We have for any $\mathbf{x} \in \mathbb{R}^{|V|}$,

$$\begin{aligned} \frac{1}{2}\mathbf{f}^\top \mathbf{R}\mathbf{f} &= \frac{1}{2}\mathbf{f}^\top \mathbf{R}\mathbf{f} - \mathbf{x}^\top \underbrace{(\mathbf{B}\mathbf{f} - \mathbf{d})}_0 \\ &\geq \min_{\mathbf{f}' \in \mathbb{R}^{|E|}} \underbrace{\frac{1}{2}\mathbf{f}'^\top \mathbf{R}\mathbf{f}' - \mathbf{x}^\top (\mathbf{B}\mathbf{f}' - \mathbf{d})}_{\doteq g(\mathbf{f}')}. \end{aligned}$$

Note that g is convex. Taking the gradient gives,

$$\nabla_{\mathbf{f}'} g(\mathbf{f}') = \mathbf{R}\mathbf{f}' - \mathbf{B}^\top \mathbf{x}.$$

So $\mathbf{f}' = \mathbf{R}^{-1}\mathbf{B}^\top \mathbf{x}$ is the minimizer of g . We obtain,

$$\frac{1}{2}\mathbf{f}^\top \mathbf{R}\mathbf{f} \geq -\frac{1}{2}\mathbf{x}^\top \mathbf{L}\mathbf{x} + \mathbf{d}^\top \mathbf{x},$$

but $\tilde{\mathbf{f}}^\top \mathbf{R}\tilde{\mathbf{f}} = \tilde{\mathbf{x}}^\top \mathbf{L}\tilde{\mathbf{x}} = \mathbf{d}^\top \tilde{\mathbf{x}}$ for electrical voltages $\tilde{\mathbf{x}}$ and electrical flow $\tilde{\mathbf{f}}$ (see eq. (1.18)). Thus,

$$\frac{1}{2}\mathbf{f}^\top \mathbf{R}\mathbf{f} \geq \frac{1}{2}\tilde{\mathbf{x}}^\top \mathbf{L}\tilde{\mathbf{x}} = \frac{1}{2}\tilde{\mathbf{f}}^\top \mathbf{R}\tilde{\mathbf{f}}.$$

Hence, $\tilde{\mathbf{f}}$ is the minimum electrical energy flow among all flows routing \mathbf{d} . □

Solution to exercise 1.18. TBD □

Solution to exercise 1.19. TBD □

Linear Algebra

Proof of claim 2.1. Because the characteristic polynomial of A is of degree n , it has n complex roots, which are the eigenvalues $\lambda_1, \dots, \lambda_n$ of A . We will first prove that the λ_i are real. Then, we will prove that the corresponding eigenvectors v_i are orthogonal.

1. Let λ be any eigenvalue of A . We denote by $\bar{\lambda}$ the complex conjugate of λ . Clearly, if $\lambda = \bar{\lambda}$, then $\lambda \in \mathbb{R}$. By the definition of the eigenvalue λ with associated eigenvector v , we have,

$$\lambda \bar{v}^\top v = \bar{v}^\top A v.$$

Taking the complex conjugate and transpose of both sides gives,

$$\bar{\lambda} \bar{v}^\top v = \bar{v}^\top \bar{A}^\top v = \bar{v}^\top A v = \lambda \bar{v}^\top v.$$

using that A is real and symmetric,
 $\bar{A}^\top = A$

We have $\lambda = \bar{\lambda}$ as desired.

2. It remains to show that for eigenvalues λ_i, λ_j with associated eigenvectors v_i, v_j and $i \neq j$, we have $v_i^\top v_j = 0$. By the definition of an eigenvalue, we have,

$$\begin{aligned} \lambda_i v_j^\top v_i &= v_j^\top A v_i = (v_i^\top A^\top v_j)^\top \\ &= (v_i^\top A v_j)^\top = \lambda_j (v_i^\top v_j)^\top = \lambda_j v_j^\top v_i. \end{aligned}$$

using that A is symmetric

We get that $v_j^\top v_i = 0$ if $\lambda_i \neq \lambda_j$. □

Proof of claim 2.7. Suppose λ is an eigenvalue of M with associated eigenvector v . Also suppose that $Tv = w$. Then,

$$TMT^{-1}w = TMv = \lambda Tv = w.$$

It is easy to check that the other direction holds too. □

Proof of claim 2.26. TBD □

Proof of claim 2.33. TBD □

Proof of claim 2.35. TBD □

Probability

Theorem A.1 (Jensen's inequality, finite form). *Let $f : S \rightarrow \mathbb{R}$ be a convex function on the convex set $S \subseteq \mathbb{R}^n$. Suppose that $\mathbf{x}_1, \dots, \mathbf{x}_k \in S$ and $\theta_1, \dots, \theta_k \geq 0$ with $\theta_1 + \dots + \theta_k = 1$. Then,*

$$f(\theta_1 \mathbf{x}_1 + \dots + \theta_k \mathbf{x}_k) \leq \theta_1 f(\mathbf{x}_1) + \dots + \theta_k f(\mathbf{x}_k). \quad (\text{A.1})$$

Proof. We prove the statement by induction on k . The base case, $k = 2$, follows trivially from the convexity of f . For the induction step, suppose that the statement holds for some $k \geq 2$. Assume w.l.o.g. that $\theta_{k+1} \in (0, 1)$. We have,

$$\begin{aligned} \sum_{i=1}^{k+1} \theta_i f(\mathbf{x}_i) &= (1 - \theta_{k+1}) \left(\sum_{i=1}^k \frac{\theta_i}{1 - \theta_{k+1}} f(\mathbf{x}_i) \right) + \theta_{k+1} f(\mathbf{x}_{k+1}) \\ &\geq (1 - \theta_{k+1}) f \left(\sum_{i=1}^k \frac{\theta_i}{1 - \theta_{k+1}} \mathbf{x}_i \right) + \theta_{k+1} f(\mathbf{x}_{k+1}) && \text{using the induction hypothesis} \\ &\geq f \left(\sum_{i=1}^{k+1} \theta_i \mathbf{x}_i \right). && \square \quad \text{using convexity of } f \end{aligned}$$

A.2 Part II

Convex Geometry

Proof of theorem 5.3. In our proof, we will use the following two theorems.

Fact A.2 (Bolzano-Weierstrass theorem). *Every bounded sequence in \mathbb{R}^n has a convergent subsequence.*

Fact A.3 (Boundedness theorem). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuous function and $\mathcal{F} \subseteq \mathbb{R}^n$ be non-empty, bounded, and closed. Then f is bounded on \mathcal{F} .*

Let α be the infimum of f over \mathcal{F} , i.e., the largest value for which any $\mathbf{x} \in \mathcal{F}$ satisfies $f(\mathbf{x}) \geq \alpha$. By the boundedness theorem, the infimum exists, as f is lower bounded and the set of lower bounds has a greatest lower bound, α .

Let $\mathcal{F}_k \doteq \{\mathbf{x} \in \mathcal{F} \mid \alpha \leq f(\mathbf{x}) \leq \alpha + 2^{-k}\}$. \mathcal{F}_k cannot be empty, since if it were, then $\alpha + 2^{-k}$ would be a strictly greater lower bound on f than α . For each k , let \mathbf{x}_k be some $\mathbf{x} \in \mathcal{F}_k$. $\{\mathbf{x}_k\}_{k=1}^{\infty}$ is a bounded sequence as $\mathcal{F}_k \subseteq \mathcal{F}$, so by the Bolzano-Weierstrass theorem, there exists a convergent subsequence, $\{\mathbf{y}_k\}_{k=1}^{\infty}$, with limit $\bar{\mathbf{y}}$. Because the

set is closed, $\bar{y} \in \mathcal{F}$. By continuity, $f(\bar{y}) = \lim_{k \rightarrow \infty} f(y_k)$, and by construction, $\lim_{k \rightarrow \infty} f(y_k) = \alpha$.

Thus, the optimal solution is \bar{y} . □

Solution to exercise 5.10. TBD □

Solution to exercise 5.12. TBD □

A.3 Part III

Introduction to Spectral Graph Theory

Solution to exercise 10.9. TBD □

Solution to exercise 10.10. TBD □

Solution to exercise 10.11. TBD □

Solution to exercise 10.12. TBD □

Conductance and Expanders

Solution to exercise 11.7. TBD □

Solution to exercise 11.8. TBD □

Summary of Notation

We follow these general rules:

- uppercase italic for constants N
- lowercase italic for indices i and scalar variables x
- lowercase italic bold for vectors \mathbf{x}
- uppercase italic bold for matrices \mathbf{M}
- uppercase italic for random variables X
- uppercase bold for random vectors \mathbf{X}
- uppercase italic for sets A

\doteq	equality by definition
iff	if and only if
\mathbb{N}	set of natural numbers $\{1, 2, \dots\}$
\mathbb{N}_0	set of natural numbers, including 0, $\mathbb{N} \cup \{0\}$
\mathbb{R}	set of real numbers
\mathbb{C}	set of complex numbers
$[m]$	set of natural numbers from 1 to m , $\{1, 2, \dots, m-1, m\}$
$(a, b]$	real interval between a and b including b but not including a
$f : A \rightarrow B$	function f from elements of set A to elements of set B
$\mathbb{1}\{\textit{predicate}\}$	indicator function ($\mathbb{1}\{\textit{predicate}\} \doteq 1$ if the <i>predicate</i> is true, else 0)
\leftarrow	assignment

LINEAR ALGEBRA

S^n	set of symmetric $n \times n$ matrices
S_+^n	set of symmetric and positive semi-definite $n \times n$ matrices
S_{++}^n	set of symmetric and positive definite $n \times n$ matrices
$A \preceq B$	Loewner order on symmetric matrices, $\forall \mathbf{x} \in \mathbb{R}^n : \mathbf{x}^\top \mathbf{A} \mathbf{x} \leq \mathbf{x}^\top \mathbf{B} \mathbf{x}$
$\mathbf{x} \perp \mathbf{y}$	\mathbf{x} and \mathbf{y} are orthogonal, i.e., $\mathbf{x}^\top \mathbf{y} = 0$
$\mathbf{x} \perp W$	\mathbf{x} is orthogonal to every vector \mathbf{y} in subspace W
W^\perp	orthogonal complement of subspace W , $\{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{x} \perp W\}$
$\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$	smallest subspace containing $\mathbf{x}_1, \dots, \mathbf{x}_n$
$\dim(W)$	number of vectors in a basis of a subspace W
$H(\mathbf{n}, \mu)$	hyperplane with normal \mathbf{n} and threshold μ

$\mathbf{1}_S$	vector such that $\mathbf{1}_S(i) = \mathbb{1}\{i \in S\}$
A^\top	transpose of matrix A
A^{-1}	inverse of matrix A
A^+	pseudoinverse of matrix A
$A^{1/2}$	square root of symmetric positive semi-definite matrix A
Π_A	orthogonal projection to $(\ker A)^\perp$
$\text{nnz } A$	number of non-zero entries of A
$\text{tr } A$	trace of A , $\sum_i A(i, i)$
$\text{diag}_{i \in I}\{a_i\}$	diagonal matrix with elements a_i , indexed according to the set I
$\ker A$	kernel (or null space) of A , $\{\mathbf{x} \in \mathbb{R}^n \mid A\mathbf{x} = \mathbf{0}\}$
$\text{im } A$	image of A , $\text{span}\{A(:, i)\}_{i \in [n]}$
$\lambda_i(A)$	i -th smallest eigenvalue of A
$\ A\ _{\alpha \rightarrow \beta}$	matrix norm of A induced by norms $\ \cdot\ _\alpha$ and $\ \cdot\ _\beta$

PROBABILITY

$\mathbb{P}[X = x]$	probability of a random variable X taking on the value x
$X \sim F$	random variable X follows the distribution F
$x \sim F$	value x is sampled according to distribution F
$X \perp Y$	random variable X is independent of random variable Y
$X \perp Y \mid Z$	random variable X is conditionally independent of random variable Y given random variable Z
$\mathbb{E}[X]$	expected value of random variable X
$\text{Var}[X]$	variance of random variable X
$W \in \mathbb{R}^{ V \times V }$	transition matrix of random walk, $AD^{-1} = I - D^{1/2}ND^{-1/2}$
$\tilde{W} \in \mathbb{R}^{ V \times V }$	transition matrix of lazy random walk, $\frac{I}{2} + \frac{W}{2} = I - \frac{1}{2}D^{1/2}ND^{-1/2}$
$p_t \in \mathbb{R}^{ V }$	probability distribution of a random walk at time t , $W^t p_0$
$H_{a,s}$	hitting time of s starting from a
$h_s \in \mathbb{R}^{ V }$	expected hitting times of s , $h_s(a) = \mathbb{E}[H_{a,s}]$
$C_{a,b}$	commute time between a and b , $H_{a,b} + H_{b,a}$

ANALYSIS

$\nabla f(\mathbf{x}) \in \mathbb{R}^{n \times 1}$	gradient of a function f at a point \mathbf{x} , $\left[\frac{\partial f(\mathbf{x})}{\partial x(1)} \quad \dots \quad \frac{\partial f(\mathbf{x})}{\partial x(n)}\right]^\top$
$[x, y]$	set of convex combinations of \mathbf{x} and \mathbf{y} , $\{\theta\mathbf{x} + (1 - \theta)\mathbf{y} \mid \theta \in [0, 1]\}$
$Df(\mathbf{x})[\mathbf{d}]$	directional derivative of f at \mathbf{x} in direction \mathbf{d} , $\lim_{\lambda \rightarrow 0} \frac{f(\mathbf{x} + \lambda\mathbf{d}) - f(\mathbf{x})}{\lambda}$
$D\mathbf{g}(\mathbf{x}) \in \mathbb{R}^{m \times n}$	Jacobian of vector-valued function $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\left[\frac{\partial \mathbf{g}(\mathbf{x})}{\partial x(1)} \quad \dots \quad \frac{\partial \mathbf{g}(\mathbf{x})}{\partial x(n)}\right]$
$H_f(\mathbf{x}) \in \mathbb{R}^{n \times n}$	Hessian of f , $(D\nabla f(\mathbf{x}))^\top$
$\text{epi}(f)$	epigraph of a function f , $\{(\mathbf{x}, y) \mid f(\mathbf{x}) \leq y\} \subseteq \mathbb{R}^{n+1}$

$S_\alpha(f)$	α -sub-level set of a function f , $\{x \in S \mid f(x) \leq \alpha\}$
$L_\alpha(f)$	α -level set of a function f , $\{x \in S \mid f(x) = \alpha\}$

GRAPHS

V	set of vertices
E	set of edges
n	number of vertices, $ V $
m	number of edges, $ E $
$G[X]$	subgraph of G induced by $X \subseteq V$
$u \sim v$	vertices u and v are adjacent
$\deg(v)$	degree of vertex v
$\mathbf{r} \in \mathbb{R}^{ E }$	resistances
$\mathbf{w} \in \mathbb{R}^{ E }$	weights, $1/r(e)$
$\mathbf{d} \in \mathbb{R}^{ V }$	weighted degrees, $\sum_{\{u,v\} \in E} \mathbf{w}(\{u,v\})$
$\tilde{\mathbf{A}} \in \mathbb{R}^{ V \times V }$	adjacency matrix
$\mathbf{A} \in \mathbb{R}^{ V \times V }$	weighted adjacency matrix
$\mathbf{B} \in \mathbb{R}^{ V \times E }$	incidence matrix
$\mathbf{R} \in \mathbb{R}^{ E \times E }$	diagonal matrix of resistances, $\text{diag}\{\mathbf{r}(e)\}_{e \in E}$
$\mathbf{W} \in \mathbb{R}^{ E \times E }$	diagonal matrix of weights, $\text{diag}\{\mathbf{w}(e)\}_{e \in E}$
$\mathbf{D} \in \mathbb{R}^{ V \times V }$	diagonal matrix of weighted degrees, $\text{diag}\{\mathbf{w}(v)\}_{v \in V}$
$\mathbf{L} \in \mathbb{R}^{ V \times V }$	Laplacian matrix, $\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^\top = \mathbf{B}\mathbf{W}\mathbf{B}^\top = \mathbf{D} - \mathbf{A}$
$\text{vol}(S)$	volume of a set of vertices S , $\sum_{v \in S} \mathbf{d}(v) = \mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S$
$c(S)$	value of a cut S , $\sum_{\{a,b\} \in E: a \in S, b \in V \setminus S} \mathbf{w}(\{a,b\}) = \mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S$
$\phi(S)$	conductance of a cut S , $\frac{c(S)}{\min\{\text{vol}(S), \text{vol}(V \setminus S)\}}$
$\phi(G)$	conductance of a graph G , $\min_{\emptyset \subset S \subset V} \phi(S)$
$\psi(S)$	sparsity of a cut S , $\frac{c(S)}{\min\{ S , V \setminus S \}}$
$\psi(G)$	sparsity of a graph G , $\min_{\emptyset \subset S \subset V} \psi(S)$
K_n	unit weight complete graph on n vertices
P_n	unit weight path graph on n vertices
T_d	unit weight complete binary tree with d levels
$G_{i,j}$	unit weight graph with single edge $\{i,j\}$
$G^{i,j}$	subgraph of G consisting of the shortest i,j path

FLOWS

$\mathbf{x} \in \mathbb{R}^{ V }$	voltages
$\mathbf{x}(e)$	voltage difference of edge $e = \{u,v\}$, $\mathbf{x}(u) - \mathbf{x}(v)$
$\mathbf{f} \in \mathbb{R}^{ E }$	flow

$\mathbf{d} \in \mathbb{R}^{ V }$	demands, modeling the net flow
$\tilde{\mathbf{x}} \in \mathbb{R}^{ V }$	electrical voltages
$\tilde{\mathbf{x}}_{a,b} \in \mathbb{R}^{ V }$	electrical voltages routing demands $\mathbf{1}_b - \mathbf{1}_a$
$\tilde{\mathbf{f}} \in \mathbb{R}^{ E }$	electrical flow
$\mathcal{E}(\tilde{\mathbf{f}}), \mathcal{E}(\tilde{\mathbf{x}}), \mathcal{E}(\mathbf{d})$	electrical energy, $\tilde{\mathbf{f}}^\top \mathbf{B}^\top \tilde{\mathbf{x}} = \tilde{\mathbf{f}}^\top \mathbf{R} \tilde{\mathbf{f}} = \tilde{\mathbf{x}}^\top \mathbf{L} \tilde{\mathbf{x}} = \mathbf{d}^\top \mathbf{x} = \mathbf{d}^\top \mathbf{L}^+ \mathbf{d}$
$\mathbf{c} \in \mathbb{R}^{ E }$	edge capacities
$\text{val}(\mathbf{f})$	value (routed units of flow) of an s - t flow \mathbf{f}
G_f	residual graph with respect to the flow \mathbf{f}
$\hat{\mathbf{f}}$	flow within the residual graph G_f
U	maximum edge capacity, $\max_e c(e)$
L_f	level graph of residual graph G_f

OPTIMIZATION

$L(\mathbf{y}, \mathbf{x}, \mathbf{s})$	Lagrangian of an optimization problem with primal variables \mathbf{y} and dual variables \mathbf{x}, \mathbf{s}
---	--

Bibliography

Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.

Henri Cartan. *Differential calculus*. Hermann, 1983.

Li Chen, Rasmus Kyng, Yang P Liu, Richard Peng, Maximilian Probst Gutenberg, and Sushant Sachdeva. Maximum flow and minimum-cost flow in almost-linear time. *arXiv preprint arXiv:2203.00671*, 2022.

Rasmus Kyng and Sushant Sachdeva. Approximate gaussian elimination for laplacians-fast, sparse, and simple. In *2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 573–582. IEEE, 2016.

James B Orlin. Max flows in $o(nm)$ time, or better. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 765–774, 2013.

Ralph Tyrell Rockafellar. Convex analysis. In *Convex analysis*. Princeton university press, 2015.

Daniel A Spielman and Nikhil Srivastava. Graph sparsification by effective resistances. *SIAM Journal on Computing*, 40(6):1913–1926, 2011.

Index

s - t path flow, 83

accelerated gradient descent, 47

adjacency matrix, 9

approximate solution, 43

approximate solution to linear system, 79

augmenting flow, 84

Bernstein concentration bound, 29

Bernstein matrix concentration bound, 29

blocking flow, 87

Bolzano-Weierstrass theorem, 103

Boundedness theorem, 103

Cheeger's inequality, 69

Chernoff bound, 29

Cholesky decomposition, 24

commute time, 28

complementary slackness, 54

concave function, 40

condition number, 46

conductance, 67

continuously differentiable, 32

convex function, 40

convex optimization, 40

convex set, 39

Courant-Fischer min-max theorem, 15

cut, 67, 84

cut value, 67

cycle flow, 83

demand, 9

diameter, 65

Dinitz's algorithm, 87

directional derivative, 33

dual feasible, 53

dual program, 55

dual variables, 53

Edmonds-Karp algorithm, 86

effective resistance, 71

electrical energy, 13

electrical energy-minimizing flow, 13

electrical flow, 9

electrical voltages, 9

epigraph, 40

expander, 68

expander decomposition, 68

extreme value theorem, 39

feasible flow, 83

feasible point, 39

feasible set, 39

Fenchel conjugate, 57

first-order expansion, 31

flow, 9, 83

flow value, 83

Ford-Fulkerson algorithm, 85

Fréchet differentiable, 31

gradient, 31

gradient condition, 54

gradient descent, 43

Hessian, 33

hitting time, 27

hyperplane, 51

idempotency, 23

image, 20

incidence matrix, 10

indefinite matrix, 15

infeasible point, 39

Jacobian, 32

Jensen's inequality, 29

Joule's law, 13

Karush-Kuhn-Tucker conditions, 54

kernel, 20

Kirchhoff's current law, 9

Lagrangian, 54

Laplacian matrix, 10

lazy random walk, 26

level, 87

level graph, 87

level set, 40

Lieb's theorem, 22

Loewner order, 19

Mahalanobis norm, 79

Markov property, 25

Markov's inequality, 29

martingale, 30

martingale difference sequence, 30

matrix approximation, 75

matrix function, 21

matrix norm, 18

metric, 73

mixing random walk, 26

Moore-Penrose inverse, 22

net flow, 9

net flow constraint, 10

normal vector, 51

normalized Laplacian matrix, 12

null space, 20

Ohm's law, 9, 10

optimal point, 39

optimization problem, 39

path inequality, 64

Polyak-Łojasiewicz inequality, 46

positive definite matrix, 15

positive semi-definite matrix, 15

primal feasible, 53

primal program, 53

primal-dual feasible, 53

projection matrix, 23

pseudoinverse, 22

quality, 68

quasiconvex function, 40

random walk, 25

range, 20

residual graph, 84

s-t cut, 84

s-t flow, 83

Schur complement, 79

second-order expansion, 34

separating hyperplane, 51

Separating hyperplane theorem, 51

signless Laplacian matrix, 63

Slater's condition, 56

smoothness, 44

sparsity, 69

spectral matrix norm, 18

spectral theorem for symmetric matrices, 15

stationary distribution, 26

stationary point, 32

strictly convex function, 40

strictly feasible, 56

strong convexity, 46

strong duality, 55

sub-level set, 40

Taylor's theorem (first-order), 32

Taylor's theorem (second-order), 34

test vector, 65

tight constraint, 53

- trace, 20
- transition matrix, 25
- twice continuously differentiable, 34
- twice Fréchet differentiable, 34

- volume, 67

- weak duality, 55
- weight (of an edge), 10
- weighted adjacency matrix, 11
- weighted degree, 11