

ETH Zurich

---

# Advanced Graph Algorithms and Optimization

---

Rasmus Kyng & Maximilian Probst Gutenberg

Spring 2024

These notes will be updated throughout the course. They are likely to contain typos, and they may have mistakes or lack clarity in places. Feedback and comments are welcome. Please send to [kyng@inf.ethz.ch](mailto:kyng@inf.ethz.ch) or, even better, submit a pull request at

[https://github.com/rjkyng/agao24\\_script.](https://github.com/rjkyng/agao24_script)

We want to thank scribes from the 2020 edition of the course who contributed to these notes: Hongjie Chen, Meher Chaitanya, Timon Knigge, and Tim Taubner – and we’re grateful to all the readers who’ve submitted corrections, including Martin Kucera, Alejandro Cassis, and Luke Volpatti.

*A important note:* If you’re a student browsing these notes to decide whether to take this course, please note that the current notes are incomplete. We will release parts later in the semester. You can take a look at last year’s notes for an impression of what the rest of the course will look like. Find them here:

[https://github.com/rjkyng/agao23\\_script/raw/main/agao23\\_script.pdf](https://github.com/rjkyng/agao23_script/raw/main/agao23_script.pdf)

There will, however, be some changes to the content compared to last year.

# Contents

<b>1 Course Introduction</b>	<b>9</b>
1.1 Overview . . . . .	9
1.2 Electrical Flows and Voltages - a Graph Problem from Middle School? . . . . .	9
1.3 Convex Optimization . . . . .	16
1.4 More Graph Optimization Problems . . . . .	18
<b>I Introduction to Convex Optimization</b>	<b>21</b>
<b>2 Some Basic Optimization, Convex Geometry, and Linear Algebra</b>	<b>22</b>
2.1 Overview . . . . .	22
2.2 Optimization Problems . . . . .	22
2.3 A Characterization of Convex Functions . . . . .	24
2.3.1 First-order Taylor Approximation . . . . .	25
2.3.2 Directional Derivatives . . . . .	26
2.3.3 Lower Bounding Convex Functions with Affine Functions . . . . .	26
2.4 Conditions for Optimality . . . . .	28
<b>3 Convexity, Second Derivatives and Gradient Descent</b>	<b>29</b>
3.1 A Review of Linear Algebra . . . . .	29
3.2 Characterizations of Convexity and Optimality via Second Derivatives . . . . .	31
3.2.1 A Necessary Condition for Local Extrema . . . . .	32
3.2.2 A sufficient condition for local extrema . . . . .	33

3.2.3	Characterization of convexity . . . . .	33
3.3	Gradient Descent - An Approach to Optimization? . . . . .	35
3.3.1	A Quantitative Bound on Changes in the Gradient . . . . .	35
3.3.2	Analyzing Gradient Descent . . . . .	36
<b>II</b>	<b>Spectral Graph Theory</b>	<b>39</b>
<b>4</b>	<b>Introduction to Spectral Graph Theory</b>	<b>40</b>
4.1	Recap: Incidence and Adjacency Matrices, the Laplacian Matrix and Electrical Energy . . . . .	40
4.2	Understanding Eigenvalues of the Laplacian . . . . .	42
4.3	The Complete Graph . . . . .	43
4.4	The Path Graph . . . . .	43
4.4.1	Upper Bounding $\lambda_2(P_n)$ via Test Vectors . . . . .	43
4.4.2	Lower Bounding $\lambda_n(P_n)$ via Test Vectors . . . . .	44
4.4.3	Upper bounding $\lambda_n(P_n)$ via the Loewner Order . . . . .	45
4.4.4	Lower bounding $\lambda_2(P_n)$ via the Loewner Order . . . . .	47
4.5	The Complete Binary Tree . . . . .	49
4.5.1	Deriving Bounds on $\lambda_n(T_n)$ . . . . .	49
4.5.2	Upper Bounding $\lambda_2$ via Test Vectors . . . . .	50
4.5.3	Lower Bounding $\lambda_2$ . . . . .	50
<b>5</b>	<b>Pseudo-inverses and Effective Resistance</b>	<b>52</b>
5.1	What is a (Moore-Penrose) Pseudoinverse? . . . . .	52
5.2	Electrical Flows Again . . . . .	53
5.3	Effective Resistance . . . . .	54
5.3.1	Effective Resistance is a Distance . . . . .	57
<b>6</b>	<b>Different Perspectives on Gaussian Elimination</b>	<b>59</b>
6.1	An Optimization View of Gaussian Elimination for Laplacians . . . . .	59

6.2	An Additive View of Gaussian Elimination . . . . .	62
<b>7</b>	<b>Conductance, Expanders and Cheeger's Inequality</b>	<b>66</b>
7.1	Conductance and Expanders . . . . .	67
7.2	A Lower Bound for Conductance via Eigenvalues . . . . .	68
7.3	An Upper Bound for Conductance via Eigenvalues . . . . .	70
7.4	Conclusion . . . . .	73
<b>8</b>	<b>Random Walks</b>	<b>74</b>
8.1	A Primer on Random Walks . . . . .	74
8.2	Convergence Results for Random Walks . . . . .	75
8.2.1	Making Random Walks Lazy . . . . .	76
8.2.2	Convergence of Lazy Random Walks . . . . .	77
8.2.3	The Rate of Convergence . . . . .	78
8.3	Properties of Random Walks . . . . .	80
8.3.1	Hitting Times . . . . .	80
8.3.2	Commute Time . . . . .	81
<b>9</b>	<b>Random Matrix Concentration and Spectral Graph Sparsification</b>	<b>83</b>
9.1	Matrix Sampling and Approximation . . . . .	83
9.2	Matrix Concentration . . . . .	86
9.2.1	Matrix Functions . . . . .	87
9.2.2	Monotonicity and Operator Monotonicity . . . . .	87
9.2.3	Some Useful Facts . . . . .	89
9.2.4	Proof of Matrix Bernstein Concentration Bound . . . . .	90
9.3	Spectral Graph Sparsification . . . . .	91
<b>10</b>	<b>Solving Laplacian Linear Equations</b>	<b>98</b>
10.1	Solving Linear Equations Approximately . . . . .	98
10.2	Preconditioning and Approximate Gaussian Elimination . . . . .	99
10.3	Approximate Gaussian Elimination Algorithm . . . . .	100

10.4 Analyzing Approximate Gaussian Elimination . . . . .	103
10.4.1 Normalization, a.k.a. Isotropic Position . . . . .	104
10.4.2 Martingales . . . . .	104
10.4.3 Martingale Difference Sequence as Edge-Samples . . . . .	106
10.4.4 Stopped Martingales . . . . .	107
10.4.5 Sample Norm Control . . . . .	108
10.4.6 Random Matrix Concentration from Trace Exponentials . . . . .	110
10.4.7 Mean-Exponential Bounds from Variance Bounds . . . . .	111
10.4.8 The Overall Mean-Trace-Exponential Bound . . . . .	112
<b>III Combinatorial Graph Algorithms</b>	<b>115</b>
<b>11 Classical Algorithms for Maximum Flow I</b>	<b>116</b>
11.1 Maximum Flow . . . . .	116
11.2 Flow Decomposition . . . . .	117
11.3 Cuts and Minimum Cuts . . . . .	119
11.4 Algorithms for Max flow . . . . .	121
<b>12 Classical Algorithms for Maximum Flow II</b>	<b>128</b>
12.1 Overview . . . . .	128
12.2 Blocking Flows . . . . .	128
12.3 Dinic's Algorithm . . . . .	129
12.4 Finding Blocking Flows . . . . .	131
12.5 Minimum Cut as a Linear Program . . . . .	133
<b>13 The Cut-Matching Game: Expanders via Max Flow</b>	<b>135</b>
13.1 Introduction . . . . .	135
13.2 Embedding Graphs into Expanders . . . . .	136
13.3 The Cut-Matching Algorithm . . . . .	137
13.4 Constructing an Expander via Random Walks . . . . .	140

## IV Further Topics in Convex Optimization 145

<b>14 Separating Hyperplanes, Lagrange Multipliers, KKT Conditions, and Convex Duality</b>	<b>146</b>
14.1 Overview . . . . .	146
14.2 Separating Hyperplane Theorem . . . . .	147
14.3 Lagrange Multipliers and Duality of Convex Problems . . . . .	149
14.3.1 Karush-Kuhn Tucker Optimality Conditions for Convex Problems . .	150
14.3.2 Slater's Condition . . . . .	152
14.3.3 The Lagrangian and The Dual Program . . . . .	153
14.3.4 Strong Duality and KKT . . . . .	155
14.3.5 Proof that Slater's Condition Implies Strong Duality . . . . .	156
<b>15 Fenchel Conjugates and Newton's Method</b>	<b>161</b>
15.1 Lagrange Multipliers and Convex Duality Recap . . . . .	161
15.2 Fenchel Conjugates . . . . .	163
15.3 Newton's Method . . . . .	166
15.3.1 Warm-up: Quadratic Optimization . . . . .	166
15.3.2 $K$ -stable Hessian . . . . .	167
15.3.3 Linearly Constrained Newton's Method . . . . .	168
<b>16 Interior Point Methods for Maximum Flow</b>	<b>171</b>
16.1 An Interior Point Method . . . . .	171
16.1.1 A Barrier Function and an Algorithm . . . . .	172
16.1.2 Updates using Divergence . . . . .	173
16.1.3 Understanding the Divergence Objective . . . . .	175
16.1.4 Quadratically Smoothing Divergence and Local Agreement . . . . .	176
16.1.5 Step size for divergence update . . . . .	179

<b>V Further Topics in Combinatorial Graph Algorithms</b>	<b>182</b>
<b>17 Low-Diameter Decompositions</b>	<b>183</b>
17.1 Low-Diameter Decomposition in Undirected Graphs . . . . .	184
17.2 Generalizing LDDs to Directed Graphs . . . . .	186
<b>18 Negative Single-Source Shortest Paths</b>	<b>191</b>
18.1 A Not-So-Brief Recap: Classic Algorithms for Handling Negative Weights . .	191
18.2 Scaling for the Negative-Weight SSSP Problem . . . . .	194
18.3 The (Normalized) Negative-Weight SSSP Problem . . . . .	195
18.4 A Near-Linear-Time Algorithm . . . . .	198
<b>19 Maximum Matching</b>	<b>199</b>
19.1 Bipartite Matching . . . . .	199
19.1.1 Unit Edge Weights . . . . .	199
19.1.2 General Edge Weights . . . . .	201
19.1.3 Faster Implementation . . . . .	204
19.2 General Graph Matching . . . . .	207
19.2.1 Unit Edge Weights . . . . .	207
19.2.2 General Edge Weights . . . . .	211

# Chapter 1

## Course Introduction

### 1.1 Overview

This course will take us quite deep into modern approaches to graph algorithms using convex optimization techniques. By studying convex optimization through the lens of graph algorithms, we'll try to develop an understanding of fundamental phenomena in optimization. Much of our time will be devoted to flow problems on graphs. We will not only be studying these problems for their own sake, but also because they often provide a useful setting for thinking more broadly about optimization.

The course will cover some traditional discrete approaches to various graph problems, especially flow problems, and then contrast these approaches with modern, asymptotically faster methods based on combining convex optimization with spectral and combinatorial graph theory.

### 1.2 Electrical Flows and Voltages - a Graph Problem from Middle School?

We will dive right into graph problems by considering how electrical current moves through a network of resistors.

First, let us recall some middle school physics. If some of these things don't make sense to you, don't worry, in less than a paragraph from here, we'll be back to safely doing math.

Recall that a typical battery that one buys from Migros has two endpoints, and produces what is called a *voltage difference* between these endpoints.

One end of the battery will have a positive charge (I think that means an excess of positrons<sup>1</sup>), and the other a negative charge. If we connect the two endpoints with a wire, then a current will flow from one end of the battery to the other in an attempt to even out this imbalance of charge.

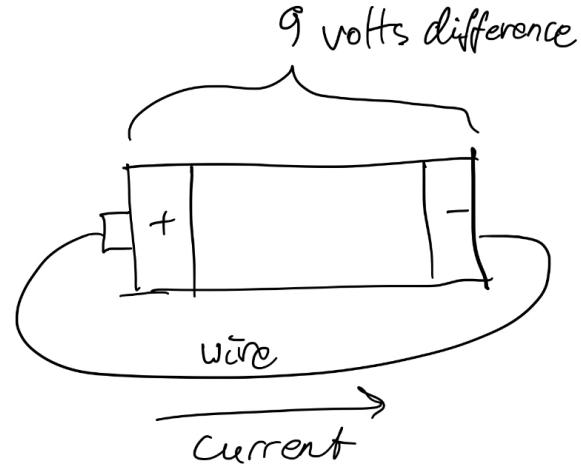


Figure 1.1: A 9 volts battery with a wire attached.

We can also imagine a kind of battery that tries to send a certain amount of current through the wires between its endpoints, e.g. 1 unit of charge per unit of time. This will be a little more convenient to work with, so let us focus on that case.

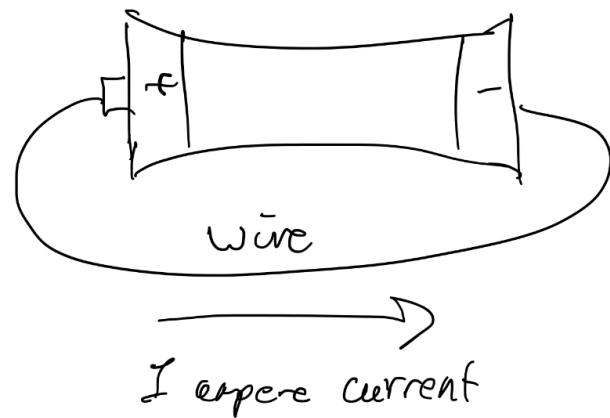


Figure 1.2: A 1 ampere battery with a wire attached.

A *resistor* is a piece of wire that connects two points  $u$  and  $v$ , and is completely described by a single number  $r$  called its *resistance*.

---

<sup>1</sup>I'm joking, of course! Try Wikipedia if you want to know more. However, you will not need it for this class.

If the voltage difference between the endpoints of the resistor is  $x$ , and the resistance is  $r$  then this will create a flow of charge per unit of time of  $f = x/r$ . This is called Ohm's Law.

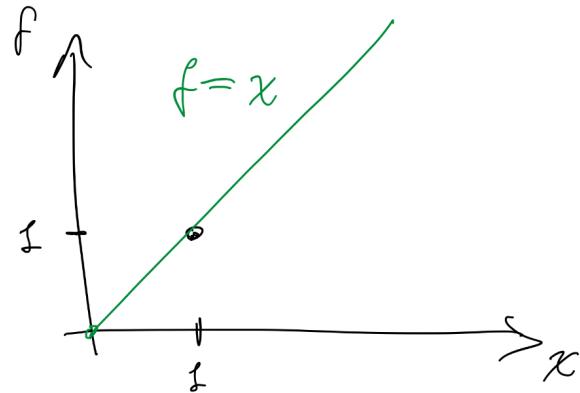


Figure 1.3: Ohm's Law for a resistor with resistance  $r = 1$ .

Suppose we set up a bunch of wires that route electricity from our current source  $s$  to our current sink  $t$  in some pattern:

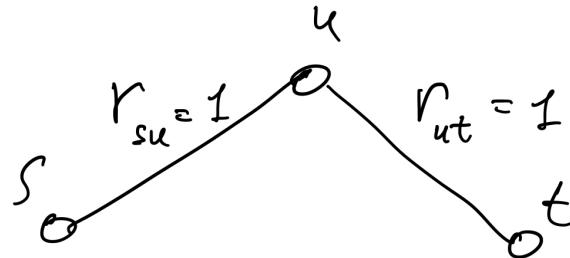


Figure 1.4: A path of two resistors.

We have one unit of charge flowing out of  $s$  per unit of time, and one unit coming into  $t$ . Because charge is conserved, the current flowing into any other point  $u$  must equal the amount flowing out of it. This is called Kirchhoff's Current Law.

To send one unit of current from  $s$  to  $t$ , we must be sending it first from  $s$  to  $u$  and then from  $u$  to  $t$ . So the current on edge  $(s, u)$  is 1 and the current on  $(u, t)$  is 1. By Ohm's Law, the voltage difference must also be 1 across each of the two wires. Thus, if the voltage is  $x$  at  $s$ , it must be  $x + 1$  at  $u$  and  $x + 2$  at  $t$ . What is  $x$ ? It turns out it doesn't matter: We only care about the differences. So let us set  $x = 0$ .

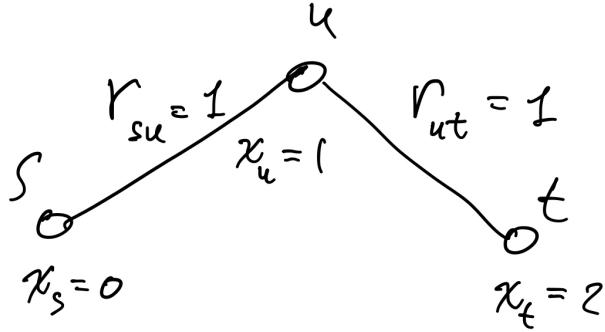


Figure 1.5: A path of two resistors.

Let us try one more example:

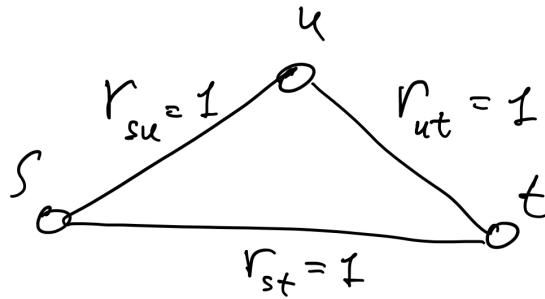


Figure 1.6: A network with three resistors.

How much flow will go directly from  $s$  to  $t$  and how much via  $u$ ?

Well, we know what the net current flowing into and out of each vertex must be, and we can use that to set up some equations. Let us say the voltage at  $s$  is  $x_s$ , at  $u$  is  $x_u$  and at  $t$  is  $x_t$ .

- Net current at  $s$ :  $-1 = (x_s - x_t) + (x_s - x_u)$
- Net current at  $u$ :  $0 = (x_u - x_s) + (x_u - x_t)$
- Net current at  $t$ :  $1 = (x_t - x_s) + (x_t - x_u)$

The following is a solution:  $x_s = 0$ ,  $x_u = \frac{1}{3}$ ,  $x_t = \frac{2}{3}$ . And as before, we can shift all the voltages by some constant  $x$  and get another solution  $x_s = x + 0$ ,  $x_u = x + \frac{1}{3}$ ,  $x_t = x + \frac{2}{3}$ . You might want to convince yourself that these are the only solutions.

**Electrical flows in general graphs.** Do we know enough to calculate the electrical flow in some other network of resistors? To answer this, let us think about the network as a graph. Consider an undirected graph  $G = (V, E)$  with  $|V| = n$  vertices and  $|E| = m$  edges, and let us assume  $G$  is connected. Let's associate a resistance  $r(e) > 0$  with every edge  $e \in E$ .

To keep track of the direction of the flow on each edge, it will be useful to assign an arbitrary direction to every edge. So let's do that, but remember that this is just a bookkeeping tool that helps us track where flow is going.

A *flow* in the graph is a vector  $\mathbf{f} : \mathbb{R}^E$ . The *net flow* of  $\mathbf{f}$  at a vertex  $u \in V$  is defined as  $\sum_{v \rightarrow u} \mathbf{f}(v, u) - \sum_{u \rightarrow v} \mathbf{f}(u, v)$ .

We say a flow routes the demands  $\mathbf{d} \in \mathbb{R}^V$  if the net flow at every vertex  $v$  is  $\mathbf{d}(v)$ .

We can assign a voltage to every vertex  $\mathbf{x} \in R^V$ . Ohm's Law says that the electrical flow induced by these voltages will be  $\mathbf{f}(u, v) = \frac{1}{r(u, v)}(\mathbf{x}(v) - \mathbf{x}(u))$ .

Say we want to route one unit of current from vertex  $s \in V$  to vertex  $t \in V$ . As before, we can write an equation for every vertex saying that the voltage differences must produce the desired net current:

- Net current at  $s$ :  $-1 = \sum_{(s, v)} \frac{1}{r(s, v)}(\mathbf{x}(v) - \mathbf{x}(s))$
- Net current at  $u \in V \setminus \{s, t\}$ :  $0 = \sum_{(u, v)} \frac{1}{r(u, v)}(\mathbf{x}(v) - \mathbf{x}(u))$
- Net current at  $t$ :  $1 = \sum_{(t, v)} \frac{1}{r(t, v)}(\mathbf{x}(v) - \mathbf{x}(t))$

This gives us  $n$  constraints, exactly as many as we have voltage variables. However we have to be a little careful when trying to conclude that a solution exists, yielding voltages  $\mathbf{x}$  that gives induce an electrical flow routing the desired demand.

You will prove in the exercises (Week 1, Exercise 3) that a solution  $\mathbf{x}$  exists. The proof requires two important observations: Firstly that the graph is connected, and secondly that summed over all vertices, the net demand is zero, i.e. as much flow is coming into the network as is leaving it.

**The incidence matrix and the Laplacian matrix.** To have a more compact notation for net flow constraints, we also introduce the *edge-vertex incidence matrix* of the graph,  $\mathbf{B} \in \mathbb{R}^{V \times E}$ .

$$\mathbf{B}(v, e) = \begin{cases} 1 & \text{if } e = (u, v) \\ -1 & \text{if } e = (v, u) \\ 0 & \text{o.w.} \end{cases}$$

Now we can express the net flow constraint that  $\mathbf{f}$  routes  $\mathbf{d}$  by

$$\mathbf{B}\mathbf{f} = \mathbf{d}.$$

This is also called a conservation constraint. In our examples so far, we have  $\mathbf{d}(s) = -1$ ,  $\mathbf{d}(t) = 1$  and  $\mathbf{d}(u) = 0$  for all  $u \in V \setminus \{s, t\}$ .

If we let  $\mathbf{R} = \text{diag}_{e \in E} r(e)$  then Ohm's law tells us that  $\mathbf{f} = \mathbf{R}^{-1} \mathbf{B}^\top \mathbf{x}$ . Putting these observations together, we have  $\mathbf{B} \mathbf{R}^{-1} \mathbf{B}^\top \mathbf{x} = \mathbf{d}$ . The voltages  $\mathbf{x}$  that induce  $\mathbf{f}$  must solve this system of linear equations, and we can use that to compute both  $\mathbf{x}$  and  $\mathbf{f}$ . It is exactly

the same linear equation as the one we considered earlier. We can show that for a connected graph, a solution  $\mathbf{x}$  exists if and only if the flow into the graph equals the net flow out, which we can express as  $\sum_v \mathbf{d}(v) = 0$  or  $\mathbf{1}^\top \mathbf{d} = 0$ . You will show this as part of Exercise 3. This also implies that an electrical flow routing  $\mathbf{d}$  exists if and only if the net flow into the graph equals the net flow out, which we can express as  $\mathbf{1}^\top \mathbf{d} = 0$ .

The matrix  $\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^\top$  is called the *Laplacian* of the graph and is usually denoted by  $\mathbf{L}$ .

**An optimization problem in disguise.** So far, we have looked at electrical voltages and flows as arising from a set of linear equations – and it might not be apparent that this has anything to do with optimization. But transporting current through a resistor requires energy, which will be dissipated as heat by the resistor (i.e. it will get hot!). If we send a current of  $f$  across a resistor with a potential drop of  $x$ , then the amount of energy spent per unit of time by the resistor will be  $f \cdot x$ . This is called Joule's Law. Applying Ohm's law to a resistor with resistance  $r$ , we can also express this energy per unit of time as  $f \cdot x = x^2/r = r \cdot f^2$ . Since we aren't bothering with units, we will even forget about time, and refer to these quantities as “energy”, even though a physicist would call them “power”.

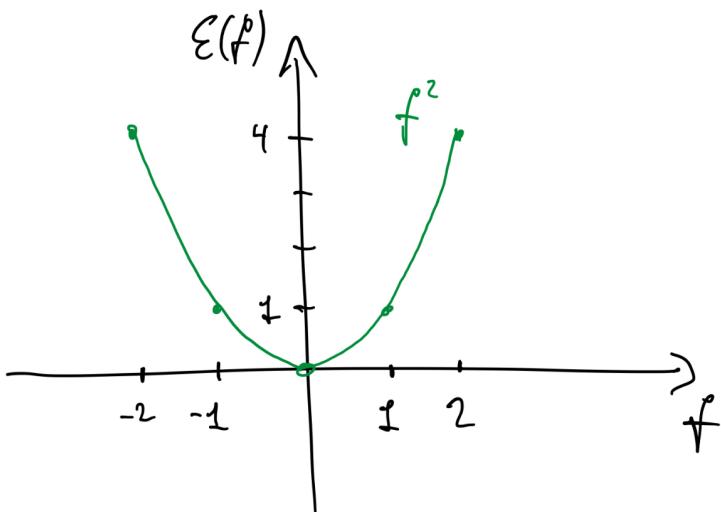


Figure 1.7: Energy as a function of flow in a resistor with resistance  $r = 1$ .

Now, another interesting question would seem to be: If we want to find a flow routing a certain demand  $\mathbf{d}$ , how should the flow behave in order to minimize the electrical energy spent routing the flow? The electrical energy of a flow vector  $\mathbf{f}$  is  $\mathcal{E}(\mathbf{f}) \stackrel{\text{def}}{=} \sum_e \mathbf{r}(e)\mathbf{f}(e)^2$ . We can phrase this as an optimization problem:

$$\begin{aligned} & \min_{\mathbf{f} \in \mathbb{R}^E} \mathcal{E}(\mathbf{f}) \\ & \text{s.t. } \mathbf{B}\mathbf{f} = \mathbf{d}. \end{aligned}$$

We call this problem *electrical energy-minimizing flow*. As we will prove later, the flow  $\mathbf{f}^*$  that minimizes the electrical energy among all flows that satisfy  $\mathbf{Bf} = \mathbf{d}$  is precisely the electrical flow.

**A pair of problems.** What about our voltages, can we also get them from some optimization problem? Well, we can work backwards from the fact that our voltages solve the equation  $\mathbf{Lx} = \mathbf{d}$ . Consider the function  $c(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Lx} - \mathbf{x}^\top \mathbf{d}$ . We should ask ourselves some questions about this function  $c : \mathbb{R}^V \rightarrow \mathbb{R}$ . Is it continuous and continuously differentiable? The answer to this is yes, and that is not hard to see. Does the function have a minimum? This is maybe not immediately clear, but the minimum does indeed exist.

When this is minimized, the derivative of  $c(\mathbf{x})$  with respect to each coordinate of  $\mathbf{x}$  must be zero. This condition yields exactly the system of linear equations  $\mathbf{Lx} = \mathbf{d}$ . You will confirm this in Exercise 4 of the first exercise sheet.

Based on our derivative condition for the optimum, we can also express the electrical voltages as the solution to an optimization problem, namely

$$\min_{\mathbf{x} \in \mathbb{R}^V} c(\mathbf{x})$$

As you are probably aware, having the derivative of each coordinate equal zero is not a sufficient condition for being at the optimum of a function<sup>2</sup>.

It is also interesting to know whether *all* solutions to  $\mathbf{Lx} = \mathbf{d}$  are in fact minimizers of  $c$ . The answer is yes, and we will see some very general tools for proving statements like this in Chapter 2.

Altogether, we can see that routing electrical current through a network of resistors leads to a *pair* of optimization problems, let's call them  $\mathbf{f}^*$  and  $\mathbf{x}^*$ , and that the solutions to the two problems are related, in our case through the equation  $\mathbf{f}^* = \mathbf{R}^{-1} \mathbf{B}^\top \mathbf{x}^*$  (Ohm's Law). But why and how are these two optimization problems related?

Instead of minimizing  $c(\mathbf{x})$ , we can equivalently think about maximizing  $-c(\mathbf{x})$ , which gives the following optimization problem:  $\max_{\mathbf{x} \in \mathbb{R}^V} -c(\mathbf{x})$ . In fact, as you will show in the exercises for Week 1, we have  $\mathcal{E}(\mathbf{f}^*) = -c(\mathbf{x}^*)$ , so the minimum electrical energy is exactly the maximum value of  $-c(\mathbf{x})$ . More generally for *any* flow that routes  $\mathbf{d}$  and *any* voltages  $\mathbf{x}$ , we have  $\mathcal{E}(\mathbf{f}) \geq -c(\mathbf{x})$ . So, for any  $\mathbf{x}$ , the value of  $-c(\mathbf{x})$  is a lower bound on the minimum energy  $\mathcal{E}(\mathbf{f}^*)$ .

This turns out to be an instance of a much broader phenomenon, known as Lagrangian duality, which allows us to learn a lot about many optimization problems by studying two related pairs of problems, a minimization problem, and a related maximization problem that gives lower bounds on the optimal value of the minimization problem.

---

<sup>2</sup>Consider the function in one variable  $c(x) = x^3$ .

**Solving  $\mathbf{L}\mathbf{x} = \mathbf{d}$ .** Given a graph  $G$  with resistances for the edges, and some net flow vector  $\mathbf{d}$ , how quickly can we compute  $\mathbf{x}$ ? Broadly speaking, there are two very different families of algorithms we could use to try to solve this problem.

We could solve the linear equation using something like *Gaussian Elimination* to compute an exact solution.

Alternatively, we could start with a guess at a solution, e.g.  $\mathbf{x}_0 = \mathbf{0}$ , and then we could try to make a change to  $\mathbf{x}_0$  to reach a new point  $\mathbf{x}_1$  with a lower value of  $c(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{L}\mathbf{x} - \mathbf{x}^\top \mathbf{d}$ , i.e.  $c(\mathbf{x}_1) < c(\mathbf{x}_0)$ . If we repeat a process like that for enough steps, say  $t$ , hopefully we eventually reach  $\mathbf{x}_t$  with  $c(\mathbf{x}_t)$  close to  $c(\mathbf{x}^*)$ , where  $\mathbf{x}^*$  is a minimizer of  $c(\mathbf{x})$  and hence  $\mathbf{L}\mathbf{x}^* = \mathbf{d}$ . Now, we also need to make sure that  $c(\mathbf{x}_t) \approx c(\mathbf{x}^*)$  implies that  $\mathbf{L}\mathbf{x}_t \approx \mathbf{d}$  in some useful sense.

One of the most basic algorithms in this framework of “guess and adjust” is called *Gradient Descent*, which we will study in Week 2. The rough idea is the following: if we make a very small step from  $\mathbf{x}$  to  $\mathbf{x} + \boldsymbol{\delta}$ , then a multivariate Taylor expansion suggests that  $c(\mathbf{x} + \boldsymbol{\delta}) - c(\mathbf{x}) \approx \sum_{v \in V} \boldsymbol{\delta}(v) \frac{\partial c(\mathbf{x})}{\partial \mathbf{x}(v)}$ .

If we are dealing with smooth convex function, this quantity is negative if we let  $\boldsymbol{\delta}(v) = -\epsilon \cdot \frac{\partial c(\mathbf{x})}{\partial \mathbf{x}(v)}$  for some small enough  $\epsilon$  so the approximation holds well. So we should be able to make progress by taking a small step in this direction. That’s Gradient Descent! The name comes from the vector of partial derivatives, which is called the gradient.

As we will see later in this course, understanding electrical problems from an optimization perspective is crucial to develop fast algorithms for computing electrical flows and voltages, but to do very well, we also need to borrow some ideas from Gaussian Elimination.

What running times do different approaches get?

1. Using Gaussian Elimination, we can find  $\mathbf{x}$  s.t.  $\mathbf{L}\mathbf{x} = \mathbf{d}$  in  $O(n^3)$  time and with asymptotically faster algorithms based on matrix multiplication, we can bring this down to roughly  $O(n^{2.372})$ .
2. Meanwhile Gradient Descent will get a running time of  $O(n^3m)$  or so – at least this is what a simple analysis suggests.
3. However, we can do much better: By combining ideas from both algorithms, and a bit more, we can get  $\mathbf{x}$  up to very high accuracy in time  $O(m \log^c n)$  where  $c$  is some small constant.

## 1.3 Convex Optimization

Recall our plot in Figure 1.7 of the energy required to route a flow  $f$  across a resistor with resistance  $r$ , which was  $\mathcal{E}(f) = r \cdot f^2$ . We see that the function has a special structure: the

graph of the function sits below the line joining any two points  $(f, \mathcal{E}(f))$  and  $(g, \mathcal{E}(g))$ . A function  $\mathcal{E} : \mathbb{R} \rightarrow \mathbb{R}$  that has this property is said to be convex.

Figure 1.8 shows the energy as a function of flow, along with two points  $(f, \mathcal{E}(f))$  and  $(g, \mathcal{E}(g))$ . We see the function sits below the line segment between these points.

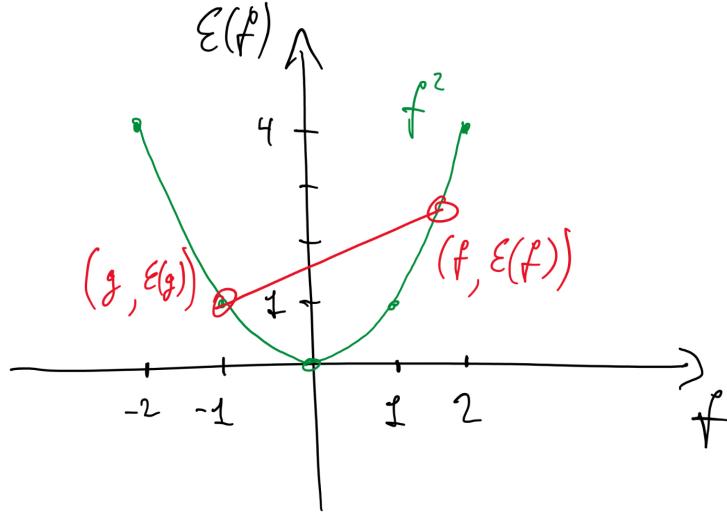


Figure 1.8: Energy as a function of flow in a resistor with resistance  $r = 1$ . The function is convex.

We can also interpret this condition as saying that for all  $\theta \in [0, 1]$

$$\mathcal{E}(\theta f + (1 - \theta)g) \leq \theta \mathcal{E}(f) + (1 - \theta) \mathcal{E}(g).$$

This immediately generalizes to functions  $\mathcal{E} : \mathbb{R}^m \rightarrow \mathbb{R}$ .

A *convex set* is a subset of  $S \subseteq \mathbb{R}^m$  s.t. if  $\mathbf{f}, \mathbf{g} \in S$  then for all  $\theta \in [0, 1]$  we have  $\theta\mathbf{f} + (1 - \theta)\mathbf{g} \in S$ .

Figure 1.9 shows some examples of sets that are and aren't convex.

Convex functions and convex sets are central to optimization, because for most problems of minimization of a convex function over a convex set, we can develop fast algorithms <sup>3</sup>.

So why convex functions and convex sets? One important reason is that for a convex function defined over a convex feasible set, any local minimum is also a global minimum, and this fact makes searching for an optimal solution computationally easier. In fact, this is closely related to why Gradient Descent works well on many convex functions.

Notice that the set  $\{\mathbf{f} : \mathbf{B}\mathbf{f} = \mathbf{d}\}$  is convex, i.e. the set of all flows that route a fixed demand

---

<sup>3</sup>There are some convex optimization problems that are NP-hard. That said, polynomial time algorithms exist for almost any convex problem you can come up with. The most general polynomial time algorithm for convex optimization is probably the Ellipsoid Method.

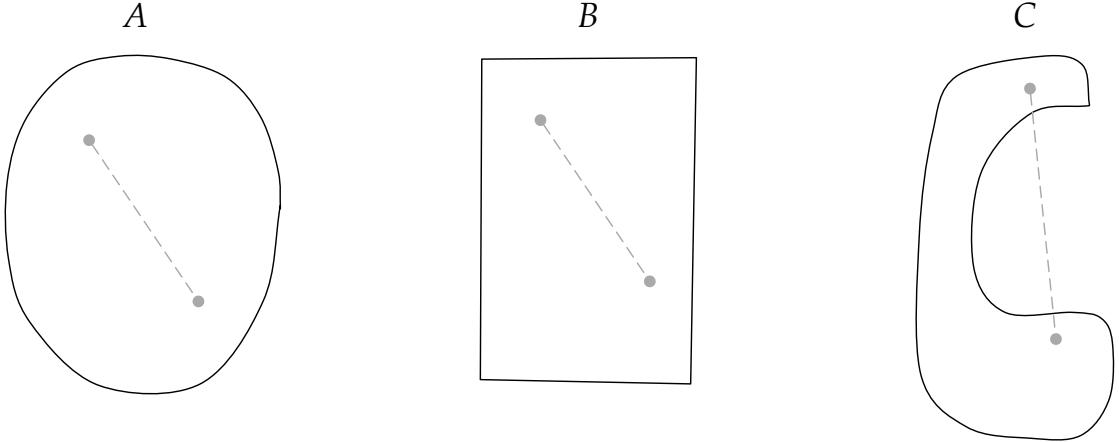


Figure 1.9: A depiction of convex and non-convex sets. The sets  $A$  and  $B$  are convex since the straight line between any two points inside them is also in the set. The set  $C$  is not convex.

$\mathbf{d}$  is convex. It is also easy to verify that  $\mathcal{E}(\mathbf{f}) = \sum_e \mathbf{r}(e)\mathbf{f}(e)^2$  is a convex function, and hence finding an electrical flow is an instance of convex minimization:

## 1.4 More Graph Optimization Problems

**Maximum flow.** Again, let  $G = (V, E)$  be an undirected, connected graph with  $n$  vertices and  $m$  edges. Suppose we want to find a flow  $\mathbf{f} \in \mathbb{R}^E$  that routes  $\mathbf{d}$ , but instead of trying to minimize electrical energy, we try to pick an  $\mathbf{f}$  that minimizes the largest amount of flow on any edge, i.e.  $\max_e |\mathbf{f}_e|$  – which we also denote by  $\|\mathbf{f}\|_\infty$ . We can write this problem as

$$\begin{aligned} & \min_{\mathbf{f} \in \mathbb{R}^E} \|\mathbf{f}\|_\infty \\ \text{s.t. } & \mathbf{Bf} = \mathbf{d} \end{aligned}$$

This problem is known as the Minimum Congested Flow Problem<sup>4</sup>. It is equivalent to the more famous Maximum Flow Problem.

The behavior of this kind of flow is very different than electrical flow. Consider the question of whether a certain demand can be routed  $\|\mathbf{f}\|_\infty \leq 1$ . Imagine sending goods from a source  $s$  to a destination  $t$  using a network of train lines that all have the same capacity and asking whether the network is able to route the goods at the rate you want: This boils down to whether routing exists with  $\|\mathbf{f}\|_\infty \leq 1$ , if we set it up right.

We have a very fast, convex optimization-based algorithm for Minimum Congested Flow: In  $m\epsilon^{-1} \log^{O(1)} n$  time, we can find a flow  $\tilde{\mathbf{f}}$  s.t.  $\mathbf{B}\tilde{\mathbf{f}} = \mathbf{d}$  and  $\|\tilde{\mathbf{f}}\|_\infty \leq (1 + \epsilon) \|\mathbf{f}^*\|_\infty$ , where  $\mathbf{f}^*$

---

<sup>4</sup>This version is called undirected, because the graph is undirected, and *uncapacitated* because we are aiming for the same bound on the flow on all edges.

is an optimal solution, i.e. an actual minimum congestion flow routing  $\mathbf{d}$ .

But what if we want  $\epsilon$  to be very small, e.g.  $1/m$ ? Then this running time isn't so good anymore. But, in this case, we can use other algorithms that find flow  $\mathbf{f}^*$  *exactly*. Unfortunately, these algorithms take time roughly  $m^{4/3+o(1)}$ .

Just as the electrical flow problem had a dual voltage problem, so maximum flow has a dual voltage problem, which is known as the  $s$ - $t$  minimum cut problem.

**Maximum flow, with directions and capacities.** We can make the maximum flow problem harder by introducing directed edges: To do so, we allow edges to exist in both directions between vertices, and we require that the flow on a directed edge is always non-negative. So now  $G = (V, E)$  is a directed graph. We can also make the problem harder by introducing capacities. We define a capacity vector  $\mathbf{c} \in \mathbb{R}^E \geq \mathbf{0}$  and try to minimize  $\|\mathbf{C}^{-1}\mathbf{f}\|_\infty$ , where  $\mathbf{C} = \text{diag}_{e \in E} \mathbf{c}(e)$ . Then our problem becomes

$$\begin{aligned} & \min_{\mathbf{f} \in \mathbb{R}^E} \|\mathbf{C}^{-1}\mathbf{f}\|_\infty \\ \text{s.t. } & \mathbf{Bf} = \mathbf{d} \\ & \mathbf{f} \geq \mathbf{0}. \end{aligned}$$

For this capacitated, directed maximum flow problem, our best algorithms run in about  $O(m\sqrt{n})$  time in sparse graphs and  $O(m^{1.483})$  in dense graphs<sup>5</sup>, even if we are willing to accept fairly low accuracy solution. If the capacities are allowed to be exponentially large, the best running time we can get is  $O(mn)$ . For this problem, we do not yet know how to improve over classical combinatorial algorithms using convex optimization.

**Multi-commodity flow.** We can make the problem even harder still, by simultaneously trying to route two types of flow (imagine pipes with Coke and Pepsi). Our problem now looks like

$$\begin{aligned} & \min_{\mathbf{f}_1, \mathbf{f}_2 \in \mathbb{R}^E} \|\mathbf{C}^{-1}(\mathbf{f}_1 + \mathbf{f}_2)\|_\infty \\ \text{s.t. } & \mathbf{Bf}_1 = \mathbf{d}_1 \\ & \mathbf{Bf}_2 = \mathbf{d}_2 \\ & \mathbf{f}_1, \mathbf{f}_2 \geq \mathbf{0}. \end{aligned}$$

Solving this problem to high accuracy is essentially as hard as solving a general linear program! We should see later in the course how to make this statement precise.

If we in the above problem additionally require that our flows must be integer valued, i.e.  $\mathbf{f}_1, \mathbf{f}_2 \in \mathbb{N}_0$ , then the problem becomes NP-complete.

---

<sup>5</sup>Provided the capacities are integers satisfying a condition like  $\mathbf{c} \leq n^{100}\mathbf{1}$ .

**Random walks in a graph.** Google famously uses<sup>6</sup> the PageRank problem to help decide how to rank their search results. This problem essentially boils down to computing the *stable distribution* of a random walk on a graph. Suppose  $G = (V, E)$  is a directed graph where each outgoing edge  $(v, u)$ , which we will define as going from  $u$  to  $v$ , has a transition probability  $p_{(v,u)} > 0$  s.t.  $\sum_{z \leftarrow u} p_{(z,u)} = 1$ . We can take a step of a random walk on the vertex set by starting at some vertex  $u_0 = u$ , and then randomly pick one of the outgoing edges  $(v, u)$  with probability  $p_{(v,u)}$  and move to the chosen vertex  $u_1 = v$ . Repeating this procedure, to take a step from the next vertex  $u_1$ , gives us a *random walk* in the graph, a sequence of vertices  $u_0, u_1, u_2, \dots, u_k$ .

We let  $\mathbf{P} \in \mathbb{R}^{V \times V}$  be the matrix of transition probabilities given by

$$\mathbf{P}_{vu} = \begin{cases} p_{(v,u)} & \text{for } (u, v) \in E \\ 0 & \text{o.w.} \end{cases}$$

Any probability distribution over the vertices can be specified by a vector  $\mathbf{p} \in \mathbb{R}^V$  where  $\mathbf{p} \geq \mathbf{0}$  and  $\sum_v \mathbf{p}(v) = 1$ . We say that probability distribution  $\boldsymbol{\pi}$  on the vertices is a *stable distribution* of the random walk if  $\boldsymbol{\pi} = \mathbf{P}\boldsymbol{\pi}$ . A strongly connected graph always has exactly one stable distribution.

How quickly can we compute the stable distribution of a general random walk? Under some mild conditions on the stable distribution<sup>7</sup>, we can find a high accuracy approximation of  $\boldsymbol{\pi}$  in time  $O(m \log^c n)$  for some constant  $c$ .

This problem does not easily fit in a framework of convex optimization, but nonetheless, our fastest algorithms for it use ideas from convex optimization.

## Topics in this Course

In this course, we will try to address the following questions.

1. What are the fundamental tools of fast convex optimization?
2. What are some problems we can solve quickly on graphs using optimization?
3. What can graphs teach us about convex optimization?
4. What algorithm design techniques are good for getting algorithms that quickly find a crude approximate solution? And what techniques are best when we need to get a highly accurate answer?
5. What is special about flow problems?

---

<sup>6</sup>At least they did at some point.

<sup>7</sup>Roughly something like  $\max_v 1/\boldsymbol{\pi}(v) \leq n^{100}$ .

## Part I

# Introduction to Convex Optimization

# Chapter 2

## Some Basic Optimization, Convex Geometry, and Linear Algebra

### 2.1 Overview

In this chapter, we will

1. Start with an overview (i.e. this list).
2. Learn some basic terminology and facts about optimization.
3. Recall our definition of convex functions and see how convex functions can also be understood in terms of a characterization based on first derivatives.
4. See how the first derivatives of a convex function can certify that we are at a global minimum.

### 2.2 Optimization Problems

Focusing for now on optimization over  $\mathbf{x} \in \mathbb{R}^n$ , we usually write optimization problems as:

$$\begin{aligned} & \min_{\mathbf{x} \in \mathbb{R}^n} (\text{or } \max) f(\mathbf{x}) \\ & \quad s.t. \quad g_1(\mathbf{x}) \leq b_1 \\ & \quad \quad \quad \vdots \\ & \quad \quad \quad \vdots \\ & \quad \quad \quad g_m(\mathbf{x}) \leq b_m \end{aligned}$$

where  $\{g_i(\mathbf{x})\}_{i=1}^m$  encode the constraints. For example, in the following optimization problem from the previous chapter

$$\begin{aligned} \min_{\mathbf{f} \in \mathbb{R}^E} & \sum_e \mathbf{r}(e) \mathbf{f}(e)^2 \\ \text{s.t. } & \mathbf{B}\mathbf{f} = \mathbf{d} \end{aligned}$$

we have the constraint  $\mathbf{B}\mathbf{f} = \mathbf{d}$ . Notice that we can rewrite this constraint as  $\mathbf{B}\mathbf{f} \leq \mathbf{d}$  and  $-\mathbf{B}\mathbf{f} \leq -\mathbf{d}$  to match the above setting. The set of points which respect the constraints is called the *feasible set*.

**Definition 2.2.1.** For a given optimization problem the set  $\mathcal{F} = \{\mathbf{x} \in \mathbb{R}^n : g_i(\mathbf{x}) \leq b_i, \forall i \in [m]\}$  is called the **feasible set**. A point  $\mathbf{x} \in \mathcal{F}$  is called a **feasible point**, and a point  $\mathbf{x}' \notin \mathcal{F}$  is called an **infeasible point**.

Ideally, we would like to find optimal solutions for the optimization problems we consider. Let's define what we mean exactly.

**Definition 2.2.2.** For a *maximization* problem  $\mathbf{x}^*$  is called an **optimal solution** if  $f(\mathbf{x}^*) \geq f(\mathbf{x}), \forall \mathbf{x} \in \mathcal{F}$ . Similarly, for a *minimization* problem  $\mathbf{x}^*$  is an optimal solution if  $f(\mathbf{x}^*) \leq f(\mathbf{x}), \forall \mathbf{x} \in \mathcal{F}$ .

What happens if there are *no feasible points*? In this case, an optimal solution cannot exist, and we say the problem is infeasible.

**Definition 2.2.3.** If  $\mathcal{F} = \emptyset$  we say that the optimization problem is **infeasible**. If  $\mathcal{F} \neq \emptyset$  we say the optimization problem is **feasible**.

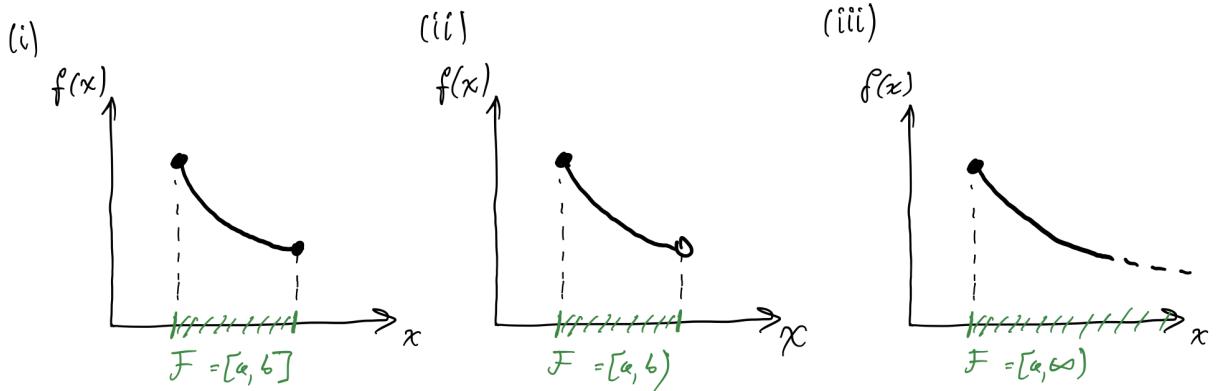


Figure 2.1

Consider three examples depicted in Figure 2.1:

(i)  $\mathcal{F} = [a, b]$

(ii)  $\mathcal{F} = (a, b)$

(iii)  $\mathcal{F} = [a, \infty)$

In the first example, the minimum of the function is attained at  $b$ . In the second case the region is open and therefore there is no minimum function value, since for every point we will choose, there will always be another point with a smaller function value. Lastly, in the third example, the region is unbounded and the function decreasing, thus again there will always be another point with a smaller function value.

**Sufficient Condition for Optimality.** The following theorem, which is a fundamental theorem in real analysis, gives us a sufficient (though not necessary) condition for optimality.

**Theorem** (Extreme Value Theorem). Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a continuous function and  $\mathcal{F} \subseteq \mathbb{R}^n$  be nonempty, bounded, and closed. Then, the optimization problem  $\min f(\mathbf{x}) : \mathbf{x} \in \mathcal{F}$  has an optimal solution.

## 2.3 A Characterization of Convex Functions

Recall the definitions of convex sets and convex functions that we introduced in Chapter 1:

**Definition 2.3.1.** A set  $S \subseteq \mathbb{R}^n$  is called a **convex set** if any two points in  $S$  contain their line, i.e. for any  $\mathbf{x}, \mathbf{y} \in S$  we have that  $\theta\mathbf{x} + (1 - \theta)\mathbf{y} \in S$  for any  $\theta \in [0, 1]$ .

**Definition 2.3.2.** For a convex set  $S \subseteq \mathbb{R}^n$ , we say that a function  $f : S \rightarrow \mathbb{R}$  is **convex on  $S$**  if for any two points  $\mathbf{x}, \mathbf{y} \in S$  and any  $\theta \in [0, 1]$  we have that:

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

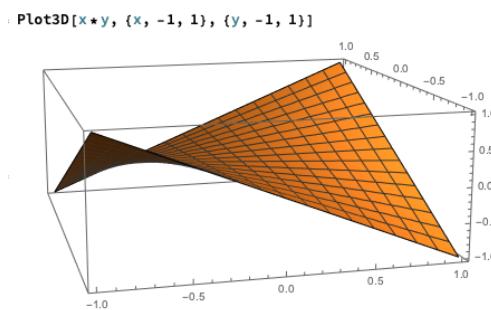


Figure 2.2: This plot shows the function  $f(x, y) = xy$ . For any fixed  $y_0$ , the function  $h(x) = f(x, y_0) = xy_0$  is linear in  $x$ , and so is a convex function in  $x$ . But is  $f$  convex?

We will first give an important characterization of convex function. To do so, we need to characterize multivariate functions via their Taylor expansion.

**Notation for this section.** In the rest of this section, we frequently consider a multivariate function  $f$  whose domain is a set  $S \subseteq \mathbb{R}^n$ , which we will require to be open. When we additionally require that  $S$  is convex, we will specify this. Note that  $S = \mathbb{R}^n$  is both open and convex and it suffices to keep this case in mind. Things sometimes get more complicated if  $S$  is not open, e.g. when the domain of  $f$  has a boundary. We will leave those complications for another time.

### 2.3.1 First-order Taylor Approximation

**Definition 2.3.3.** The **gradient** of a function  $f : S \rightarrow \mathbb{R}$  at point  $\mathbf{x} \in S$  is denoted  $\nabla f(\mathbf{x})$  and is defined as

$$\nabla f(\mathbf{x}) = \left[ \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}(1)}, \dots, \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}(n)} \right]^\top$$

**First-order Taylor expansion.** For a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  of a single variable, differentiable at  $x \in \mathbb{R}$

$$f(x + \delta) = f(x) + f'(x)\delta + o(|\delta|)$$

where by definition:

$$\lim_{\delta \rightarrow 0} \frac{o(|\delta|)}{|\delta|} = 0.$$

Similarly, a multivariate function  $f : S \rightarrow \mathbb{R}$  is said to be (*Fréchet*) *differentiable* at  $\mathbf{x} \in S$  when there exists  $\nabla f(\mathbf{x}) \in \mathbb{R}^n$  s.t.

$$\lim_{\boldsymbol{\delta} \rightarrow 0} \frac{\|f(\mathbf{x} + \boldsymbol{\delta}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top \boldsymbol{\delta}\|_2}{\|\boldsymbol{\delta}\|_2} = 0.$$

Note that this is equivalent to saying that  $f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \boldsymbol{\delta} + o(\|\boldsymbol{\delta}\|_2)$ .

We say that  $f$  is *continuously differentiable* on a set  $S \subseteq \mathbb{R}^n$  if it is differentiable and in addition the gradient is continuous on  $S$ . A differentiable convex function whose domain is an open convex set  $S \subseteq \mathbb{R}^n$  is always continuously differentiable<sup>1</sup>.

**Remark.** In this course, we will generally err on the side of being informal about functional analysis when we can afford to, and we will not worry too much about the details of different notions of differentiability (e.g. Fréchet and Gateaux differentiability), except when it turns out to be important.

**Theorem 2.3.4** (Taylor's Theorem, multivariate first-order remainder form). *If  $f : S \rightarrow \mathbb{R}$  is continuously differentiable over  $[\mathbf{x}, \mathbf{y}]$ , then for some  $\mathbf{z} \in [\mathbf{x}, \mathbf{y}]$ ,*

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

---

<sup>1</sup>See p. 248, Corollary 25.5.1 in *Convex Analysis* by Rockafellar (my version is the Second print, 1972). Rockefellar's corollary concerns finite convex functions, because he otherwise allows convex functions that may take on the values  $\pm\infty$ .

This theorem is useful for showing that the function  $f$  can be approximated by the affine function  $\mathbf{y} \rightarrow f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$  when  $\mathbf{y}$  is “close to”  $\mathbf{x}$  in some sense.

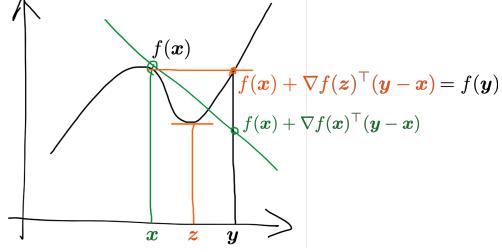


Figure 2.3: The function  $f(\mathbf{y})$  equals its linear approximation based at  $\mathbf{x}$ , with gradient coming from an intermediate point  $\mathbf{z}$ , i.e.  $f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{x})$ .

### 2.3.2 Directional Derivatives

**Definition 2.3.5.** Let  $f : S \rightarrow \mathbb{R}$  be a function differentiable at  $\mathbf{x} \in S$  and let us consider  $\mathbf{d} \in \mathbb{R}^n$ . We define the **derivative of  $f$  at  $\mathbf{x}$  in direction  $\mathbf{d}$**  as:

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

**Proposition 2.3.6.**  $Df(\mathbf{x})[\mathbf{d}] = \nabla f(\mathbf{x})^\top \mathbf{d}$ .

*Proof.* Using the first order expansion of  $f$  at  $\mathbf{x}$ :

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

hence, dividing by  $\lambda$  (and noticing that  $\|\lambda \mathbf{d}\|_2 = \lambda \|\mathbf{d}\|_2$ ):

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

letting  $\lambda$  go to 0 concludes the proof. □

### 2.3.3 Lower Bounding Convex Functions with Affine Functions

In order to prove the characterization of convex functions in the next section we will need the following lemma. This lemma says that any differentiable convex function can be lower bounded by an affine function.

**Theorem 2.3.7.** Let  $S$  be an open convex subset of  $\mathbb{R}^n$ , and let  $f : S \rightarrow \mathbb{R}$  be a differentiable function. Then,  $f$  is convex if and only if for any  $\mathbf{x}, \mathbf{y} \in S$  we have that  $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$ .

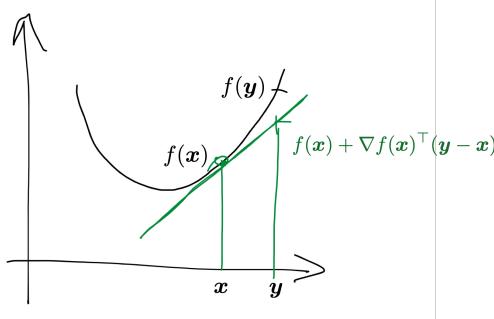


Figure 2.4: The convex function  $f(\mathbf{y})$  sits above the linear function in  $\mathbf{y}$  given by  $f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$ .

*Proof.* [  $\implies$  ] Assume  $f$  is convex, then for all  $\mathbf{x}, \mathbf{y} \in S$  and  $\theta \in [0, 1]$ , if we let  $\mathbf{z} = \theta\mathbf{y} + (1 - \theta)\mathbf{x}$ , we have that

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

and therefore by subtracting  $f(\mathbf{x})$  from both sides we get:

$$\begin{aligned} f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})) - f(\mathbf{x}) &\leq \theta f(\mathbf{y}) + (1 - \theta)f(\mathbf{x}) - f(\mathbf{x}) \\ &= \theta f(\mathbf{y}) - \theta f(\mathbf{x}). \end{aligned}$$

Thus we get that (for  $\theta > 0$ ):

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

Applying Proposition 2.3.6 with  $\mathbf{d} = \mathbf{y} - \mathbf{x}$  we have that:

$$\nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) = \lim_{\theta \rightarrow 0^+} \frac{f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\theta} \leq f(\mathbf{y}) - f(\mathbf{x}).$$

[  $\Leftarrow$  ] Assume that  $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$  for all  $\mathbf{x}, \mathbf{y} \in S$  and show that  $f$  is convex. Let  $\mathbf{x}, \mathbf{y} \in S$  and  $\mathbf{z} = \theta\mathbf{y} + (1 - \theta)\mathbf{x}$ . By our assumption we have that:

$$f(\mathbf{y}) \geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{z}) \tag{2.1}$$

$$f(\mathbf{x}) \geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{z}) \tag{2.2}$$

Observe that  $\mathbf{y} - \mathbf{z} = (1 - \theta)(\mathbf{y} - \mathbf{x})$  and  $\mathbf{x} - \mathbf{z} = \theta(\mathbf{x} - \mathbf{y})$ . Thus adding  $\theta$  times (2.1) to  $(1 - \theta)$  times (2.2) gives cancellation of the vectors multiplying the gradient, yielding

$$\begin{aligned} \theta f(\mathbf{y}) + (1 - \theta)f(\mathbf{x}) &\geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top \mathbf{0} \\ &= f(\theta\mathbf{y} + (1 - \theta)\mathbf{x}) \end{aligned}$$

This is exactly the definition of convexity.  $\square$

## 2.4 Conditions for Optimality

We now want to find necessary and sufficient conditions for local optimality.

**Definition 2.4.1.** Consider a differentiable function  $f : S \rightarrow \mathbb{R}$ . A point  $\mathbf{x} \in S$  at which  $\nabla f(\mathbf{x}) = \mathbf{0}$  is called a **stationary point**.

**Proposition 2.4.2.** *If  $\mathbf{x}$  is a local extremum of a differentiable function  $f : S \rightarrow \mathbb{R}$  then  $\nabla f(\mathbf{x}) = \mathbf{0}$ .*

*Proof.* Let us assume that  $\mathbf{x}$  is a local minimum for  $f$ . Then for all  $\mathbf{d} \in \mathbb{R}^n$ ,  $f(\mathbf{x}) \leq f(\mathbf{x} + \lambda \mathbf{d})$  for  $\lambda$  small enough. Hence:

$$0 \leq f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x}) = \lambda \nabla f(\mathbf{x})^\top \mathbf{d} + o(\|\lambda \mathbf{d}\|)$$

dividing by  $\lambda > 0$  and letting  $\lambda \rightarrow 0^+$ , we obtain  $0 \leq \nabla f(\mathbf{x})^\top \mathbf{d}$ . But, taking  $\mathbf{d} = -\nabla f(\mathbf{x})$ , we get  $0 \leq -\|\nabla f(\mathbf{x})\|_2^2$ . This implies that  $\nabla f(\mathbf{x}) = \mathbf{0}$ .

The case where  $\mathbf{x}$  is a local maximum can be dealt with similarly.  $\square$

**Remark 2.4.3.** For this proposition to hold, it is important that  $S$  is open.

For convex functions however it turns out that a stationary point necessarily implies that the function is at its minimum. Together with the proposition above, this says that for a convex function on  $\mathbb{R}^n$  a point is optimal if and only if it is stationary.

**Proposition 2.4.4.** *Let  $S \subseteq \mathbb{R}^n$  be an open convex set and let  $f : S \rightarrow \mathbb{R}$  be a differentiable and convex function. If  $\mathbf{x}$  is a stationary point then  $\mathbf{x}$  is a global minimum.*

*Proof.* From Theorem 2.3.7 we know that for all  $\mathbf{x}, \mathbf{y} \in S$ :  $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})(\mathbf{y} - \mathbf{x})$ . Since  $\nabla f(\mathbf{x}) = \mathbf{0}$  this implies that  $f(\mathbf{y}) \geq f(\mathbf{x})$ . As this holds for any  $\mathbf{y} \in S$ ,  $\mathbf{x}$  is a global minimum.

$\square$

# Chapter 3

## Convexity, Second Derivatives and Gradient Descent

**Notation for this chapter.** In this chapter, we sometimes consider a multivariate function  $f$  whose domain is a set  $S \subseteq \mathbb{R}^n$ , which we will require to be open. When we additionally require that  $S$  is convex, we will specify this. Note that  $S = \mathbb{R}^n$  is both open and convex and it suffices to keep this case in mind. Things sometimes get more complicated if  $S$  is not open, e.g. when the domain of  $f$  has a boundary. We will leave those complications for another time.

### 3.1 A Review of Linear Algebra

**Semi-definiteness of a matrix.** The following classification of symmetric matrices will be useful.

**Definition 3.1.1.** Let  $\mathbf{A}$  by a symmetric matrix in  $\mathbb{R}^{n \times n}$ . We say that  $\mathbf{A}$  is:

1. *positive definite* iff  $\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$  for all  $\mathbf{x} \in \mathbb{R}^n \setminus \{0\}$ ;
2. *positive semidefinite* iff  $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$  for all  $\mathbf{x} \in \mathbb{R}^n$ ;
3. If neither  $\mathbf{A}$  nor  $-\mathbf{A}$  is positive semi-definite, we say that  $\mathbf{A}$  is *indefinite*.

**Example: indefinite matrix.** Consider the following matrix  $\mathbf{A}$ :

$$\mathbf{A} := \begin{bmatrix} +4 & -1 \\ -1 & -2 \end{bmatrix}$$

For  $\mathbf{x} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ , we have  $\mathbf{x}^\top \mathbf{A} \mathbf{x} = 4 > 0$ . For  $\mathbf{x} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$  we have  $\mathbf{x}^\top \mathbf{A} \mathbf{x} = -2 < 0$ .  $\mathbf{A}$  is therefore indefinite.

The following theorem gives a useful characterization of (semi)definite matrices.

**Theorem 3.1.2.** *Let  $\mathbf{A}$  be a symmetric matrix in  $\mathbb{R}^{n \times n}$ .*

1.  $\mathbf{A}$  is positive definite iff all its eigenvalues are positive;
2.  $\mathbf{A}$  is positive semidefinite iff all its eigenvalues are non-negative;

In order to prove this theorem, let us first recall the Spectral Theorem for symmetric matrices.

**Theorem 3.1.3** (The Spectral Theorem for Symmetric Matrices). *For all symmetric  $\mathbf{A} \in \mathbb{R}^{n \times n}$  there exist  $\mathbf{V} \in \mathbb{R}^{n \times n}$  and a diagonal matrix  $\Lambda \in \mathbb{R}^{n \times n}$  s.t.*

1.  $\mathbf{A} = \mathbf{V} \Lambda \mathbf{V}^\top$ .
2.  $\mathbf{V}^\top \mathbf{V} = \mathbf{I}$  (the  $n \times n$  identity matrix). I.e. the columns of  $\mathbf{V}$  form an orthonormal basis. Furthermore,  $\mathbf{v}_i$  is an eigenvector of  $\lambda_i(\mathbf{A})$ , the  $i$ th eigenvalue of  $\mathbf{A}$ .
3.  $\Lambda_{ii} = \lambda_i(\mathbf{A})$ .

Using the Spectral Theorem, we can show the following result:

**Theorem 3.1.4** (The Courant-Fischer Theorem). *Let  $\mathbf{A}$  be a symmetric matrix in  $\mathbb{R}^{n \times n}$ , with eigenvalues  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ . Then*

1.

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

2.

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

Theorem 3.1.2 is an immediate corollary of Theorem 3.1.4, since we can see that the minimum value of the quadratic form  $\mathbf{x}^\top \mathbf{A} \mathbf{x}$  over  $\mathbf{x} \in W = \mathbb{R}^n$  is  $\lambda_1(\mathbf{A}) \|\mathbf{x}\|_2^2$ .

*Proof of Theorem 3.1.4.* We start by showing Part 1.

Consider letting  $W = \text{span} \{ \mathbf{v}_1, \dots, \mathbf{v}_i \}$ , and normalize  $\mathbf{x} \in W$  so that  $\|\mathbf{x}\|_2 = 1$ . Then  $\mathbf{x} = \sum_{j=1}^i \mathbf{c}(j) \mathbf{v}_j$  for some vector  $\mathbf{c} \in \mathbb{R}^i$  with  $\|\mathbf{c}\|_2 = 1$ .

Using the decomposition from Theorem 3.1.3  $\mathbf{A} = \mathbf{V} \Lambda \mathbf{V}^\top$  where  $\Lambda$  is a diagonal matrix of eigenvalues of  $\mathbf{A}$ , which we take to be sorted in increasing order. Then  $\mathbf{x}^\top \mathbf{A} \mathbf{x} = \mathbf{x}^\top \mathbf{V}^\top \Lambda \mathbf{V} \mathbf{x} = (\mathbf{V} \mathbf{x})^\top \Lambda (\mathbf{V} \mathbf{x}) = \sum_{j=1}^i \lambda_j \mathbf{c}(j)^2 \leq \lambda_i \|\mathbf{c}\|_2^2 = \lambda_i$ . So this choice of  $W$  ensures the maximizer cannot achieve a value above  $\lambda_i$ .

But is it possible that the “minimizer” can do better by choosing a different  $W$ ? Let  $T = \text{span}\{\mathbf{v}_i, \dots, \mathbf{v}_n\}$ . As  $\dim(T) = n+1-i$  and  $\dim(W) = i$ , we must have  $\dim(W \cap T) \geq 1$ , by a standard property of subspaces. Hence for any  $W$  of this dimension,

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

where the last equality follows from a similar calculation to our first one. Thus,  $\lambda_i$  can always be achieved by the “maximizer” for all  $W$  of this dimension.

Part 2 can be dealt with similarly. □

**Example: a positive semidefinite matrix.** Consider the following matrix  $\mathbf{A}$ :

$$\mathbf{A} := \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$

For  $\mathbf{x} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ , we have  $\mathbf{A}\mathbf{x} = \mathbf{0}$ , so  $\lambda = 0$  is an eigenvalue of  $\mathbf{A}$ . For  $\mathbf{x} = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ , we have  $\mathbf{A}\mathbf{x} = \begin{pmatrix} 2 \\ -2 \end{pmatrix} = 2\mathbf{x}$ , so  $\lambda = 2$  is the other eigenvalue of  $\mathbf{A}$ . As both are non-negative, by the theorem above,  $\mathbf{A}$  is positive semidefinite.

Since we are learning about symmetric matrices, there is one more fact that everyone should know about them. We’ll use  $\lambda_{\max}(\mathbf{A})$  denote maximum eigenvalue of a matrix  $\mathbf{A}$ , and  $\lambda_{\min}(\mathbf{A})$  the minimum.

**Claim 3.1.5.** *For a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\|\mathbf{A}\| = \max(|\lambda_{\max}(\mathbf{A})|, |\lambda_{\min}(\mathbf{A})|)$ .*

## 3.2 Characterizations of Convexity and Optimality via Second Derivatives

We will now use the second derivatives of a function to obtain characterizations of convexity and optimality. We will begin by introducing the *Hessian*, the matrix of pairwise second derivatives of a function. We will see that it plays a role in approximating a function via a second-order Taylor expansion. We will then use *semi-definiteness* of the Hessian matrix to characterize both conditions of optimality as well as the convexity of a function.

**Definition 3.2.1.** Given a function  $f : S \rightarrow \mathbb{R}$  its **Hessian** matrix at point  $\mathbf{x} \in S$  denoted  $\mathbf{H}_f(\mathbf{x})$  (also sometimes denoted  $\nabla^2 f(\mathbf{x})$ ) is:

$$\mathbf{H}_f(\mathbf{x}) := \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x(1)^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x(1)\partial x(2)} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x(1)\partial x(n)} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x(2)\partial x(1)} & \frac{\partial^2 f(\mathbf{x})}{\partial x(2)^2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x(2)\partial x(n)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x(n)\partial x(1)} & \frac{\partial^2 f(\mathbf{x})}{\partial x(n)\partial x(2)} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x(n)^2} \end{bmatrix}$$

**Second-order Taylor expansion.** When  $f$  is twice differentiable it is possible to obtain an approximation of  $f$  by quadratic functions. Our definition of  $f : S \rightarrow \mathbb{R}$  being twice (Fréchet) differentiable at  $\mathbf{x} \in S$  is that there exists  $\nabla f(\mathbf{x}) \in \mathbb{R}^n$  and  $\mathbf{H}_f(\mathbf{x}) \in R^{n \times n}$  s.t.

$$\lim_{\delta \rightarrow 0} \frac{\|f(\mathbf{x} + \boldsymbol{\delta}) - f(\mathbf{x}) - (\nabla f(\mathbf{x})^\top \boldsymbol{\delta} + \frac{1}{2}\boldsymbol{\delta}^\top \mathbf{H}_f(\mathbf{x})\boldsymbol{\delta})\|_2}{\|\boldsymbol{\delta}\|_2^2} = 0.$$

This is equivalent to saying that for all  $\boldsymbol{\delta}$

$$f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \boldsymbol{\delta} + \frac{1}{2}\boldsymbol{\delta}^\top \mathbf{H}_f(\mathbf{x})\boldsymbol{\delta} + o(\|\boldsymbol{\delta}\|_2^2).$$

where by definition:

$$\lim_{\boldsymbol{\delta} \rightarrow 0} \frac{o(\|\boldsymbol{\delta}\|_2^2)}{\|\boldsymbol{\delta}\|_2^2} = 0$$

We say that  $f$  is *twice continuously differentiable* on a set  $S \subseteq \mathbb{R}^n$  if it is twice differentiable and in addition the gradient and Hessian are continuous on  $S$ .

As for first order expansions, we have a Taylor's Theorem, which we state in the so-called remainder form.

**Theorem 3.2.2** (Taylor's Theorem, multivariate second-order remainder form). *If  $f : S \rightarrow \mathbb{R}$  is twice continuously differentiable over  $[\mathbf{x}, \mathbf{y}]$ , then for some  $\mathbf{z} \in [\mathbf{x}, \mathbf{y}]$ ,*

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

### 3.2.1 A Necessary Condition for Local Extrema

Recall that in the previous chapter, we show the following proposition.

**Proposition 3.2.3.** *If  $\mathbf{x}$  is a local extremum of a differentiable function  $f : S \rightarrow \mathbb{R}$  then  $\nabla f(\mathbf{x}) = \mathbf{0}$ .*

We can now give the second-order necessary conditions for local extrema via the Hessian.

**Theorem 3.2.4.** Let  $f : S \rightarrow \mathbb{R}$  be a function twice differentiable at  $\mathbf{x} \in S$ . If  $\mathbf{x}$  is a local minimum, then  $\mathbf{H}_f(\mathbf{x})$  is positive semidefinite.

*Proof.* Let us assume that  $\mathbf{x}$  is a local minimum. We know from Proposition 3.2.3 that  $\nabla f(\mathbf{x}) = \mathbf{0}$ , hence the second-order expansion at  $\mathbf{x}$  takes the form:

$$f(\mathbf{x} + \lambda \mathbf{d}) = f(\mathbf{x}) + \lambda^2 \frac{1}{2} \mathbf{d}^\top \mathbf{H}_f(\mathbf{x}) \mathbf{d} + o(\lambda^2 \|\mathbf{d}\|_2^2)$$

Because  $\mathbf{x}$  is a local minimum, we can then derive

$$0 \leq \lim_{\lambda \rightarrow 0^+} \frac{f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x})}{\lambda^2} = \frac{1}{2} \mathbf{d}^\top \mathbf{H}_f(\mathbf{x}) \mathbf{d}$$

This is true for any  $\mathbf{d}$ , hence  $\mathbf{H}_f(\mathbf{x})$  is positive semidefinite.  $\square$

**Remark 3.2.5.** Again, for this proposition to hold, it is important that  $S$  is open.

### 3.2.2 A sufficient condition for local extrema

A local minimum thus is a stationary point and has a positive semi-definite Hessian. The converse is almost true, but we need to strengthen the Hessian condition slightly.

**Theorem 3.2.6.** Let  $f : S \rightarrow \mathbb{R}$  be a function twice differentiable at a stationary point  $\mathbf{x} \in S$ . If  $\mathbf{H}_f(\mathbf{x})$  is positive definite then  $\mathbf{x}$  is a local minimum.

*Proof.* Let us assume that  $\mathbf{H}_f(\mathbf{x})$  is positive definite. We know that  $\mathbf{x}$  is a stationary point. We can write the second-order expansion at  $\mathbf{x}$ :

$$f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \frac{1}{2} \boldsymbol{\delta}^\top \mathbf{H}_f(\mathbf{x}) \boldsymbol{\delta} + o(\|\boldsymbol{\delta}\|_2^2)$$

Because the Hessian is positive definite, it has a strictly positive minimum eigenvalue  $\lambda_{\min}$ , we can conclude that  $\boldsymbol{\delta}^\top \mathbf{H}_f(\mathbf{x}) \boldsymbol{\delta} \geq \lambda_{\min} \|\boldsymbol{\delta}\|_2^2$ . From this, we conclude that when  $\|\boldsymbol{\delta}\|_2^2$  is small enough,  $f(\mathbf{x} + \boldsymbol{\delta}) - f(\mathbf{x}) \geq \frac{1}{4} \lambda_{\min} \|\boldsymbol{\delta}\|_2^2 > 0$ . This proves that  $\mathbf{x}$  is a local minimum.  $\square$

**Remark 3.2.7.** When  $\mathbf{H}_f(\mathbf{x})$  is indefinite at a stationary point  $\mathbf{x}$ , we have what is known as a *saddle point*:  $\mathbf{x}$  will be a minimum along the eigenvectors of  $\mathbf{H}_f(\mathbf{x})$  for which the eigenvalues are positive and a maximum along the eigenvectors of  $\mathbf{H}_f(\mathbf{x})$  for which the eigenvalues are negative.

### 3.2.3 Characterization of convexity

**Definition 3.2.8.** For a convex set  $S \subseteq \mathbb{R}^n$ , we say that a function  $f : S \rightarrow \mathbb{R}$  is **strictly convex on  $S$**  if for any two points  $\mathbf{x}_1, \mathbf{x}_2 \in S$  and any  $\theta \in (0, 1)$  we have that:

$$f(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2) < \theta f(\mathbf{x}_1) + (1 - \theta) f(\mathbf{x}_2).$$

**Theorem 3.2.9.** Let  $S \subseteq \mathbb{R}^n$  be open and convex, and let  $f : S \rightarrow \mathbb{R}$  be twice continuously differentiable.

1. If  $H_f(\mathbf{x})$  is positive semi-definite for any  $\mathbf{x} \in S$  then  $f$  is convex on  $S$ .
2. If  $H_f(\mathbf{x})$  is positive definite for any  $\mathbf{x} \in S$  then  $f$  is **strictly** convex on  $S$ .
3. If  $f$  is convex, then  $H_f(\mathbf{x})$  is positive semi-definite  $\forall \mathbf{x} \in S$ .

*Proof.*

1. By applying Theorem 3.2.2, we find that for some  $\mathbf{z} \in [\mathbf{x}, \mathbf{y}]$ :

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

If  $H_f(\mathbf{z})$  is positive semi-definite, this necessarily implies that:

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

and from Theorem 2.3.7 we get that  $f$  is convex.

2. if  $H_f(\mathbf{x})$  is positive definite, we have that:

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

Applying the same idea as in Theorem 2.3.7 we can show that in this case  $f$  is **strictly** convex.

3. Let  $f$  be a convex function. For  $\mathbf{x} \in S$ , and some small  $\lambda > 0$ , for any  $\mathbf{d} \in \mathbb{R}^n$  we have that  $\mathbf{x} + \lambda \mathbf{d} \in S$ . From the Taylor expansion of  $f$  we get:

$$f(\mathbf{x} + \lambda \mathbf{d}) = f(\mathbf{x}) + \lambda \nabla f(\mathbf{x})^\top \mathbf{d} + \frac{\lambda^2}{2} \mathbf{d}^\top H_f(\mathbf{x}) \mathbf{d} + o(\lambda^2 \|\mathbf{d}\|_2^2).$$

From Lemma 2.3.7 we get that if  $f$  is convex then:

$$f(\mathbf{x} + \lambda \mathbf{d}) \geq f(\mathbf{x}) + \lambda \nabla f(\mathbf{x})^\top \mathbf{d}.$$

Therefore, we have that for any  $\mathbf{d} \in \mathbb{R}^n$ :

$$\frac{\lambda^2}{2} \mathbf{d}^\top H_f(\mathbf{x}) \mathbf{d} + o(\|\lambda \mathbf{d}\|^2) \geq 0$$

Dividing by  $\lambda^2$  and taking  $\lambda \rightarrow 0^+$  gives us that for any  $\mathbf{d} \in \mathbb{R}^n$ :  $\mathbf{d}^\top H_f(\mathbf{x}) \mathbf{d} \geq 0$ .  $\square$

**Remark 3.2.10.** It is important to note that if  $S$  is open and  $f$  is strictly convex, then  $H_f(\mathbf{x})$  may still (only) be positive semi-definite  $\forall \mathbf{x} \in S$ . Consider  $f(x) = x^4$  which is strictly convex, then the Hessian is  $H_f(x) = 12x^2$  which equals 0 at  $x = 0$ .

### 3.3 Gradient Descent - An Approach to Optimization?

We have begun to develop an understanding of convex functions, and what we have learned already suggests a way for us to try to find an approximate minimizer of a given convex function.

Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex and differentiable, and we want to solve

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

We would like to find  $\mathbf{x}^*$ , a global minimizer of  $f$ . Suppose we start with some initial guess  $\mathbf{x}_0$ , and we want to update it to  $\mathbf{x}_1$  with  $f(\mathbf{x}_1) < f(\mathbf{x}_0)$ . If we can repeatedly make updates like this, maybe we eventually find a point with nearly minimum function value, i.e. some  $\tilde{\mathbf{x}}$  with  $f(\tilde{\mathbf{x}}) \approx f(\mathbf{x}^*)$ ?

Recall that  $f(\mathbf{x}_0 + \boldsymbol{\delta}) = f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^\top \boldsymbol{\delta} + o(\|\boldsymbol{\delta}\|_2)$ . This means that if we choose  $\mathbf{x}_1 = \mathbf{x}_0 - \lambda \nabla f(\mathbf{x}_0)$ , we get

$$f(\mathbf{x}_0 - \lambda \nabla f(\mathbf{x}_0)) = f(\mathbf{x}_0) - \lambda \|\nabla f(\mathbf{x}_0)\|_2^2 + o(\lambda \|\nabla f(\mathbf{x}_0)\|_2)$$

And because  $f$  is convex, we know that  $\nabla f(\mathbf{x}_0) \neq \mathbf{0}$  unless we are already at a global minimum. So, for some small enough  $\lambda > 0$ , we should get  $f(\mathbf{x}_1) < f(\mathbf{x}_0)$  unless we're already at a global minimizer. This idea of taking a step in the direction of  $-\nabla f(\mathbf{x}_0)$  is what is called *Gradient Descent*. But how do we choose  $\lambda$  each time? And does this lead to an algorithm that quickly reaches a point with close to minimal function value? To get good answers to these questions, we need to assume more about the function  $f$  that we are trying to minimize.

In the following subsection, we will see some conditions that suffice. But there are also many other settings where one can show that some form of gradient descent converges.

#### 3.3.1 A Quantitative Bound on Changes in the Gradient

**Definition 3.3.1.** Let  $f : S \rightarrow \mathbb{R}$  be a differentiable function, where  $S \subseteq \mathbb{R}^n$  is convex and open. We say that  $f$  is  $\beta$ -gradient Lipschitz iff for all  $\mathbf{x}, \mathbf{y} \in S$

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

We also refer to this as  $f$  being  $\beta$ -smooth.

**Proposition 3.3.2.** Consider a twice continuously differentiable  $f : S \rightarrow \mathbb{R}$ . Then  $f$  is  $\beta$ -gradient Lipschitz if and only if for all  $\mathbf{x} \in S$ ,  $\|\mathbf{H}_f(\mathbf{x})\| \leq \beta$ .

You will prove this in Exercise 13 (Week 2) of the first exercise set.

**Proposition 3.3.3.** Let  $f : S \rightarrow \mathbb{R}$  be a  $\beta$ -gradient Lipschitz function. Then for all  $\mathbf{x}, \mathbf{y} \in S$ ,

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

To prove this proposition, we need the following result from multi-variate calculus. This is a restricted form of the fundamental theorem of calculus for line integrals.

**Proposition 3.3.4.** Let  $f : S \rightarrow \mathbb{R}$  be a differentiable function, and consider  $\mathbf{x}, \mathbf{y}$  such that  $[\mathbf{x}, \mathbf{y}] \in S$ . Let  $\mathbf{x}_\theta = \mathbf{x} + \theta(\mathbf{y} - \mathbf{x})$ . Then

$$f(\mathbf{y}) = f(\mathbf{x}) + \int_{\theta=0}^1 \nabla f(\mathbf{x}_\theta)^\top (\mathbf{y} - \mathbf{x}) d\theta$$

Now, we're in a position to show Proposition 3.3.3

*Proof of Proposition 3.3.3.* Let  $f : S \rightarrow \mathbb{R}$  be a  $\beta$ -gradient Lipschitz function. Consider arbitrary  $\mathbf{x}, \mathbf{y} \in S$  such that  $[\mathbf{x}, \mathbf{y}] \in S$

$$\begin{aligned} f(\mathbf{y}) &= f(\mathbf{x}) + \int_{\theta=0}^1 \nabla f(\mathbf{x}_\theta)^\top (\mathbf{y} - \mathbf{x}) d\theta \\ &= f(\mathbf{x}) + \int_{\theta=0}^1 \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) d\theta + \int_{\theta=0}^1 (\nabla f(\mathbf{x}_\theta) - \nabla f(\mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) d\theta \end{aligned}$$

Next we use Cauchy-Schwarz pointwise.

We also evaluate the first integral.

$$\leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \int_{\theta=0}^1 \|\nabla f(\mathbf{x}_\theta) - \nabla f(\mathbf{x})\| \|\mathbf{y} - \mathbf{x}\| d\theta$$

Then we apply  $\beta$ -gradient Lipschitz and note  $\mathbf{x}_\theta - \mathbf{x} = \theta(\mathbf{y} - \mathbf{x})$ .

$$\begin{aligned} &\leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \int_{\theta=0}^1 \beta \theta \|\mathbf{y} - \mathbf{x}\|^2 d\theta. \\ &= f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|^2. \end{aligned}$$

□

### 3.3.2 Analyzing Gradient Descent

It turns out that just knowing a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex and  $\beta$ -gradient Lipschitz is enough to let us figure out a reasonable step size for Gradient Descent and let us analyze its convergence.

We start at a point  $\mathbf{x}_0 \in \mathbb{R}^n$ , and we try to find a point  $\mathbf{x}_1 = \mathbf{x}_0 + \boldsymbol{\delta}$  with lower function value. We will let our upper bound from Proposition 3.3.3 guide us, in fact, we could ask, what is the *best* update for minimizing this upper bound, i.e. a  $\boldsymbol{\delta}$  solving

$$\min_{\boldsymbol{\delta} \in \mathbb{R}^n} f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^\top \boldsymbol{\delta} + \frac{\beta}{2} \|\boldsymbol{\delta}\|^2$$

We can compute the best according to this upper bound by noting first that function is convex and continuously differentiable, and hence will be minimized at any point where the gradient is zero. Thus we want  $\mathbf{0} = \nabla_{\boldsymbol{\delta}} (f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^\top \boldsymbol{\delta} + \frac{\beta}{2} \|\boldsymbol{\delta}\|^2) = \nabla f(\mathbf{x}_0) + \beta \boldsymbol{\delta}$ , which occurs at  $\boldsymbol{\delta} = -\frac{1}{\beta} \nabla f(\mathbf{x}_0)$ .

Plugging in this value into the upper bound, we get that  $f(\mathbf{x}_1) \leq f(\mathbf{x}_0) - \|\nabla f(\mathbf{x}_0)\|_2^2 / 2\beta$ .

Now, as our algorithm, we will start with some guess  $\mathbf{x}_0$ , and then at every step we will update our guess using the best step based on our Proposition 3.3.3 upper bound on  $f$  at  $\mathbf{x}_i$ , and so we get

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \frac{1}{\beta} \nabla f(\mathbf{x}_i) \text{ and } f(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_i) - \frac{\|\nabla f(\mathbf{x}_i)\|_2^2}{2\beta}. \quad (3.1)$$

Let us try to prove that our algorithm converges toward an  $\mathbf{x}$  with low function value.

We will measure this by looking at

$$\text{gap}_i = f(\mathbf{x}_i) - f(\mathbf{x}^*)$$

where  $\mathbf{x}^*$  is a global minimizer of  $f$  (note that there may not be a unique minimizer of  $f$ ). We will try to show that this function value gap grows small. Using  $f(\mathbf{x}_{i+1}) - f(\mathbf{x}_i) = \text{gap}_{i+1} - \text{gap}_i$ , we get

$$\text{gap}_{i+1} - \text{gap}_i \leq -\frac{\|\nabla f(\mathbf{x}_i)\|_2^2}{2\beta} \quad (3.2)$$

If the  $\text{gap}_i$  value is never too much bigger than  $\frac{\|\nabla f(\mathbf{x}_i)\|_2^2}{2\beta}$ , then this should help us show we are making progress. But how much can they differ? We will now try to show a limit on this.

Recall that in the previous chapter we showed the following theorem.

**Theorem 3.3.5.** *Let  $S$  be an open convex subset of  $\mathbb{R}^n$ , and let  $f : S \rightarrow \mathbb{R}$  be a differentiable function. Then,  $f$  is convex if and only if for any  $\mathbf{x}, \mathbf{y} \in S$  we have that  $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$ .*

Using the convexity of  $f$  and the lower bound on convex functions given by Theorem 3.3.5, we have that

$$f(\mathbf{x}^*) \geq f(\mathbf{x}_i) + \nabla f(\mathbf{x}_i)^\top (\mathbf{x}^* - \mathbf{x}_i) \quad (3.3)$$

Rearranging gets us

$$\begin{aligned} \text{gap}_i &\leq \nabla f(\mathbf{x}_i)^\top (\mathbf{x}_i - \mathbf{x}^*) \\ &\leq \|\nabla f(\mathbf{x}_i)\|_2 \|\mathbf{x}_i - \mathbf{x}^*\|_2 \end{aligned} \tag{3.4}$$

by Cauchy-Schwarz.

At this point, we are essentially ready to connect Equation (3.2) with Equation (3.4) and analyze the convergence rate of our algorithm.

However, at the moment, we see that the change  $\text{gap}_{i+1} - \text{gap}_i$  in how close we are to the optimum function value is governed by the norm of the gradient  $\|\nabla f(\mathbf{x}_i)\|_2$ , while the size of the gap is related to *both* this quantity and the distance  $\|\mathbf{x}_i - \mathbf{x}^*\|_2$  between the current solution  $\mathbf{x}_i$  and an optimum  $\mathbf{x}^*$ . Do we need both or can we get rid of, say, the distance? Unfortunately, with this algorithm and for this class of functions, a dependence on the distance is necessary. However, we can simplify things considerably using the following observation, which you will prove in the exercises (Exercise 2):

**Claim 3.3.6.** *When running Gradient Descent as given by the step in Equation (3.1), for all  $i$   $\|\mathbf{x}_i - \mathbf{x}^*\|_2 \leq \|\mathbf{x}_0 - \mathbf{x}^*\|_2$ .*

Combining this Claim with Equation (3.2) and Equation (3.4),

$$\text{gap}_{i+1} - \text{gap}_i \leq -\frac{1}{2\beta} \cdot \left( \frac{\text{gap}_i}{\|\mathbf{x}_0 - \mathbf{x}^*\|_2} \right)^2 \tag{3.5}$$

At this point, a simple induction will complete the proof of following result.

**Theorem 3.3.7.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a  $\beta$ -gradient Lipschitz, convex function. Let  $\mathbf{x}_0$  be a given starting point, and let  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$  be a minimizer of  $f$ . The Gradient Descent algorithm given by*

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \frac{1}{\beta} \nabla f(\mathbf{x}_i)$$

*ensures that the  $k$ th iterate satisfies*

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{2\beta \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{k + 1}.$$

Carrying out this induction is one of the exercises in problem set 2.

# Part II

## Spectral Graph Theory

# Chapter 4

## Introduction to Spectral Graph Theory

In this chapter, we will study graphs through linear algebra. This approach is known as Spectral Graph Theory and turns out to be surprisingly powerful. An in-depth treatment of many topics in this area can be found in [Spi19].

### 4.1 Recap: Incidence and Adjacency Matrices, the Laplacian Matrix and Electrical Energy

To exploit tools from linear algebra, we first have to introduce some interesting matrices that capture important properties of the graph  $G$ . In Chapter 1, we have already seen some of these matrices. Let us recall them and also revisit some additional material from the lectures before. For the rest of the lecture, we let  $G = (V, E, \mathbf{w})$  be an undirected graph, with  $n = |V|$  vertices and  $m = |E|$  edges, where  $\mathbf{w} \in \mathbb{R}_+^E$  assigns a positive weight for every edge. We also assume that  $G$  is connected!

**Fundamental Matrices.** We start by re-introducing the *edge-vertex incidence matrix*  $\mathbf{B}$  of the graph  $G$ . Therefore, we first have to associate an arbitrary direction to every edge. We then let  $\mathbf{B} \in \mathbb{R}^{V \times E}$ .

$$\mathbf{B}(v, e) = \begin{cases} 1 & \text{if } e = (u, v) \\ -1 & \text{if } e = (v, u) \\ 0 & \text{o.w.} \end{cases}$$

The edge directions are only there to help us track the meaning of signs of quantities defined on edges: The math we do should not depend on the choice of sign.

We define the *weight matrix*  $\mathbf{W} \in \mathbb{R}^{E \times E}$  be the diagonal matrix given by  $\mathbf{W} = \text{diag}(\mathbf{w})$ , i.e  $\mathbf{W}(e, e) = \mathbf{w}(e)$ . The weighted degree of a vertex is defined as  $\mathbf{d}(v) = \sum_{\{u,v\} \in E} \mathbf{w}(u, v)$ .

Again we treat the edges as undirected. Let  $\mathbf{D} = \text{diag}(\mathbf{d})$  be the *degree matrix* that is a diagonal matrix in  $\mathbb{R}^{V \times V}$  with weighted degrees on the diagonal.

We define the *weighted adjacency matrix*  $\mathbf{A} \in \mathbb{R}^{V \times V}$  of a graph by

$$\mathbf{A}(u, v) = \begin{cases} \mathbf{w}(u, v) & \text{if } \{u, v\} \in E \\ 0 & \text{otherwise.} \end{cases}$$

Note that we treat the edges as undirected here, so  $\mathbf{A}^\top = \mathbf{A}$ .

Finally, we define the *Laplacian* of the graph as  $\mathbf{L} = \mathbf{B} \mathbf{W} \mathbf{B}^\top$ . Note that in the first chapter, we defined the Laplacian as  $\mathbf{B} \mathbf{R}^{-1} \mathbf{B}^\top$ , where  $\mathbf{R}$  is the diagonal matrix with edge resistances on the diagonal. We want to think of high *weight* on an edge as expressing that two vertices are highly connected, whereas we think of high resistance on an edge as expressing that the two vertices are poorly connected, so we let  $\mathbf{w}(e) = 1/\mathbf{R}(e, e)$ . In Problem Set 1, you showed that  $\mathbf{L} = \mathbf{D} - \mathbf{A}$ , and that for  $\mathbf{x} \in \mathbb{R}^V$ ,

$$\mathbf{x}^\top \mathbf{L} \mathbf{x} = \sum_{\{a,b\} \in E} \mathbf{w}(a, b) (\mathbf{x}(a) - \mathbf{x}(b))^2.$$

**Flows and Voltages.** Let us next briefly review flows and voltages. Given the edge-vertex incidence matrix  $\mathbf{B}$ , we can now express the net flow constraint that  $\mathbf{f}$  routes a demand  $\mathbf{d}$  by

$$\mathbf{B} \mathbf{f} = \mathbf{d}.$$

This is also called a conservation constraint. In our examples so far, we have  $\mathbf{d}(s) = -1$ ,  $\mathbf{d}(t) = 1$  and  $\mathbf{d}(u) = 0$  for all  $u \in V \setminus \{s, t\}$ . Note that, unfortunately, the demand vector is typically denoted by  $\mathbf{d}$ , but so is the degree vector. Since it should always be clear from context which is which, we will also overload  $\mathbf{d}$  in this way!

We then define for any flow  $\mathbf{f} \in \mathbb{R}^E$  the electrical energy by

$$\mathcal{E}(\mathbf{f}) = \sum_e \mathbf{r}(e) \mathbf{f}(e)^2 = \mathbf{f}^\top \mathbf{R} \mathbf{f}$$

and for any electrical voltages  $\mathbf{x} \in \mathbb{R}^E$

$$\mathcal{E}(\mathbf{f}) = \mathbf{f}^\top \mathbf{R} \mathbf{f} = \mathbf{x}^\top \mathbf{L} \mathbf{x}.$$

While the former definition is standard, the second definition is motivated by the fact that by Ohm's law, where we use  $\mathbf{R} = \text{diag}_{e \in E} \mathbf{r}(e)$ , we have that electrical voltages  $\mathbf{x}$  will induce an electrical flow  $\mathbf{f} = \mathbf{R}^{-1} \mathbf{B}^\top \mathbf{x}$ . And for such a pair  $\mathbf{f}$  and  $\mathbf{x}$ , we then have that the energy associated with  $\mathbf{x}$  is equal to the energy consumed by the flow  $\mathbf{f}$  induced by  $\mathbf{x}$ :

$$\mathcal{E}(\mathbf{f}) = \mathbf{f}^\top \mathbf{R} \mathbf{f} = \mathbf{x}^\top \mathbf{L} \mathbf{x} = \mathcal{E}(\mathbf{x}).$$

**The Courant-Fisher Theorem.** Let us also recall the Courant-Fischer theorem, which we proved in Chapter 3 (Theorem 3.1.4). Here we add an additional characterization of  $\lambda_i$ . Note that the second equality for both versions to obtain  $\lambda_i$  can be extracted from our proof of the Courant-Fischer theorem in Chapter 3.

**Theorem 4.1.1** (The Courant-Fischer Theorem, eigenbasis version). *Let  $\mathbf{A}$  be a symmetric matrix in  $\mathbb{R}^{n \times n}$ , with eigenvalues  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ , and corresponding eigenvectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  which form an orthonormal basis. Then*

1.

$$\lambda_i = \min_{\substack{\text{subspace } W \subseteq \mathbb{R}^n \\ \dim(W) = i}} \max_{\substack{\mathbf{x} \in W, \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \min_{\substack{\mathbf{x} \perp \mathbf{x}_1, \dots, \mathbf{x}_{i-1} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$$

2.

$$\lambda_i = \max_{\substack{\text{subspace } W \subseteq \mathbb{R}^n \\ \dim(W) = n+1-i}} \min_{\substack{\mathbf{x} \in W, \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \max_{\substack{\mathbf{x} \perp \mathbf{x}_{i+1}, \dots, \mathbf{x}_n \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$$

Of course, we also have  $\lambda_i(\mathbf{A}) = \frac{\mathbf{x}_i^\top \mathbf{A} \mathbf{x}_i}{\mathbf{x}_i^\top \mathbf{x}_i}$ .

## 4.2 Understanding Eigenvalues of the Laplacian

We would like to understand the eigenvalues of the Laplacian matrix of a graph.

But first, why should we care? It turns out that Laplacian eigenvalues can help us understand many properties of a graph. But we are going to start with a simple motivating observation: Electrical voltages  $\mathbf{x} \in \mathbb{R}^V$  consume electrical energy  $\mathcal{E}(\mathbf{x}) = \mathbf{x}^\top \mathbf{L} \mathbf{x}$ . This means that by the Courant-Fischer Theorem

$$\mathcal{E}(\mathbf{x}) = \mathbf{x}^\top \mathbf{L} \mathbf{x} \leq \lambda_n(L) \mathbf{x}^\top \mathbf{x}$$

And, for any voltages  $\mathbf{x} \perp \mathbf{1}$ ,

$$\mathcal{E}(\mathbf{x}) = \mathbf{x}^\top \mathbf{L} \mathbf{x} \geq \lambda_2(L) \mathbf{x}^\top \mathbf{x}.$$

Thus, we can use the eigenvalues to give upper and lower bounds on how much electrical energy will be consumed by the flow induced by  $\mathbf{x}$ , in terms compared to  $\mathbf{x}^\top \mathbf{x} = \|\mathbf{x}\|_2^2$ .

In a couple of chapters, we will also prove the following claim, which shows that the Laplacian eigenvalues can directly tell us about the electrical energy that is required to route a given demand.

**Claim 4.2.1.** *Given a demand vector  $\mathbf{d} \in \mathbb{R}^V$  such that  $\mathbf{d} \perp \mathbf{1}$ , the electrical voltages  $\mathbf{x}$  that route  $\mathbf{d}$  satisfy  $\mathbf{L} \mathbf{x} = \mathbf{d}$  and the electrical energy of these voltages satisfies*

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

While the techniques in this section can be used to analyze various Laplacian matrices, we focus on understanding the eigenvalues of certain graph families that are particularly instructive. In this section, we obtain the following results.

Graph Family	$\lambda_2$ lower bound	$\lambda_2$ upper bound	$\lambda_n$ lower bound	$\lambda_n$ upper bound
Complete Graph $K_n$	$n$	$n$	$n$	$n$
Path Graph $P_n$	$\frac{1}{n^2}$	$\frac{12}{n^2}$	1	4
Binary Tree $T_n$	$\frac{1}{2n \log_2(n)}$	$\frac{2}{n-1}$	1	6

## 4.3 The Complete Graph

To get a sense of how Laplacian eigenvalues behave, let us start by considering the  $n$ -vertex complete graph with unit weights, which we denote by  $K_n$ . The adjacency matrix of  $K_n$  is  $\mathbf{A} = \mathbf{1}\mathbf{1}^\top - \mathbf{I}$ , since it has ones everywhere, except for the diagonal, where entries are zero. The degree matrix  $\mathbf{D} = (n-1)\mathbf{I}$ . Thus the Laplacian is  $\mathbf{L} = \mathbf{D} - \mathbf{A} = n\mathbf{I} - \mathbf{1}\mathbf{1}^\top$ .

Thus for any  $\mathbf{y} \perp \mathbf{1}$ , we have  $\mathbf{y}^\top \mathbf{L} \mathbf{y} = n\mathbf{y}^\top \mathbf{y} - (\mathbf{1}^\top \mathbf{y})^2 = n\mathbf{y}^\top \mathbf{y}$ .

From this, we can conclude that any  $\mathbf{y} \perp \mathbf{1}$  is an eigenvector of eigenvalue  $n$ , and that all  $\lambda_2 = \lambda_3 = \dots = \lambda_n = n$ .

## 4.4 The Path Graph

Next, let us try to understand  $\lambda_2$  and  $\lambda_n$  for  $P_n$ , the  $n$  vertex path graph with unit weight edges. I.e. the graph has edges  $E = \{\{i, i+1\} \text{ for } i = 1 \text{ to } (n-1)\}$ . This is in a sense the least well-connected unit weight graph on  $n$  vertices, whereas  $K_n$  is the most well-connected.

### 4.4.1 Upper Bounding $\lambda_2(P_n)$ via Test Vectors

We can use the eigenbasis version of the Courant-Fisher theorem to observe that the second-smallest eigenvalue of the Laplacian is given by

$$\lambda_2(\mathbf{L}) = \min_{\substack{\mathbf{x} \neq \mathbf{0} \\ \mathbf{x}^\top \mathbf{1} = 0}} \frac{\mathbf{x}^\top \mathbf{L} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}. \quad (4.1)$$

We can get a better understanding of this particular case through a couple of simple observations. Suppose  $\mathbf{x} = \mathbf{y} + \alpha \mathbf{1}$ , where  $\mathbf{y} \perp \mathbf{1}$ . Then  $\mathbf{x}^\top \mathbf{L} \mathbf{x} = \mathbf{y}^\top \mathbf{L} \mathbf{y}$ , and  $\|\mathbf{x}\|_2^2 = \|\mathbf{y}\|_2^2 + \alpha^2 \|\mathbf{1}\|_2^2$ .

So for any given vector, you can increase the value of  $\frac{\mathbf{x}^\top \mathbf{Lx}}{\mathbf{x}^\top \mathbf{x}}$ , by replacing  $\mathbf{x}$  with its component orthogonal to  $\mathbf{1}$ , which we denoted by  $\mathbf{y}$ .

We can conclude from Equation (4.1) that for *any* vector  $\mathbf{y} \perp \mathbf{1}$ ,

$$\lambda_2 \leq \frac{\mathbf{y}^\top \mathbf{Ly}}{\mathbf{y}^\top \mathbf{y}}$$

When we use a vector  $\mathbf{y}$  in this way to prove a bound on an eigenvalue, we call it a *test vector*.

Now, we'll use a test vector to give an upper bound on  $\lambda_2(\mathbf{L}_{P_n})$ . Let  $\mathbf{x} \in \mathbb{R}^V$  be given by  $\mathbf{x}(i) = (n+1) - 2i$ , for  $i \in [n]$ . This vector satisfies  $\mathbf{x} \perp \mathbf{1}$ . We picked this because we wanted a sequence of values growing linearly along the path, while also making sure that the vector is orthogonal to  $\mathbf{1}$ . Now

$$\begin{aligned} \lambda_2(\mathbf{L}_{P_n}) &\leq \frac{\sum_{i \in [n-1]} (\mathbf{x}(i) - \mathbf{x}(i+1))^2}{\sum_{i=1}^n \mathbf{x}(i)^2} \\ &= \frac{\sum_{i=1}^{n-1} 2^2}{\sum_{i=1}^n (n+1-2i)^2} \\ &= \frac{4(n-1)}{(n+1)n(n-1)/3} \\ &= \frac{12}{n(n+1)} \leq \frac{12}{n^2}. \end{aligned}$$

Later, we will prove a lower bound that shows this value is right up to a constant factor. But the test vector approach based on the Courant-Fischer theorem doesn't immediately work when we want to prove lower bounds on  $\lambda_2(\mathbf{L})$ .

#### 4.4.2 Lower Bounding $\lambda_n(P_n)$ via Test Vectors

We can see from either version of the Courant-Fischer theorem that

$$\lambda_n(\mathbf{L}) = \max_{\mathbf{v} \neq \mathbf{0}} \frac{\mathbf{v}^\top \mathbf{Lv}}{\mathbf{v}^\top \mathbf{v}}. \quad (4.2)$$

Thus for *any* vector  $\mathbf{y} \neq 0$ ,

$$\lambda_n \geq \frac{\mathbf{y}^\top \mathbf{Ly}}{\mathbf{y}^\top \mathbf{y}}.$$

This means we can get a test vector-based lower bound on  $\lambda_n$ . Let us apply this to the Laplacian of  $P_n$ . We'll try the vector  $\mathbf{x} \in \mathbb{R}^V$  given by  $\mathbf{x}(1) = -1$ , and  $\mathbf{x}(n) = 1$  and  $\mathbf{x}(i) = 0$  for  $i \neq 1, n$ .

Here we get

$$\lambda_n(\mathbf{L}_{P_n}) \geq \frac{\mathbf{x}^\top \mathbf{L} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{2}{2} = 1. \quad (4.3)$$

Again, it's not clear how to use the Courant-Fischer theorem to prove an upper bound on  $\lambda_n(\mathbf{L})$ . But, later we'll see how to prove an upper that shows that for  $P_n$ , the lower bound we obtained is right up to constant factors.

### 4.4.3 Upper bounding $\lambda_n(P_n)$ via the Loewner Order

In the previous sections, we first saw a complete characterization of the eigenvalues and eigenvectors of the unit weight complete graph on  $n$  vertices,  $K_n$ . Namely,  $\mathbf{L}_{K_n} = n\mathbf{I} - \mathbf{1}\mathbf{1}^\top$ , and this means that *every* vector  $\mathbf{y} \perp \mathbf{1}$  is an eigenvector of eigenvalue  $n$ . Ideally, we would like to prove an almost matching lower bound on  $\lambda_2(P_n)$  and an almost matching upper bound on  $\lambda_n(P_n)$ , but it is not clear how to get that from the Courant-Fischer theorem.

To get there, we need to introduce some more tools. We now introduce an ordering on symmetric matrices called the *Loewner order*, which I also like to just call the positive semi-definite order. As we will see in a moment, it is a partial order on symmetric matrices, we denote it by “ $\preceq$ ”. For convenience, we allow ourselves to both write  $\mathbf{A} \preceq \mathbf{B}$  and equivalently  $\mathbf{B} \succeq \mathbf{A}$ .

**The Loewner Order, aka. the Positive Semi-Definite Order.** For a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  we define that

$$\mathbf{A} \succeq \mathbf{0}$$

if and only if  $\mathbf{A}$  is positive semi-definite.

More generally, when we have two symmetric matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ , we will write

$$\mathbf{A} \preceq \mathbf{B} \text{ if and only if for all } \mathbf{x} \in \mathbb{R}^n \text{ we have } \mathbf{x}^\top \mathbf{A} \mathbf{x} \leq \mathbf{x}^\top \mathbf{B} \mathbf{x} \quad (4.4)$$

This is a partial order, because it satisfies the three requirements of

1. Reflexivity:  $\mathbf{A} \preceq \mathbf{A}$ .
2. Anti-symmetry:  $\mathbf{A} \preceq \mathbf{B}$  and  $\mathbf{B} \preceq \mathbf{A}$  implies  $\mathbf{A} = \mathbf{B}$
3. Transitivity:  $\mathbf{A} \preceq \mathbf{B}$  and  $\mathbf{B} \preceq \mathbf{C}$  implies  $\mathbf{A} \preceq \mathbf{C}$

Check for yourself that these properties hold!

The PSD order has other very useful properties:  $\mathbf{A} \preceq \mathbf{B}$  implies  $\mathbf{A} + \mathbf{C} \preceq \mathbf{B} + \mathbf{C}$  for any symmetric matrix  $\mathbf{C}$ . Convince yourself of this too!

And, combining this observation with transitivity, we can see that  $\mathbf{A} \preceq \mathbf{B}$  and  $\mathbf{C} \preceq \mathbf{D}$  implies  $\mathbf{A} + \mathbf{C} \preceq \mathbf{B} + \mathbf{D}$ .

Here is another useful property: If  $\mathbf{0} \preceq \mathbf{A}$  then for all  $\alpha \geq 1$

$$\frac{1}{\alpha} \mathbf{A} \preceq \mathbf{A} \preceq \alpha \mathbf{A}.$$

Here is another one:

**Claim 4.4.1.** *If  $\mathbf{A} \preceq \mathbf{B}$ , then for all  $i$*

$$\lambda_i(\mathbf{A}) \leq \lambda_i(\mathbf{B}).$$

*Proof.* We can prove this Claim by applying the subspace version of the Courant-Fischer theorem.

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

□

Note that the converse of Clam 4.4.1 is very much false, for example the matrices  $\mathbf{A} = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}$  and  $\mathbf{B} = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$  have equal eigenvalues, but both  $\mathbf{A} \not\preceq \mathbf{B}$  and  $\mathbf{B} \not\preceq \mathbf{A}$ .

**Remark 4.4.2.** It's useful to get used to and remember some of the properties of the Loewner order, but all the things we have established so far are almost immediate from the basic characterization in Equation (4.4). So, ideally, don't memorize all these facts, instead, try to see that they are simple consequences of the definition.

**Upper bounding  $\lambda_n(L_G)$  Using Degrees.** In an earlier chapter, we observed that for any graph  $G = (V, E, \mathbf{w})$ ,  $\mathbf{L} = \mathbf{D} - \mathbf{A} \succeq \mathbf{0}$ . We can see this from  $\mathbf{x}^\top (\mathbf{D} - \mathbf{A}) \mathbf{x} = \sum_{(u,v) \in E} \mathbf{w}(u,v)(\mathbf{x}(u) - \mathbf{x}(v))^2 \geq 0$ . Similarly  $\mathbf{D} + \mathbf{A} \succeq \mathbf{0}$  because  $\mathbf{x}^\top (\mathbf{D} + \mathbf{A}) \mathbf{x} = \sum_{(u,v) \in E} \mathbf{w}(u,v)(\mathbf{x}(u) + \mathbf{x}(v))^2 \geq 0$ . But this means that  $-\mathbf{A} \preceq \mathbf{D}$  and hence  $\mathbf{L} = \mathbf{D} - \mathbf{A} \preceq 2\mathbf{D}$ .

So, for the path graph  $P_n$ , we have  $\mathbf{L}_{P_n} \preceq \mathbf{D} - \mathbf{A} \preceq 2\mathbf{D} \preceq 4\mathbf{I}$ . So by Claim 4.4.1

$$\lambda_n(\mathbf{L}_{P_n}) \leq 4. \tag{4.5}$$

We can see that our test vector-based lower bound on  $\lambda_n(\mathbf{L}_{P_n})$  from Equation (4.3) is tight up to a factor 4.

Since this type of argument works for any unit weight graph, it proves the following claim.

**Claim 4.4.3.** *For any unit weight graph  $G$ ,  $\lambda_n(\mathbf{L}_G) \leq 2 \max_{v \in V} \mathbf{d}(v)$ .*

This is tight on a graph consisting of a single edge.

#### 4.4.4 Lower bounding $\lambda_2(P_n)$ via the Loewner Order

We will now conclude our analysis of the eigenvalues  $\lambda_2$  and  $\lambda_n$  of the path graph  $P_n$  by finally deriving a lower bound on  $\lambda_2(P_n)$ . Again, this requires us to use new tools. As in the last section, we first introduce the new tool for general graphs and then exploit the tool to derive the lower bound.

**The Loewner Order for Graphs.** It's sometimes convenient to overload the PSD order to also apply to graphs. We will write

$$G \preceq H$$

if  $\mathbf{L}_G \preceq \mathbf{L}_H$ .

For example, given two unit weight graphs  $G = (V, E)$  and  $H = (V, F)$ , if  $H$  is a subgraph of  $G$ , then

$$\mathbf{L}_H \preceq \mathbf{L}_G.$$

We can see this from the Laplacian quadratic form:

$$\mathbf{x}^\top \mathbf{L}_G \mathbf{x} = \sum_{(u,v) \in E} \mathbf{w}(u,v)(\mathbf{x}(u) - \mathbf{x}(v))^2.$$

Dropping edges will only decrease the value of the quadratic form. The same is for decreasing the weights of edges. The graph order notation is especially useful when we allow for scaling a graph by a constant, say  $c > 0$ ,

$$c \cdot H \preceq G$$

What is  $c \cdot H$ ? It is the same graph as  $H$ , but the weight of every edge is multiplied by  $c$ . Now we can make statements like  $\frac{1}{2}H \preceq G \preceq 2H$ , which turn out to be useful notion of the two graphs approximating each other.

**The Path Inequality.** Now, we'll see a general tool for comparing two graphs  $G$  and  $H$  to prove inequalities like  $cH \preceq G$  for some constant  $c$ . Our tools won't necessarily work well for all cases, but we'll see some examples where they do.

In the rest of the chapter, we will often need to compare two graphs defined on the same vertex set  $V = \{1, \dots, n\} = [n]$ .

We use  $G_{i,j}$  to denote the unit weight graph on vertex set  $[n]$  consisting of a single edge between vertices  $i$  and  $j$ .

**Lemma 4.4.4** (The Path Inequality).

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

*Proof.* We want to show that for every  $\mathbf{x} \in \mathbb{R}^n$ ,

$$(n-1) \cdot \sum_{i=1}^{n-1} (\mathbf{x}(i+1) - \mathbf{x}(i))^2 \geq (\mathbf{x}(n) - \mathbf{x}(1))^2.$$

For  $i \in [n-1]$ , set

$$\Delta(i) = \mathbf{x}(i+1) - \mathbf{x}(i).$$

The inequality we want to prove then becomes

$$(n-1) \sum_{i=1}^{n-1} (\Delta(i))^2 \geq \left( \sum_{i=1}^{n-1} \Delta(i) \right)^2.$$

But, this is immediate from the Cauchy-Schwarz inequality  $\mathbf{a}^\top \mathbf{b} \leq \|\mathbf{a}\|_2 \|\mathbf{b}\|_2$ :

$$\begin{aligned} (n-1) \sum_{i=1}^{n-1} (\Delta(i))^2 &= \|\mathbf{1}_{n-1}\|^2 \cdot \|\Delta\|^2 \\ &= (\|\mathbf{1}_{n-1}\| \cdot \|\Delta\|)^2 \\ &\geq (\mathbf{1}_{n-1}^\top \Delta)^2 \\ &= \left( \sum_{i=1}^{n-1} \Delta(i) \right)^2 \end{aligned}$$

□

**Lower Bounding  $\lambda_2$  of a Path Graph.** We will now use Lemma 4.4.4 to prove a lower bound on  $\lambda_2(\mathbf{L}_{P_n})$ . Our strategy will be to prove that the path  $P_n$  is at least some multiple of the complete graph  $K_n$ , measured by the Loewner order, i.e.  $K_n \preceq f(n)P_n$  for some function  $f : \mathbb{N} \rightarrow \mathbb{R}$ . We can combine this with our observation earlier that  $\lambda_2(\mathbf{L}_{K_n}) = n$  to show that

$$f(n)\lambda_2(\mathbf{L}_{P_n}) \geq \lambda_2(\mathbf{L}_{K_n}) = n, \quad (4.6)$$

and this will give our lower bound on  $\lambda_2(\mathbf{L}_{P_n})$ . When establishing the inequality between  $P_n$  and  $K_n$ , we can treat each edge of the complete graph separately, by first noting that

$$\mathbf{L}_{K_n} = \sum_{i < j} \mathbf{L}_{G_{i,j}}$$

For every edge  $(i, j)$  in the complete graph, we apply the Path Inequality, Lemma 4.4.4:

$$\begin{aligned} G_{i,j} &\preceq (j-i) \sum_{k=i}^{j-1} G_{k,k+1} \\ &\preceq (j-i)P_n \end{aligned}$$

This inequality says that  $G_{i,j}$  is at most  $(j - i)$  times the part of the path connecting  $i$  to  $j$ , and that this part of the path is less than the whole.

Summing inequality (4.3) over all edges  $(i, j) \in K_n$  gives

$$K_n = \sum_{i < j} G_{i,j} \preceq \sum_{i < j} (j - i) P_n.$$

To finish the proof, we compute

$$\sum_{i < j} (j - i) \leq \sum_{i < j} n \leq n^3$$

So

$$L_{K_n} \preceq n^3 \cdot L_{P_n}.$$

Plugging this into Equation (4.6) we obtain

$$\frac{1}{n^2} \leq \lambda_2(P_n).$$

This only differs from our test vector-based upper bound in Equation (4.3) by a factor 12.

We could make this considerably tighter by being more careful about the sums.

## 4.5 The Complete Binary Tree

To consolidate our new expertise, let us finally consider the complete binary tree  $T_n$  with unit weight edges on  $n$  vertices. Note that in order to have  $T_n$  be a complete binary tree, we require that  $n$  is such that there is an integer  $d$  for which  $n = 2^{d+1} - 1$ .

$T_n$  is the balanced binary tree on this many vertices, i.e. it consists of a root node, which has two children, each of those children have two children and so on until we reach a depth of  $d$  from the root, at which point the child vertices have no more children. A simple induction shows that indeed  $n = 2^{d+1} - 1$ .

We can also describe the edge set by saying that each node  $i$  has edges to its children  $2i$  and  $2i + 1$  whenever the node labels do not exceed  $n$ . We emphasize that we still think of the graph as undirected.

### 4.5.1 Deriving Bounds on $\lambda_n(T_n)$

We'll start by above bounding  $\lambda_n(L_{T_n})$  using a test vector.

We let  $\mathbf{x}(i) = 0$  for all nodes that have a child node, and  $\mathbf{x}(i) = -1$  for even-numbered leaf nodes and  $\mathbf{x}(i) = +1$  for odd-numbered leaf nodes. Note that there are  $(n+1)/2$  leaf nodes, and every leaf node has a single edge, connecting it to a parent with value 0. Thus

$$\lambda_n(\mathbf{L}) = \max_{\mathbf{v} \neq \mathbf{0}} \frac{\mathbf{v}^\top \mathbf{L} \mathbf{v}}{\mathbf{v}^\top \mathbf{v}} \geq \frac{\mathbf{x}^\top \mathbf{L} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{(n+1)/2}{(n+1)/2} = 1. \quad (4.7)$$

Meanwhile, every vertex has degree at most 3, so by Claim 4.4.3,  $\lambda_n(\mathbf{L}) \leq 6$ . So we can bound the largest eigenvalue above and below by a constant.

### 4.5.2 Upper Bounding $\lambda_2$ via Test Vectors

Let us give an upper bound on  $\lambda_2$  of the tree using a test vector. Let  $\mathbf{x} \in \mathbb{R}^v$  have  $\mathbf{x}(1) = 0$  and  $\mathbf{x}(i) = -1$  for  $i$  in the left subtree and  $\mathbf{x}(i) = +1$  in the right subtree. Then

$$\lambda_2(\mathbf{L}_{T_n}) = \min_{\substack{\mathbf{v} \neq \mathbf{0} \\ \mathbf{v}^\top \mathbf{1} = 0}} \frac{\mathbf{v}^\top \mathbf{L} \mathbf{v}}{\mathbf{v}^\top \mathbf{v}} \leq \frac{\mathbf{x}^\top \mathbf{L} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{2}{n-1}.$$

So, we have shown  $\frac{1}{2n \log_2(n)} \leq \lambda_2(\mathbf{L}_{T_n}) \leq \frac{2}{n-1}$ , and unlike the previous examples, the gap is more than a constant.

### 4.5.3 Lower Bounding $\lambda_2$

The following lemma gives a simple lower bound on  $\lambda_2$  for any graph.

**Lemma 4.5.1.** *For any unit weight graph  $G$  with diameter  $D$ ,*

$$\lambda_2(\mathbf{L}_G) \geq \frac{1}{nD}.$$

*Proof.* We will again prove a lower bound comparing  $G$  to the complete graph. For each edge  $(i, j) \in K_n$ , let  $G^{i,j}$  denote a shortest path in  $G$  from  $i$  to  $j$ . This path will have length at most  $D$ . So, we have

$$\begin{aligned} K_n &= \sum_{i < j} G_{i,j} \\ &\preceq \sum_{i < j} DG^{i,j} \\ &\preceq \sum_{i < j} DG \\ &\preceq n^2 DG. \end{aligned}$$

So, we obtain the bound

$$n^2 D \lambda_2(G) \geq n,$$

which implies our desired statement.  $\square$

Since the complete binary tree  $T_n$  has diameter  $2d \leq 2\log_2(n)$ , by Lemma 4.5.1,  $\lambda_2(\mathbf{L}_{T_n}) \geq \frac{1}{2n \log_2(n)}$ .

In the exercises for Week 3, I will ask you to improve the lower bound to  $1/(cn)$  for some constant  $c$ .

# Chapter 5

## Pseudo-inverses and Effective Resistance

### 5.1 What is a (Moore-Penrose) Pseudoinverse?

Recall that for a connected graph  $G$  with Laplacian  $\mathbf{L}$ , we have  $\ker(\mathbf{L}) = \text{span}\{\mathbf{1}\}$ , which means  $\mathbf{L}$  is not invertible. However, we still want some matrix which behaves like a real inverse. To be more specific, given a Laplacian  $\mathbf{L} \in \mathbb{R}^{V \times V}$ , we want some matrix  $\mathbf{L}^+ \in \mathbb{R}^{V \times V}$  s.t.

- 1)  $(\mathbf{L}^+)^T = \mathbf{L}^+$  (symmetric)
- 2)  $\mathbf{L}^+ \mathbf{1} = \mathbf{0}$ , or more generally,  $\mathbf{L}^+ \mathbf{v} = \mathbf{0}$  for  $\mathbf{v} \in \ker(\mathbf{L})$
- 3)  $\mathbf{L}^+ \mathbf{L} \mathbf{v} = \mathbf{L} \mathbf{L}^+ \mathbf{v} = \mathbf{v}$  for  $\mathbf{v} \perp \mathbf{1}$ , or more generally, for  $\mathbf{v} \in \ker(\mathbf{L})^\perp$

Under the above conditions,  $\mathbf{L}^+$  is uniquely defined and we call it the pseudoinverse of  $\mathbf{L}$ . Note that there are many other equivalent definitions of the pseudoinverse of some matrix  $\mathbf{A}$ , and we can also generalize the concept to matrices that aren't symmetric or even square.

Let  $\lambda_i, \mathbf{v}_i$  be the  $i$ -th pair of eigenvalue and eigenvector of  $\mathbf{L}$ , with  $\{\mathbf{v}_i\}_{i=1}^n$  forming a orthogonal basis. Then by the spectral theorem,

$$\mathbf{L} = \mathbf{V} \mathbf{A} \mathbf{V}^\top = \sum_i \lambda_i \mathbf{v}_i \mathbf{v}_i^\top,$$

where  $\mathbf{V} = [\mathbf{v}_1 \ \cdots \ \mathbf{v}_n]$  and  $\mathbf{A} = \text{diag}\{\lambda_1, \dots, \lambda_n\}$ . And we can show that its pseudoinverse is exactly

$$\mathbf{L}^+ = \sum_{i, \lambda_i \neq 0} \lambda_i^{-1} \mathbf{v}_i \mathbf{v}_i^\top.$$

Checking conditions 1), 2), 3) is immediate. We can also prove uniqueness, but this takes slightly more work.

## 5.2 Electrical Flows Again

Recall the incidence matrix  $\mathbf{B} \in \mathbb{R}^{V \times E}$  of a graph  $G = (V, E)$ .

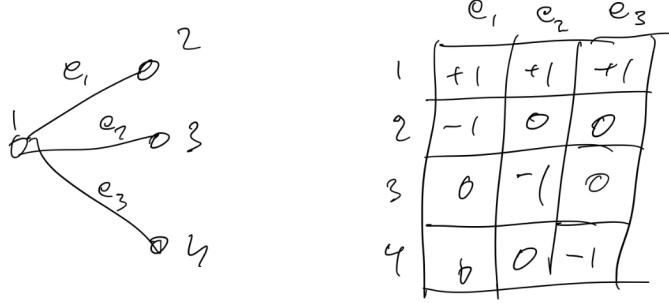


Figure 5.1: An example of a graph and its incidence matrix  $B$ .

In Chapter 1, we introduced the electrical flow routing demand  $\mathbf{d} \in \mathbb{R}^V$ . Let's call the electrical flow  $\tilde{\mathbf{f}} \in \mathbb{R}^E$ . The net flow constraint requires  $\mathbf{B}\tilde{\mathbf{f}} = \mathbf{d}$ . By Ohm's Law,  $\tilde{\mathbf{f}} = \mathbf{R}^{-1}\mathbf{B}^\top \mathbf{x}$  for some voltage  $\mathbf{x} \in \mathbb{R}^V$  where  $\mathbf{R} = \text{diag}(\mathbf{r})$  and  $\mathbf{r}(e) = \text{resistance of edge } e$ . We showed (in the exercises) that when  $\mathbf{d} \perp \mathbf{1}$ , there exists a voltage  $\tilde{\mathbf{x}} \perp \mathbf{1}$  s.t.  $\tilde{\mathbf{f}} = \mathbf{R}^{-1}\mathbf{B}^\top \tilde{\mathbf{x}}$  and  $\mathbf{B}\tilde{\mathbf{f}} = \mathbf{d}$ . This  $\tilde{\mathbf{x}}$  solves  $\mathbf{L}\mathbf{x} = \mathbf{d}$  where  $\mathbf{L} = \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^\top$ .

And we also made the following claim.

**Claim 5.2.1.**

$$\tilde{\mathbf{f}} = \arg \min_{\mathbf{B}\mathbf{f}=\mathbf{d}} \mathbf{f}^\top \mathbf{R}\mathbf{f} \text{ where } \mathbf{f}^\top \mathbf{R}\mathbf{f} = \sum_e \mathbf{r}(e)\mathbf{f}(e)^2, \quad (5.1)$$

You proved this in the exercises for Week 1. Let's recap the proof briefly, just to get back into thinking about electrical flows.

*Proof.* Consider any  $\mathbf{f} \in \mathbb{R}^E$  s.t.  $\mathbf{B}\mathbf{f} = \mathbf{d}$ . For any  $\mathbf{x} \in \mathbb{R}^V$ , we have

$$\begin{aligned} \frac{1}{2}\mathbf{f}^\top \mathbf{R}\mathbf{f} &= \frac{1}{2}\mathbf{f}^\top \mathbf{R}\mathbf{f} - \mathbf{x}^\top (\underbrace{\mathbf{B}\mathbf{f} - \mathbf{d}}_0) \\ &\geq \min_{\mathbf{f} \in \mathbb{R}^E} \underbrace{\frac{1}{2}\mathbf{f}^\top \mathbf{R}\mathbf{f} - \mathbf{x}^\top \mathbf{B}\mathbf{f} + \mathbf{d}^\top \mathbf{x}}_{g(\mathbf{f})} \\ &= \mathbf{d}^\top \mathbf{x} - \frac{1}{2}\mathbf{x}^\top \mathbf{L}\mathbf{x} \end{aligned}$$

since  $\nabla_{\mathbf{f}} g(\mathbf{f}) = \mathbf{0}$  gives us  $\mathbf{f} = \mathbf{R}^{-1}\mathbf{B}^\top \mathbf{x}$ . Thus, for all  $\mathbf{f} \in \mathbb{R}^E$  s.t.  $\mathbf{B}\mathbf{f} = \mathbf{d}$  and all  $\mathbf{x} \in \mathbb{R}^V$ ,

$$\frac{1}{2}\mathbf{f}^\top \mathbf{R}\mathbf{f} \geq \mathbf{d}^\top \mathbf{x} - \frac{1}{2}\mathbf{x}^\top \mathbf{L}\mathbf{x}. \quad (5.2)$$

But for the electrical flow  $\tilde{\mathbf{f}}$  and electrical voltage  $\tilde{\mathbf{x}}$ , we have  $\tilde{\mathbf{f}} = \mathbf{R}^{-1}\mathbf{B}^\top \tilde{\mathbf{x}}$  and  $\mathbf{L}\tilde{\mathbf{x}} = \mathbf{d}$ .

So

$$\tilde{\mathbf{f}}^\top \mathbf{R}\tilde{\mathbf{f}} = (\mathbf{R}^{-1}\mathbf{B}^\top \tilde{\mathbf{x}})^\top \mathbf{R}(\mathbf{R}^{-1}\mathbf{B}^\top \tilde{\mathbf{x}}) = \tilde{\mathbf{x}}^\top \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^\top \tilde{\mathbf{x}} = \tilde{\mathbf{x}}^\top \mathbf{L}\tilde{\mathbf{x}} = \tilde{\mathbf{x}}^\top \mathbf{d}.$$

Therefore,

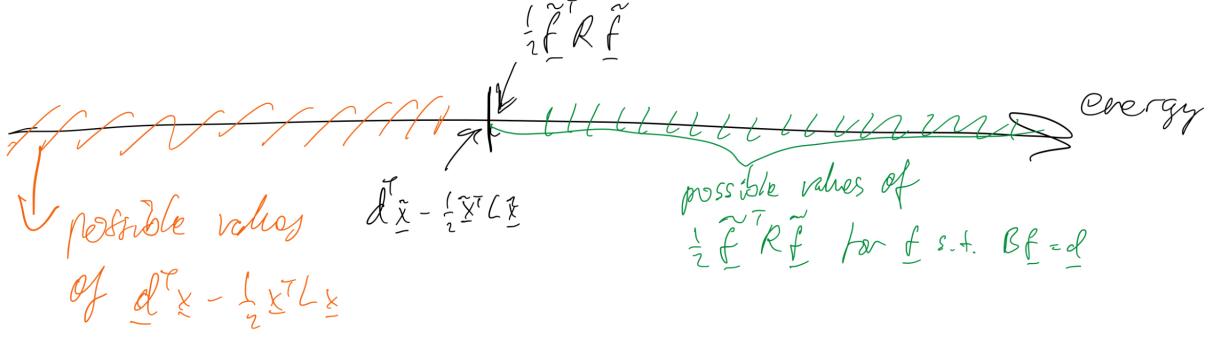
$$\frac{1}{2}\tilde{\mathbf{f}}^\top \mathbf{R}\tilde{\mathbf{f}} = \mathbf{d}^\top \tilde{\mathbf{x}} - \frac{1}{2}\tilde{\mathbf{x}}^\top \mathbf{L}\tilde{\mathbf{x}}. \quad (5.3)$$

By combining Equation (5.2) and Equation (5.3), we see that for all  $\mathbf{f}$  s.t.  $\mathbf{B}\mathbf{f} = \mathbf{d}$ ,

$$\frac{1}{2}\mathbf{f}^\top \mathbf{R}\mathbf{f} \geq \mathbf{d}^\top \tilde{\mathbf{x}} - \frac{1}{2}\tilde{\mathbf{x}}^\top \mathbf{L}\tilde{\mathbf{x}} = \frac{1}{2}\tilde{\mathbf{f}}^\top \mathbf{R}\tilde{\mathbf{f}}.$$

Thus  $\tilde{\mathbf{f}}$  is the minimum electrical energy flow among all flows that route demand  $\mathbf{d}$ , proving Equation (5.1) holds.

The drawing below shows how the quantities line up:



□

### 5.3 Effective Resistance

Given a graph  $G = (V, E)$ , for any pair of vertices  $(a, b) \in V$ , we want to compute the cost (or energy) of routing 1 unit of current from  $a$  to  $b$ . We call such cost the effective resistance between  $a$  and  $b$ , denoted by  $R_{\text{eff}}(a, b)$ . Recall for a single resistor  $r(a, b)$ ,

$$\text{energy} = r(a, b)f^2(a, b) = r(a, b).$$

So when we have a graph consisting of just one edge  $(a, b)$ , the effective resistance is just  $R_{\text{eff}}(a, b) = r(a, b)$ .

In a general graph, we can also consider the energy required to route one unit of current between two vertices. For any pair  $a, b \in V$ , we have

$$R_{\text{eff}}(a, b) = \min_{\mathbf{B}\mathbf{f} = \mathbf{e}_b - \mathbf{e}_a} \mathbf{f}^\top \mathbf{R}\mathbf{f},$$

where  $\mathbf{e}_v \in \mathbb{R}^V$  is the indicator vector of  $v$ . Note that the cost of routing  $F$  units of flow from  $a$  to  $b$  will be  $R_{\text{eff}}(a, b) \cdot F^2$ .

Since  $(\mathbf{e}_b - \mathbf{e}_a)^\top \mathbf{1} = 0$ , we know from the previous section that  $R_{\text{eff}}(a, b) = \tilde{\mathbf{f}}^\top \mathbf{R} \tilde{\mathbf{f}}$  where  $\tilde{\mathbf{f}}$  is the electrical flow. Now we can write  $\mathbf{L}\tilde{\mathbf{x}} = \mathbf{e}_b - \mathbf{e}_a$  and  $\tilde{\mathbf{x}} = \mathbf{L}^+(\mathbf{e}_b - \mathbf{e}_a)$  for the electrical voltages routing 1 unit of current from  $a$  to  $b$ . Now the energy of routing 1 unit of current from  $a$  to  $b$  is

$$R_{\text{eff}}(a, b) = \tilde{\mathbf{f}}^\top \mathbf{R} \tilde{\mathbf{f}} = \tilde{\mathbf{x}}^\top \mathbf{L} \tilde{\mathbf{x}} = (\mathbf{e}_b - \mathbf{e}_a)^\top \mathbf{L}^+ \mathbf{L} \mathbf{L}^+ (\mathbf{e}_b - \mathbf{e}_a) = (\mathbf{e}_b - \mathbf{e}_a)^\top \mathbf{L}^+ (\mathbf{e}_b - \mathbf{e}_a),$$

where the last equality is due to  $\mathbf{L}^+ \mathbf{L} \mathbf{L}^+ = \mathbf{L}^+$ .

**Remark 5.3.1.** We have now seen several different expressions that all take on the same value: the energy of the electrical flow. It's useful to remind yourself what these are. Consider an electrical flow  $\tilde{\mathbf{f}}$  routes demand  $\mathbf{d}$ , and associated electrical voltages  $\tilde{\mathbf{x}}$ . We know that  $\mathbf{B}\tilde{\mathbf{f}} = \mathbf{d}$ , and  $\mathbf{f} = \mathbf{R}^{-1}\mathbf{B}^\top \tilde{\mathbf{x}}$ , and  $\mathbf{L}\tilde{\mathbf{x}} = \mathbf{d}$ , where  $\mathbf{L} = \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^\top$ . And we have seen how to express the electrical energy using many different quantities:

$$\tilde{\mathbf{f}}^\top \mathbf{R} \tilde{\mathbf{f}} = \tilde{\mathbf{x}}^\top \mathbf{L} \tilde{\mathbf{x}} = \mathbf{d}^\top \mathbf{L}^+ \mathbf{d} = \mathbf{d}^\top \tilde{\mathbf{x}} = \tilde{\mathbf{f}}^\top \mathbf{B}^\top \tilde{\mathbf{x}}$$

**Claim 5.3.2.** Any PSD matrix  $\mathbf{A}$  has a PSD square root  $\mathbf{A}^{1/2}$  s.t.  $\mathbf{A}^{1/2} \mathbf{A}^{1/2} = \mathbf{A}$ .

*Proof.* By the spectral theorem,  $\mathbf{A} = \sum_i \lambda_i \mathbf{v}_i \mathbf{v}_i^\top$  where  $\{\mathbf{v}_i\}$  are orthonormal. Let  $\mathbf{A}^{1/2} = \sum_i \lambda_i^{1/2} \mathbf{v}_i \mathbf{v}_i^\top$ . Then

$$\begin{aligned} \mathbf{A}^{1/2} \mathbf{A}^{1/2} &= \left( \sum_i \lambda_i^{1/2} \mathbf{v}_i \mathbf{v}_i^\top \right)^2 \\ &= \sum_i \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \mathbf{v}_i \mathbf{v}_i^\top + \sum_{i \neq j} \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \mathbf{v}_j \mathbf{v}_j^\top \\ &= \sum_i \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \end{aligned}$$

where the last equality is due to  $\mathbf{v}_i^\top \mathbf{v}_j = \delta_{ij}$ . It's easy to see that  $\mathbf{A}^{1/2}$  is also PSD.  $\square$

Let  $\mathbf{L}^{+/2}$  be the square root of  $\mathbf{L}^+$ . So

$$R_{\text{eff}}(a, b) = (\mathbf{e}_b - \mathbf{e}_a)^\top \mathbf{L}^+ (\mathbf{e}_b - \mathbf{e}_a) = \|\mathbf{L}^{+/2}(\mathbf{e}_b - \mathbf{e}_a)\|^2.$$

**Example: Effective resistance in a path.** Consider a path graph on vertices  $V = \{1, 2, 3, \dots, k+1\}$ , with resistances  $\mathbf{r}(1), \mathbf{r}(2), \dots, \mathbf{r}(k)$  on the edges of the path.



Figure 5.2: A path graph with  $k$  edges.

The effective resistance between the endpoints is

$$R_{\text{eff}}(1, k+1) = \sum_{i=1}^k \mathbf{r}(i)$$

To see this, observe that to have 1 unit of flow going from vertex 1 to vertex  $k+1$ , we must have one unit flowing across each edge  $i$ . Let  $\Delta(i)$  be the voltage difference across edge  $i$ , and  $\mathbf{f}(i)$  the flow on the edge. Then  $1 = \mathbf{f}(i) = \frac{\Delta(i)}{r(i)}$ , so that  $\Delta(i) = \mathbf{r}(i)$ . The electrical voltages are then  $\tilde{\mathbf{x}} \in \mathbb{R}^V$  where  $\tilde{\mathbf{x}}(i) = \tilde{\mathbf{x}}(1) + \sum_{j < i} \Delta(j)$ . Hence the effective resistance is

$$R_{\text{eff}}(1, k+1) = \mathbf{d}^\top \tilde{\mathbf{x}} = (\mathbf{e}_{k+1} - \mathbf{e}_1)^\top \tilde{\mathbf{x}} = \tilde{\mathbf{x}}(k+1) - \tilde{\mathbf{x}}(1) = \sum_{i=1}^k \mathbf{r}(i).$$

This behavior is sometimes known as the fact that the resistance of resistors adds up when they are connected in series.

**Example: Effective resistance of parallel edges.** So far, we have only considered graphs with at most one edge between any two vertices. But that math also works if we allow a pair of vertices to have multiple distinct edges connecting them. We refer to this as *multi-edges*. Suppose we have a graph on just two vertices,  $V = \{1, 2\}$ , and these are connected by  $k$  parallel multi-edges with resistances  $\mathbf{r}(1), \mathbf{r}(2), \dots, \mathbf{r}(k)$ .

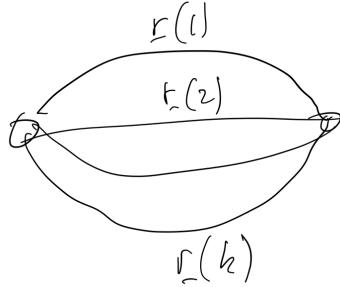


Figure 5.3: A graph on just two vertices with  $k$  parallel multiedges.

The effective resistance between the endpoints is

$$R_{\text{eff}}(1, 2) = \frac{1}{\sum_{i=1}^k 1/\mathbf{r}(i)}.$$

Let's see why. Our electrical voltages  $\tilde{\mathbf{x}} \in \mathbb{R}^V$  can be described by just the voltage difference  $\Delta \in \mathbb{R}$  between vertex 1 and vertex 2, i.e.  $\tilde{\mathbf{x}}(2) - \tilde{\mathbf{x}}(1) = \Delta$ , which creates a flow on edge  $i$  of  $\tilde{\mathbf{f}}(i) = \Delta/\mathbf{r}(i)$ . Thus the total flow from vertex 1 to vertex 2 is  $1 = \sum_i \Delta/\mathbf{r}(i)$ , so that  $\Delta = \frac{1}{\sum_{i=1}^k 1/\mathbf{r}(i)}$ . Meanwhile, the effective resistance is also

$$R_{\text{eff}}(1, 2) = (\mathbf{e}_2 - \mathbf{e}_1)^\top \tilde{\mathbf{x}} = \Delta = \frac{1}{\sum_{i=1}^k 1/\mathbf{r}(i)}$$

### 5.3.1 Effective Resistance is a Distance

**Definition 5.3.3.** Consider a weighted undirected graph  $G$  with vertex set  $V$ . We say function  $d : V \times V \rightarrow \mathbb{R}$ , which takes a pair of vertices and returns a real number, is a *distance* if it satisfies

1.  $d(a, a) = 0$  for all  $a \in V$
2.  $d(a, b) \geq 0$  for all  $a, b \in V$ .
3.  $d(a, b) = d(b, a)$  for all  $a, b \in V$ .
4.  $d(a, b) \leq d(a, c) + d(c, b)$  for all  $a, b, c \in V$ .

**Lemma 5.3.4.**  $R_{\text{eff}}$  is a distance.

Before proving this lemma, let's see a claim that will help us finish the proof.

**Claim 5.3.5.** Let  $\mathbf{L}\tilde{\mathbf{x}} = \mathbf{e}_b - \mathbf{e}_a$ . Then for all  $c \in V$ , we have  $\tilde{\mathbf{x}}(b) \geq \tilde{\mathbf{x}}(c) \geq \tilde{\mathbf{x}}(a)$ .

We only sketch a proof of this claim:

*Proof sketch.* Consider any  $c \in V$ , where  $c \neq a, b$ . Now  $(\mathbf{L}\tilde{\mathbf{x}})(c) = 0$ , i.e.

$$\left( \sum_{(u,c)} \mathbf{w}(u,c) \right) \tilde{\mathbf{x}}(c) - \left( \sum_{(u,c)} \mathbf{w}(u,c) \tilde{\mathbf{x}}(u) \right) = 0$$

Rearranging  $\tilde{\mathbf{x}}(c) = \frac{\sum_{(u,c)} \mathbf{w}(u,c) \tilde{\mathbf{x}}(u)}{\sum_{(u,c)} \mathbf{w}(u,c)}$ . This tells us that  $\tilde{\mathbf{x}}(c)$  is a weighted average of the voltages of its neighbors. From this, we can show that  $\tilde{\mathbf{x}}(a)$  and  $\tilde{\mathbf{x}}(b)$  are the extreme values.  $\square$

*Proof.* It is easy to check that conditions 1, 2, and 3 of Definition 5.3.3 are satisfied by  $R_{\text{eff}}$ . Let us confirm condition 4.

For any  $u, v$ , let  $\tilde{\mathbf{x}}_{u,v} = \mathbf{L}^+(-\mathbf{e}_u + \mathbf{e}_v)$ . Then

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

Thus,

$$\begin{aligned}
R_{\text{eff}}(a, b) &= (-\mathbf{e}_a + \mathbf{e}_b)^\top \tilde{\mathbf{x}}_{a,b} = (-\mathbf{e}_a + \mathbf{e}_b)^\top (\tilde{\mathbf{x}}_{a,c} + \tilde{\mathbf{x}}_{c,b}) \\
&= -\tilde{\mathbf{x}}_{a,c}(a) + \tilde{\mathbf{x}}_{a,c}(b) - \tilde{\mathbf{x}}_{c,b}(a) + \tilde{\mathbf{x}}_{c,b}(b) \\
&\leq -\tilde{\mathbf{x}}_{a,c}(a) + \tilde{\mathbf{x}}_{a,c}(c) - \tilde{\mathbf{x}}_{c,b}(c) + \tilde{\mathbf{x}}_{c,b}(b).
\end{aligned}$$

where in the last line we applied Claim 5.3.5 to show that  $\tilde{\mathbf{x}}_{a,c}(b) \leq \tilde{\mathbf{x}}_{a,c}(c)$  and  $-\tilde{\mathbf{x}}_{c,b}(a) \leq -\tilde{\mathbf{x}}_{c,b}(c)$ .  $\square$

# Chapter 6

## Different Perspectives on Gaussian Elimination

### 6.1 An Optimization View of Gaussian Elimination for Laplacians

In this section, we will explore how to exactly minimize a Laplacian quadratic form by minimizing over one variable at a time. It turns out that this is in fact Gaussian Elimination in disguise – or, more precisely, the variant of Gaussian elimination that we tend to use on symmetric matrices, which is called Cholesky factorization.

Consider a Laplacian  $\mathbf{L}$  of a connected graph  $G = (V, E, \mathbf{w})$ , where  $\mathbf{w} \in \mathbb{R}^E$  is a vector of positive edge weights. Let  $\mathbf{W} \in \mathbb{R}^{E \times E}$  be the diagonal matrix with the edge weights on the diagonal, i.e.  $\mathbf{W} = \text{diag}(\mathbf{w})$  and  $\mathbf{L} = \mathbf{B} \mathbf{W} \mathbf{B}^\top$ . Let  $\mathbf{d} \in \mathbb{R}^V$  be a demand vector s.t.  $\mathbf{d} \perp \mathbf{1}$ .

Let us define an energy

$$\mathcal{E}(\mathbf{x}) = -\mathbf{d}^\top \mathbf{x} + \frac{1}{2} \mathbf{x}^\top \mathbf{L} \mathbf{x}$$

Note that this function is convex and is minimized at  $\mathbf{x}$  s.t.  $\mathbf{L}\mathbf{x} = \mathbf{d}$ .

We will now explore an approach to solving the minimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^V} \mathcal{E}(\mathbf{x})$$

Let  $\mathbf{x} = \begin{pmatrix} y \\ \mathbf{z} \end{pmatrix}$  where  $y \in \mathbb{R}$  and  $\mathbf{z} \in \mathbb{R}^{V \setminus \{1\}}$ .

We will now explore how to minimize over  $y$ , given any  $\mathbf{z}$ . Once we find an expression for  $y$  in terms of  $\mathbf{z}$ , we will be able to reduce it to a new quadratic minimization problem in  $\mathbf{z}$ ,

$$\mathcal{E}'(\mathbf{z}) = -\mathbf{d}'^\top \mathbf{z} + \frac{1}{2} \mathbf{z}^\top \mathbf{L}' \mathbf{z}$$

where  $\mathbf{d}'$  is a demand vector on the remaining vertices, with  $\mathbf{d} \perp \mathbf{1}$  and  $\mathbf{L}'$  is a Laplacian of a graph on the remaining vertices  $V' = V \setminus \{1\}$ . We can then repeat the procedure to eliminate another variable and so on. Eventually, we can then find all the solution to our original minimization problem.

To help us understand how to minimize over the first variable, we introduce some notation for the first row and column of the Laplacian:

$$\mathbf{L} = \begin{pmatrix} W & -\mathbf{a}^\top \\ -\mathbf{a} & \text{diag}(\mathbf{a}) + \mathbf{L}_{-1} \end{pmatrix}. \quad (6.1)$$

Note that  $W$  is the weighted degree of vertex 1, and that

$$\begin{pmatrix} W & -\mathbf{a}^\top \\ -\mathbf{a} & \text{diag}(\mathbf{a}) \end{pmatrix} \quad (6.2)$$

is the Laplacian of the subgraph of  $G$  containing only the edges incident on vertex 1, while  $\mathbf{L}_{-1}$  is the Laplacian of the subgraph of  $G$  containing all edges *not* incident on vertex 1.

Let us also write  $\mathbf{d} = \begin{pmatrix} b \\ \mathbf{c} \end{pmatrix}$  where  $b \in \mathbb{R}$  and  $\mathbf{c} \in \mathbb{R}^{V \setminus \{1\}}$ .

Now,

$$\begin{aligned} \mathcal{E}(\mathbf{x}) &= -\mathbf{d}^\top \mathbf{x} + \frac{1}{2} \mathbf{x}^\top \mathbf{L} \mathbf{x} = -\begin{pmatrix} b \\ \mathbf{c} \end{pmatrix}^\top \begin{pmatrix} y \\ \mathbf{z} \end{pmatrix} + \frac{1}{2} \begin{pmatrix} y \\ \mathbf{z} \end{pmatrix}^\top \begin{pmatrix} W & -\mathbf{a}^\top \\ -\mathbf{a} & \text{diag}(\mathbf{a}) + \mathbf{L}_{-1} \end{pmatrix} \begin{pmatrix} y \\ \mathbf{z} \end{pmatrix} \\ &= -by - \mathbf{c}^\top \mathbf{z} + \frac{1}{2} (y^2 W - 2y \mathbf{a}^\top \mathbf{z} + \mathbf{z}^\top \text{diag}(\mathbf{a}) \mathbf{z} + \mathbf{z}^\top \mathbf{L}_{-1} \mathbf{z}). \end{aligned}$$

Now, to minimize over  $y$ , we set  $\frac{\partial}{\partial y} \mathcal{E}(\mathbf{x}) = 0$  and get

$$-b + yW - \mathbf{a}^\top \mathbf{z} = 0.$$

Solving for  $y$ , we get that the minimizing  $y$  is

$$y = \frac{1}{W}(b + \mathbf{a}^\top \mathbf{z}). \quad (6.3)$$

Observe that

$$\begin{aligned} \mathcal{E}(\mathbf{x}) &= -by - \mathbf{c}^\top \mathbf{z} + \frac{1}{2} (y^2 W - 2y \mathbf{a}^\top \mathbf{z} + \mathbf{z}^\top \text{diag}(\mathbf{a}) \mathbf{z} + \mathbf{z}^\top \mathbf{L}_{-1} \mathbf{z}) \\ &= -by - \mathbf{c}^\top \mathbf{z} + \frac{1}{2} \left( \frac{1}{W} (yW - \mathbf{a}^\top \mathbf{z})^2 - \underbrace{\frac{1}{W} \mathbf{z}^\top \mathbf{a} \mathbf{a}^\top \mathbf{z} + \mathbf{z}^\top \text{diag}(\mathbf{a}) \mathbf{z} + \mathbf{z}^\top \mathbf{L}_{-1} \mathbf{z}}_{\text{Let } S = \text{diag}(\mathbf{a}) - \frac{1}{W} \mathbf{a} \mathbf{a}^\top + \mathbf{L}_{-1}} \right) \\ &= -by - \mathbf{c}^\top \mathbf{z} + \frac{1}{2} \left( \frac{1}{W} (yW - \mathbf{a}^\top \mathbf{z})^2 + \mathbf{z}^\top S \mathbf{z} \right), \end{aligned}$$

where we simplified the expression by defining  $\mathbf{S} = \text{diag}(\mathbf{a}) - \frac{1}{W}\mathbf{a}\mathbf{a}^\top + \mathbf{L}_{-1}$ . Plugging in  $y = \frac{1}{W}(b + \mathbf{a}^\top \mathbf{z})$ , we get

$$\min_y \mathcal{E} \begin{pmatrix} y \\ \mathbf{z} \end{pmatrix} = -\left( \mathbf{c} + b \frac{1}{W} \mathbf{a} \right)^\top \mathbf{z} - \frac{b^2}{2W} + \frac{1}{2} \mathbf{z}^\top \mathbf{S} \mathbf{z}.$$

Now, we define  $\mathbf{d}' = \mathbf{c} + b \frac{1}{W} \mathbf{a}$  and  $\mathcal{E}'(\mathbf{z}) = -\mathbf{d}'^\top \mathbf{z} + \frac{1}{2} \mathbf{z}^\top \mathbf{S} \mathbf{z}$ . And, we can see that

$$\arg \min_z \min_y \mathcal{E} \begin{pmatrix} y \\ \mathbf{z} \end{pmatrix} = \arg \min_z \mathcal{E}'(\mathbf{z}),$$

since dropping the constant term  $-\frac{b^2}{2W}$  does not change what the minimizing  $\mathbf{z}$  values are.

### Claim 6.1.1.

1.  $\mathbf{d}' \perp \mathbf{1}$

2.  $\mathbf{S} = \text{diag}(\mathbf{a}) - \frac{1}{W}\mathbf{a}\mathbf{a}^\top + \mathbf{L}_{-1}$  is a Laplacian of a graph on the vertex set  $V \setminus \{1\}$ .

We will prove Claim 6.1.1 in a moment. From the Claim, we see that the problem of finding  $\arg \min_z \mathcal{E}'(\mathbf{z})$ , is exactly of the same form as finding  $\arg \min_x \mathcal{E}(\mathbf{x})$ , but with one fewer variables.

We can get a minimizing  $\mathbf{x}$  that solves  $\arg \min_x \mathcal{E}(\mathbf{x})$  by repeating the variable elimination procedure until we get down to a single variable and finding its value. We then have to work back up to getting a solution for  $\mathbf{z}$ , and then substitute that into Equation (6.3) to get the value for  $y$ .

**Remark 6.1.2.** In fact, this perspective on Gaussian elimination also makes sense for any positive definite matrix. In this setting, minimizing over one variable will leave us with another positive definite quadratic minimization problem.

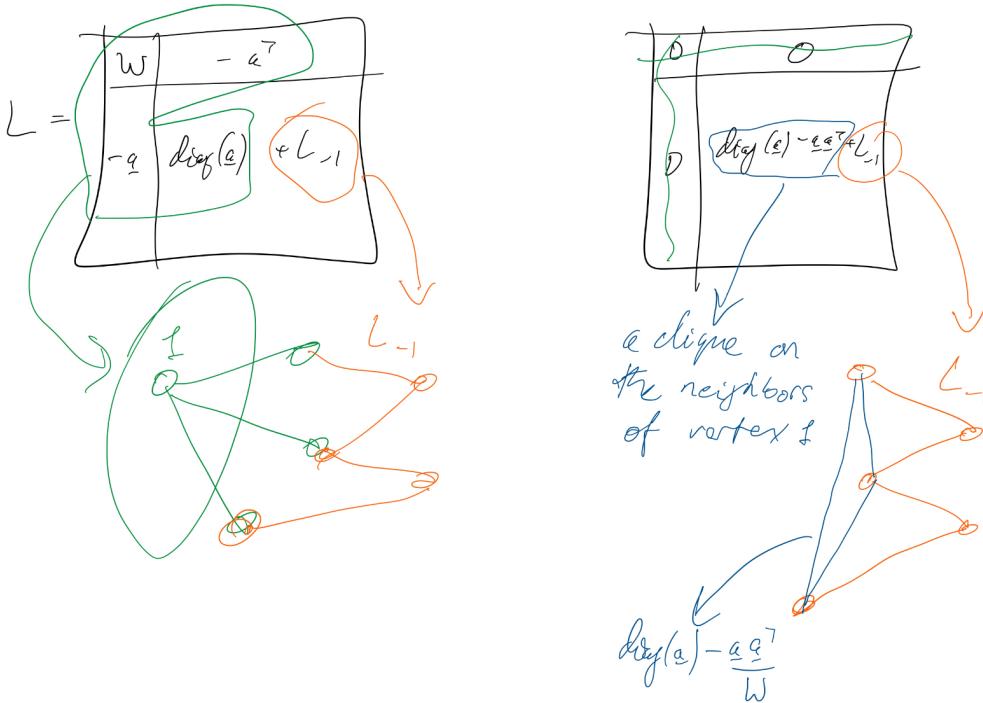
*Proof of Claim 6.1.1.* To establish the first part, we note that  $\mathbf{1}^\top \mathbf{d}' = \mathbf{1}^\top \mathbf{c} + b \frac{\mathbf{1}^\top \mathbf{a}}{W} = \mathbf{1}^\top \mathbf{c} + b = \mathbf{1}^\top \mathbf{d} = 0$ . To establish the second part, we notice that  $\mathbf{L}_{-1}$  is a graph Laplacian by definition. Since the sum of two graph Laplacians is another graph Laplacian, it now suffices to show that  $\mathbf{S}$  is a graph Laplacian.

**Claim 6.1.3.** A matrix  $\mathbf{M}$  is a graph Laplacian if and only if it satisfies the following conditions:

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

Let's see that Claim 6.1.3 is true. Firstly, when the conditions hold we can write  $\mathbf{M} = \mathbf{D} - \mathbf{A}$  where  $\mathbf{D}$  is diagonal and non-negative, and  $\mathbf{A}$  is non-negative, symmetric, and zero on the diagonal, and from the last condition  $\mathbf{D}(i, i) = \sum_{j \neq i} \mathbf{A}(i, j)$ . Thus we can view  $\mathbf{A}$  as a graph adjacency matrix and  $\mathbf{D}$  as the corresponding diagonal matrix of weighted degrees. Secondly, it is easy to check that the conditions hold for any graph Laplacian, so the conditions indeed hold if and only if. Now we have to check that the claim applies to  $\mathbf{S}$ . We leave this as an exercise for the reader.

Finally, we want to argue that the graph corresponding to  $\mathbf{S}$  is connected. Consider any  $i, j \in V \setminus \{1\}$ . Since  $G$ , the graph of  $\mathbf{L}$ , is connected, there exists a simple path in  $G$  connecting  $i$  and  $j$ . If this path does not use vertex 1, it is a path in the graph of  $\mathbf{L}_{-1}$  and hence in the graph of  $\mathbf{S}$ . If the path does use vertex 1, it must do so by reaching the vertex on some edge  $(v, 1)$  and leaving on a different edge  $(1, u)$ . Replace this pair of edges with edge  $(u, v)$ , which appears in the graph of  $\mathbf{S}$  because  $\mathbf{S}(u, v) < 0$ . Now we have a path in the graph of  $\mathbf{S}$ .  $\square$



## 6.2 An Additive View of Gaussian Elimination

**Cholesky decomposition basics.** Again we consider a graph Laplacian  $\mathbf{L} \in \mathbb{R}^{n \times n}$  of a connected graph  $G = (V, E, \mathbf{w})$ , where as usual  $|V| = n$  and  $|E| = m$ .

In this Section, we'll study how to decompose a graph Laplacian as  $\mathbf{L} = \mathcal{L}\mathcal{L}^\top$ , where  $\mathcal{L} \in \mathbb{R}^{n \times n}$  is a lower triangular matrix, i.e.  $\mathcal{L}(i, j) = 0$  for  $i < j$ . Such a factorization is called a Cholesky decomposition. It is essentially the result of Gaussian elimination with a slight twist to ensure the matrices maintained at intermediate steps of the algorithm remain symmetric.

We use  $\text{nnz}(\mathbf{A})$  to denote the number of non-zero entries of matrix  $\mathbf{A}$ .

**Lemma 6.2.1.** *Given an invertible square lower triangular matrix  $\mathcal{L}$ , we can solve the linear equation  $\mathcal{L}\mathbf{y} = \mathbf{b}$  in time  $O(\text{nnz}(\mathcal{L}))$ . Similarly, given an upper triangular matrix  $\mathcal{U}$ , we can solve linear equations  $\mathcal{U}\mathbf{z} = \mathbf{c}$  in time  $O(\text{nnz}(\mathcal{U}))$ .*

We omit the proof, which is a straight-forward exercise. The algorithms for solving linear equations in upper and lower triangular matrices are known as forward and back substitution respectively.

**Remark 6.2.2.** Strictly speaking, the lemma requires us to have access an adjacency list representation of  $\mathcal{L}$  so that we can quickly tell where the non-zero entries are.

Using forward and back substitution, if we have a decomposition of an invertible matrix  $\mathbf{M}$  into  $\mathbf{M} = \mathcal{L}\mathcal{L}^\top$ , we can now solve linear equations in  $\mathbf{M}$  in time  $O(\text{nnz}(\mathcal{L}))$ .

**Remark 6.2.3.** We have learned about decompositions using a lower triangular matrix, and later we will see an algorithm for computing these. In fact, we can have more flexibility than that. From an algorithmic perspective, it is sufficient that there exists a permutation matrix  $\mathbf{P}$  s.t.  $\mathbf{P}\mathcal{L}\mathbf{P}^\top$  is lower triangular. If we know the ordering under which the matrix becomes lower triangular, we can perform substitution according to that order to solve linear equations in the matrix without having to explicitly apply a permutation to the matrix.

**Dealing with pseudoinverses.** But how can we solve a linear equation in  $\mathbf{L} = \mathcal{L}\mathcal{L}^\top$ , where  $\mathbf{L}$  is not invertible? For graph Laplacians we have a simple characterization of the kernel, and because of this, dealing with the lack of invertibility turns out to be fairly easy.

We can use the following lemma which you will prove in an exercise next week.

**Lemma 6.2.4.** *Consider a real symmetric matrix  $\mathbf{M} = \mathbf{X}\mathbf{Y}\mathbf{X}^\top$ , where  $\mathbf{X}$  is real and invertible and  $\mathbf{Y}$  is real symmetric. Let  $\Pi_{\mathbf{M}}$  denote the orthogonal projection to the image of  $\mathbf{M}$ . Then  $\mathbf{M}^+ = \Pi_{\mathbf{M}}(\mathbf{X}^\top)^{-1}\mathbf{Y}^+\mathbf{X}^{-1}\Pi_{\mathbf{M}}$ .*

The factorizations  $\mathbf{L} = \mathcal{L}\mathcal{L}^\top$  that we produce will have the property that all diagonal entries of  $\mathcal{L}$  are strictly non-zero, except that  $\mathcal{L}(n, n) = 0$ . Let  $\widehat{\mathcal{L}}$  be the matrix whose entries agree with  $\mathcal{L}$ , except that  $\widehat{\mathcal{L}}(n, n) = 1$ . Let  $\mathcal{D}$  be the diagonal matrix with  $\mathcal{D}(i, i) = 1$  for  $i < n$  and  $\mathcal{D}(n, n) = 0$ . Then  $\mathcal{L}\mathcal{L}^\top = \widehat{\mathcal{L}}\mathcal{D}\widehat{\mathcal{L}}^\top$ , and  $\widehat{\mathcal{L}}$  is invertible, and  $\mathcal{D}^+ = \mathcal{D}$ . Finally,  $\Pi_{\mathbf{L}} = \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$ , because this matrix acts like identity on vectors orthogonal to  $\mathbf{1}$  and ensures  $\Pi_{\mathbf{L}}\mathbf{1} = \mathbf{0}$ , and this matrix can be applied to a vector in  $O(n)$  time. Thus  $\mathbf{L}^+ = \Pi_{\mathbf{L}}(\widehat{\mathcal{L}}^\top)^{-1}\mathcal{D}\widehat{\mathcal{L}}^{-1}\Pi_{\mathbf{L}}$ , and this matrix can be applied in time  $O(\text{nnz}(\mathcal{L}))$ .

**An additive view of Gaussian Elimination.** The following theorem describes Gaussian Elimination / Cholesky decomposition of a graph Laplacian.

**Theorem 6.2.5** (Cholesky Decomposition on graph Laplacians). *Let  $\mathbf{L} \in \mathbb{R}^{n \times n}$  be a graph Laplacian of a connected graph  $G = (V, E, \mathbf{w})$ , where  $|V| = n$ . Using Gaussian Elimination, we can compute in  $O(n^3)$  time a factorization  $\mathbf{L} = \mathbf{L}\mathbf{L}^\top$  where  $\mathbf{L}$  is lower triangular, and has positive diagonal entries except  $\mathbf{L}(n, n) = 0$ .*

*Proof.* Let  $\mathbf{L}^{(0)} = \mathbf{L}$ . We will use  $\mathbf{A}(:, i)$  to denote the  $i$ th column of a matrix  $\mathbf{A}$ . Now, for  $i = 1$  to  $i = n - 1$  we define

$$\mathbf{l}_i = \frac{1}{\sqrt{\mathbf{L}^{(i-1)}(i, i)}} \mathbf{L}^{(i-1)}(:, i) \text{ and } \mathbf{L}^{(i)} = \mathbf{L}^{(i-1)} - \mathbf{l}_i \mathbf{l}_i^\top$$

Finally, we let  $\mathbf{l}_n = \mathbf{0}_{n \times 1}$ . We will show later that

$$\mathbf{L}^{(n-1)} = \mathbf{0}_{n \times n}. \quad (6.4)$$

It follows that  $\mathbf{L} = \sum_i \mathbf{l}_i \mathbf{l}_i^\top$ , provided this procedure is well-defined, i.e.  $\mathbf{L}^{(i-1)}(i, i) \neq 0$  for all  $i < n$ . We will sketch a proof of this later, while also establishing several other properties of the procedure.

Given a matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $U \subseteq [n]$ , we will use  $\mathbf{A}(U, U)$  to denote the principal submatrix of  $\mathbf{A}$  obtained by restricting to the rows and columns with index in  $U$ , i.e. all entries  $\mathbf{A}(i, j)$  where  $i, j \in U$ .

**Claim 6.2.6.** *Fix some  $i < n$ . Let  $U = \{i + 1, \dots, n\}$ . Then  $\mathbf{L}^{(i)}(i, j) = 0$  if  $i \notin U$  or  $j \notin U$ . And  $\mathbf{L}^{(i)}(U, U)$  is a graph Laplacian of a connected graph on the vertex set  $U$ .*

From this claim, it follows that  $\mathbf{L}^{(i-1)}(i, i) \neq 0$  for  $i < n$ , since a connected graph Laplacian on a graph with  $|U| > 1$  vertices cannot have a zero on the diagonal, and it follows that  $\mathbf{L}^{(n-1)}(i, i) = 0$ , because the only graph we allow on one vertex is the empty graph. This shows Equation (6.4) holds.  $\square$

*Sketch of proof of Claim 6.2.6.* We will focus on the first elimination, as the remaining are similar. Adopting the same notation as in Equation (6.1), we write

$$\mathbf{L}^{(0)} = \mathbf{L} = \begin{pmatrix} W & -\mathbf{a}^\top \\ -\mathbf{a} & \text{diag}(\mathbf{a}) + \mathbf{L}_{-1} \end{pmatrix}$$

and, noting that

$$\mathbf{l}_1 \mathbf{l}_1^\top = \begin{pmatrix} W & -\mathbf{a}^\top \\ -\mathbf{a} & \frac{1}{W} \mathbf{a} \mathbf{a}^\top \end{pmatrix}$$

we see that

$$\mathbf{L}^{(1)} = \mathbf{L}^{(0)} - \mathbf{l}_1 \mathbf{l}_1^\top = \begin{pmatrix} 0 & \mathbf{0} \\ \mathbf{0} & \text{diag}(\mathbf{a}) - \frac{1}{W} \mathbf{a} \mathbf{a}^\top + \mathbf{L}_{-1} \end{pmatrix}.$$

Thus the first row and column of  $\mathbf{L}^{(1)}$  are zero claimed. It also follows by Claim 6.1.1 that  $\mathbf{L}^{(1)}(\{2, \dots, n\}, \{2, \dots, n\})$  is the Laplacian of a connected graph. This proves Claim 6.2.6 for the case  $i = 1$ . An induction following the same pattern can be used to prove the claim for all  $i < n$ .  $\square$

# Chapter 7

## Conductance, Expanders and Cheeger's Inequality

A common algorithmic problem that arises is the problem of partitioning the vertex set  $V$  of a graph  $G$  into clusters  $X_1, X_2, \dots, X_k$  such that

- for each  $i$ , the *induced* graph  $G[X_i] = (X_i, E \cap (X_i \times X_i))$  is "well-connected", and
- only an  $\epsilon$ -fraction of edges  $e$  are not contained in any induced graph  $G[X_i]$  (where  $\epsilon$  is a very small constant).

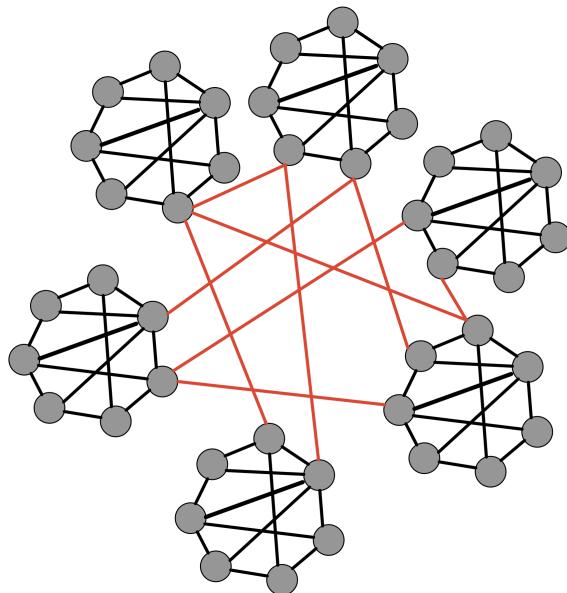


Figure 7.1: After removing the red edges (of which there are few in relation to the total number of edges), each connected component in  $G$  is "well-connected".

In this lecture, we make precise what "well-connected" means by introducing the notions of *conductance* and *expanders*.

Building on the chapter, we show that the second eigenvalue of the Laplacian  $L$  associated with graph  $G$  can be used to certify that a graph is "well-connected" (more precisely the second eigenvalue of a normalized version of the Laplacian). This result, called Cheeger's inequality, is one of the key tools in Spectral Graph Theory. Moreover, it can be turned into an algorithm that computes the partition efficiently!

## 7.1 Conductance and Expanders

**Graph Definitions.** In this lecture, we let  $G = (V, E)$  be unweighted<sup>1</sup> and always be *connected*, and let  $\mathbf{d}(v)$  be the degree of a vertex  $v$  in  $G$ . We define the *volume*  $\text{vol}(S)$  for any vertex subset  $S \subseteq V$ , to be the sum of degrees, i.e.  $\text{vol}(S) = \sum_{v \in S} \mathbf{d}(v)$ .

For any  $A, B \subseteq V$ , we define  $E(A, B)$  to be the set of edges in  $E \cap (A \times B)$ , i.e. with one endpoint in  $A$  and one endpoint in  $B$ . We let  $G[A]$  be the *induced* graph  $G$  by  $A \subseteq V$ , which is the graph  $G$  restricted to the vertices  $A$ , i.e. an edge  $e$  in  $G$  is in  $G[A]$  iff both endpoints are in  $A$ .

**Conductance.** Given set  $\emptyset \subset S \subset V$ , then we define the conductance  $\phi(S)$  of  $S$  by

$$\phi(S) = \frac{|E(S, V \setminus S)|}{\min\{\text{vol}(S), \text{vol}(V \setminus S)\}}.$$

It can be seen that  $\phi(\cdot)$  is symmetric in the sense that  $\phi(S) = \phi(V \setminus S)$ . We define the conductance of the graph  $G$  denoted  $\phi(G)$  by

$$\phi(G) = \min_{\emptyset \subset S \subset V} \phi(S).$$

We note that finding the conductance of a graph  $G$  is NP-hard. However, good approximations can be found as we will see today (and in a later lecture).

**Expander and Expander Decomposition.** For any  $\phi \in (0, 1]$ , we say that a graph  $G$  is a  $\phi$ -expander if  $\phi(G) \geq \phi$ . We say that the partition  $X_1, X_2, \dots, X_k$  of the vertex set  $V$  is a  $\phi$ -expander decomposition of quality  $q$  if

- each induced graph  $G[X_i]$  is a  $\phi$ -expander, and
- the number of edges not contained in any  $G[X_i]$  is at most  $q \cdot \phi \cdot m$ .

---

<sup>1</sup>Everything we present here also works for weighted graphs, however, we focus on unweighted graphs for simplicity.

Today, we obtain a  $\phi$ -expander decomposition of quality  $q = O(\phi^{-1/2} \cdot \log n)$ . In a few lectures, we revisit the problem and obtain quality  $q = O(\log^c n)$  for some small constant  $c$ . In practice, we mostly care about values  $\phi \approx 1$ .

**An Algorithm to Compute Conductance and Expander Decomposition.** In this lecture, the main focus is *not* to obtain an algorithm to compute conductance but rather only to show that the conductance of a graph can be approximated using the eigenvalues of the "normalized" Laplacian.

However, this proof gives then rise to an algorithm  $\text{CERTIFYORCUT}(G, \phi)$  that given a graph  $G$  and a parameter  $\phi$  either:

- Certifies that  $G$  is a  $\phi$ -expander, or
- Presents a *cut*  $S$  such that  $\phi(S) \leq \sqrt{2\phi}$ .

Careful but rather straight-forward recursive application of this procedure then yields an algorithm to compute a  $\phi$ -expander decomposition.

To implement this algorithm above, we are required to find the second eigenvalue of the 'normalized' Laplacian  $\mathbf{N}$  that we introduce in this lecture. We can find the eigenvalues using the power method<sup>2</sup> which multiplies a random vector with the matrix multiple times to find the largest eigenvalue. In our case, we want to use the inverse of  $\mathbf{N}$  and pick a random vector orthogonal to the trivial first eigenvector which yields the 2nd eigenvalue of  $\mathbf{N}$ . The inverse-matrix vector multiplications can be done extremely fast given a solver for Laplacian systems  $\mathbf{L}\mathbf{x} = \mathbf{d}$  which we learn about in the next chapters.

## 7.2 A Lower Bound for Conductance via Eigenvalues

**An Alternative Characterization of Conductance.** Let us now take a closer look at the definition of conductance and observe that if a set  $S$  has  $\text{vol}(S) \leq \text{vol}(V)/2$  then

$$\phi(S) = \frac{|E(S, V \setminus S)|}{\min\{\text{vol}(S), \text{vol}(V \setminus S)\}} = \frac{|E(S, V \setminus S)|}{\text{vol}(S)} = \frac{\mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S}{\mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S}.$$

To see this, observe that we can rewrite the numerator above using the Laplacian of  $G$  as

$$|E(S, V \setminus S)| = \sum_{(u,v) \in E} (\mathbf{1}_S(u) - \mathbf{1}_S(v))^2 = \mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S$$

where  $\mathbf{1}_S$  is the characteristic vector of  $S$ . Further, we can rewrite the denominator as

$$\text{vol}(S) = \mathbf{1}_S^\top \mathbf{d} = \mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S$$

---

<sup>2</sup>See [https://en.wikipedia.org/wiki/Power\\_iteration](https://en.wikipedia.org/wiki/Power_iteration).

where  $\mathbf{D} = \text{diag}(\mathbf{d})$  is the degree-matrix. We can now alternatively define the graph conductance of  $G$  by

$$\phi(G) = \min_{\substack{\emptyset \subset S \subset V, \\ \text{vol}(S) \leq \text{vol}(V)/2}} \frac{\mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S}{\mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S} \quad (7.1)$$

where we use that  $\phi(S) = \phi(V \setminus S)$  such that the objective value is unchanged as long as for each set  $\emptyset \subset S \subset V$  either  $S$  or  $V \setminus S$  is in the set that we minimize over.

**The Normalized Laplacian.** Let us next define the *normalized* Laplacian

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

To learn a bit about this new matrix, let us first look at the first eigenvalue where we use the test vector  $y = \mathbf{D}^{1/2} \mathbf{1}$ , to get by Courant-Fischer (see Theorem 3.1.4) that

$$\lambda_1(\mathbf{N}) = \min_{x \neq 0} \frac{x^\top \mathbf{N} x}{x^\top x} \leq \frac{y^\top \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2} y}{y^\top y} = \frac{\mathbf{1}^\top \mathbf{L} \mathbf{1}}{\mathbf{1}^\top \mathbf{y}} = 0 \quad (7.2)$$

because  $\mathbf{D}^{-1/2} \mathbf{D}^{1/2} = I$  and  $\mathbf{L} \mathbf{1} = 0$  (for the former we use the assumption that  $G$  is connected). Since  $\mathbf{N}$  is PSD (as you will show in the exercises), we also know  $\lambda_1(\mathbf{N}) \geq 0$ , so  $\lambda_1(\mathbf{N}) = 0$ .

Let us use Courant-Fischer again to reason a bit about the second eigenvalue of  $\mathbf{N}$ :

$$\lambda_2(\mathbf{N}) = \min_{\substack{x \perp \mathbf{D}^{1/2} \mathbf{1} \\ x \neq 0}} \frac{x^\top \mathbf{N} x}{x^\top x} = \min_{\substack{z \perp \mathbf{d} \\ z \neq 0}} \frac{z^\top \mathbf{D}^{1/2} \mathbf{N} \mathbf{D}^{1/2} z}{z^\top z} = \min_{\substack{z \perp \mathbf{d} \\ z \neq 0}} \frac{z^\top \mathbf{L} z}{z^\top z}. \quad (7.3)$$

**Relating Conductance to the Normalized Laplacian.** At this point, it might become clearer why  $\mathbf{N}$  is a natural matrix to consider when arguing about conductance: if we could argue that for every  $\emptyset \subset S \subset V$ ,  $\text{vol}(S) \leq \text{vol}(V)/2$ , we have  $\mathbf{1}_S \perp \mathbf{d}$ , then it would be easy to see that taking the second eigenvalue of  $\mathbf{N}$  in equation 7.3 is a relaxation of the minimization problem 7.1 defining  $\phi(G)$ .

While this is clearly not true, we can still argue along these lines.

**Theorem 7.2.1** (Cheeger's Inequality, Lower Bound). *We have  $\frac{\lambda_2(\mathbf{N})}{2} \leq \phi(G)$ .*

*Proof.* Instead of using  $\mathbf{1}_S$  directly, we shift  $\mathbf{1}_S$  by  $\mathbf{1}$  such that it is orthogonal to  $\mathbf{d}$ : we define  $\mathbf{z}_S = \mathbf{1}_S - \alpha \mathbf{1}$  where  $\alpha$  is the scalar that solves

$$\begin{aligned} 0 &= \mathbf{d}^\top \mathbf{z}_S \\ \iff 0 &= \mathbf{d}^\top (\mathbf{1}_S - \alpha \mathbf{1}) \\ \iff 0 &= \mathbf{d}^\top \mathbf{1}_S - \alpha \mathbf{d}^\top \mathbf{1} \\ \iff \alpha &= \frac{\mathbf{d}^\top \mathbf{1}_S}{\mathbf{d}^\top \mathbf{1}} = \frac{\text{vol}(S)}{\text{vol}(V)}. \end{aligned}$$

To conclude the proof, it remains to argue that  $\frac{\mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S}{\mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S} \geq \frac{1}{2} \cdot \frac{\mathbf{z}_S^\top \mathbf{L} \mathbf{z}_S}{\mathbf{z}_S^\top \mathbf{D} \mathbf{z}_S}$ .

- Numerator: since  $\mathbf{1}^\top \mathbf{L} \mathbf{1} = 0$ , we have that  $\mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S = \mathbf{z}_S^\top \mathbf{L} \mathbf{z}_S$ .
- Denominator: observe by straight-forward calculations that

$$\begin{aligned}\mathbf{z}_S^\top \mathbf{D} \mathbf{z}_S &= \text{vol}(S) \cdot (1 - \alpha)^2 + \text{vol}(V \setminus S) \cdot (-\alpha)^2 \\ &= \text{vol}(S) - 2 \text{vol}(S) \cdot \alpha + \text{vol}(V) \cdot \alpha^2 \\ &= \text{vol}(S) - \frac{\text{vol}(S)^2}{\text{vol}(V)} \\ &= \text{vol}(S) - \text{vol}(S) \cdot \frac{\text{vol}(S)}{\text{vol}(V)} \\ &\geq \frac{1}{2} \text{vol}(S) = \frac{1}{2} \mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S\end{aligned}$$

where we use the assumption that  $\text{vol}(S) \leq \text{vol}(V)/2$ .

□

### 7.3 An Upper Bound for Conductance via Eigenvalues

Slightly more surprisingly, we can also show that the second eigenvalue  $\lambda_2(\mathbf{N})$  can be used to upper bound the conductance.

**Theorem 7.3.1** (Cheeger's Inequality, Upper Bound). *We have  $\phi(G) \leq \sqrt{2 \cdot \lambda_2(\mathbf{N})}$ .*

*Proof.* To prove the theorem, we want to show that for *any*  $\mathbf{z} \perp \mathbf{d}$ , we can find a set  $\emptyset \subset S \subset V$  with  $\text{vol}(S) \leq \text{vol}(V)/2$ , such that

$$\frac{\mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S}{\mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S} \leq \sqrt{2 \cdot \frac{\mathbf{z}^\top \mathbf{L} \mathbf{z}}{\mathbf{z}^\top \mathbf{D} \mathbf{z}}}. \quad (7.4)$$

As a first step, we would like to change  $\mathbf{z}$  slightly to make it more convenient to work with:

- we *relabel* the vertices in  $V$  such that we have

$$\mathbf{z}(1) \leq \mathbf{z}(2) \leq \dots \leq \mathbf{z}(n).$$

- we *center*  $\mathbf{z}$ , that is we let  $\mathbf{z}_c = \mathbf{z} - \alpha \mathbf{1}$  where  $\alpha$  is chosen such that

$$\sum_{z_c(i) < 0} \mathbf{d}(i) < \text{vol}(V)/2 \text{ and } \sum_{z_c(i) \leq 0} \mathbf{d}(i) \geq \text{vol}(V)/2$$

i.e.  $\sum_{z_c(i) > 0} \mathbf{d}(i) \leq \text{vol}(V)/2$ .

- we *scale*, let  $\mathbf{z}_{sc} = \beta \mathbf{z}_c$  for some scalar  $\beta$  such that  $\mathbf{z}_{sc}(1)^2 + \mathbf{z}_{sc}(n)^2 = 1$ .

In the exercises, you will show that changing  $\mathbf{z}$  to  $\mathbf{z}_{sc}$  can only make the ratio we are interested in smaller, i.e. that  $\frac{\mathbf{z}^\top \mathbf{L} \mathbf{z}}{\mathbf{z}^\top \mathbf{D} \mathbf{z}} \geq \frac{\mathbf{z}_{sc}^\top \mathbf{L} \mathbf{z}_{sc}}{\mathbf{z}_{sc}^\top \mathbf{D} \mathbf{z}_{sc}}$ . Thus, if we can show that equation 7.4 holds for  $\mathbf{z}_{sc}$  in place of  $\mathbf{z}$ , then it also follows for  $\mathbf{z}$  itself.

We now arrive at the main idea of the proof: we define the set  $S_\tau = \{i \in V \mid \mathbf{z}_{sc}(i) < \tau\}$  for some random variable  $\tau$  with probability density function

$$p(t) = \begin{cases} 2 \cdot |t| & t \in [\mathbf{z}_{sc}(1), \mathbf{z}_{sc}(n)], \\ 0 & \text{otherwise.} \end{cases} \quad (7.5)$$

So, we have probability  $\mathbb{P}[\alpha < \tau < \beta] = \int_{t=\alpha}^{\beta} p(t) dt$ .

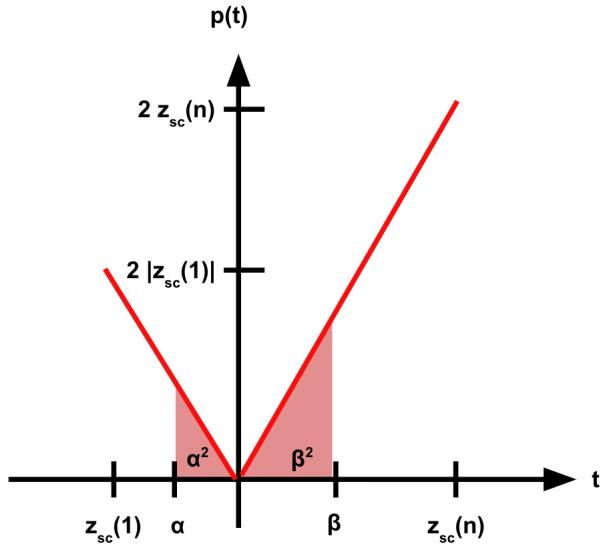


Figure 7.2: The probability density function  $p(t)$ . The area between  $\alpha$  and  $\beta$  for  $\alpha < 0 < \beta$  is  $\alpha^2 + \beta^2$ . This is the probability of the event  $\mathbb{P}[\alpha < \tau < \beta]$ . Note in particular that the total area is  $\mathbf{z}_{sc}(1)^2 + \mathbf{z}_{sc}(n)^2$  which is equal to 1 by assumption and thus  $p(t)$  is a valid pdf.

Since the volume incident to  $S_\tau$  might be quite large, let us define  $S$  for convenience by

$$S = \begin{cases} S_\tau & \text{vol}(S_\tau) < \text{vol}(V)/2, \\ V \setminus S_\tau & \text{otherwise.} \end{cases}$$

**Claim 7.3.2.** We have  $\frac{\mathbb{E}_\tau[\mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S]}{\mathbb{E}_\tau[\mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S]} \leq \sqrt{2 \cdot \frac{\mathbf{z}_{sc}^\top \mathbf{L} \mathbf{z}_{sc}}{\mathbf{z}_{sc}^\top \mathbf{D} \mathbf{z}_{sc}}}$ .

*Proof.* Recall  $\mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S = E(S_\tau, V \setminus S_\tau)$ , and by choice of  $\tau$ , we have for any edge  $e = \{i, j\} \in E$

where  $\mathbf{z}_{sc}(i) \leq \mathbf{z}_{sc}(j)$ ,

$$\begin{aligned}\mathbb{P}_\tau[e \in E(S_\tau, V \setminus S_\tau)] &= \mathbb{P}_\tau[\mathbf{z}_{sc}(i) < \tau \leq \mathbf{z}_{sc}(j)] \\ &= \int_{t=i}^j 2|t| dt \\ &= \begin{cases} \mathbf{z}_{sc}(i)^2 + \mathbf{z}_{sc}(j)^2 & \mathbf{z}(i) < 0 < \mathbf{z}(j), \\ |\mathbf{z}_{sc}(i)^2 - \mathbf{z}_{sc}(j)^2| & \text{otherwise.} \end{cases}\end{aligned}$$

The last equality can be seen from Figure 7.2. We can further upper bound either case by  $|\mathbf{z}_{sc}(i) - \mathbf{z}_{sc}(j)| \cdot (|\mathbf{z}_{sc}(i)| + |\mathbf{z}_{sc}(j)|)$  (we leave this as an exercise).

Using our new upper bound, we can sum over all edges  $e \in E$  to conclude that

$$\begin{aligned}\mathbb{E}_\tau[|E(S_\tau, V \setminus S_\tau)|] &\leq \sum_{i \sim j} |\mathbf{z}_{sc}(i) - \mathbf{z}_{sc}(j)| \cdot (|\mathbf{z}_{sc}(i)| + |\mathbf{z}_{sc}(j)|) \\ &\leq \sqrt{\sum_{i \sim j} (\mathbf{z}_{sc}(i) - \mathbf{z}_{sc}(j))^2 \cdot \sum_{i \sim j} (|\mathbf{z}_{sc}(i)| + |\mathbf{z}_{sc}(j)|)^2}\end{aligned}$$

where the last line follows from  $\langle \mathbf{x}, \mathbf{y} \rangle^2 \leq \langle \mathbf{x}, \mathbf{x} \rangle \cdot \langle \mathbf{y}, \mathbf{y} \rangle$  (i.e. Cauchy-Schwarz).

The first sum should look familiar by now: it is simply the Quadratic Laplacian Form  $\sum_{i \sim j} (\mathbf{z}_{sc}(i) - \mathbf{z}_{sc}(j))^2 = \mathbf{z}_{sc}^\top \mathbf{L} \mathbf{z}_{sc}$ .

It is not hard to reason about the second term either

$$\sum_{i \sim j} (|\mathbf{z}_{sc}(i)| + |\mathbf{z}_{sc}(j)|)^2 \leq 2 \sum_{i \sim j} \mathbf{z}_{sc}(i)^2 + \mathbf{z}_{sc}(j)^2 = 2 \sum_{i \in V} \mathbf{d}(i) \mathbf{z}_{sc}(i)^2 = 2 \mathbf{z}_{sc}^\top \mathbf{D} \mathbf{z}_{sc}.$$

Putting everything together, we obtain

$$\mathbb{E}_\tau[|E(S_\tau, V \setminus S_\tau)|] \leq \sqrt{\mathbf{z}_{sc}^\top \mathbf{L} \mathbf{z}_{sc} \cdot 2 \mathbf{z}_{sc}^\top \mathbf{D} \mathbf{z}_{sc}} = \sqrt{2 \cdot \frac{\mathbf{z}_{sc}^\top \mathbf{L} \mathbf{z}_{sc}}{\mathbf{z}_{sc}^\top \mathbf{D} \mathbf{z}_{sc}}} \cdot \mathbf{z}_{sc}^\top \mathbf{D} \mathbf{z}_{sc} \quad (7.6)$$

While this almost looks like what we want, we still have to argue that  $\mathbf{z}_{sc}^\top \mathbf{D} \mathbf{z}_{sc} = \mathbb{E}_\tau[\mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S]$  to finish the proof.

To this end, when unrolling the expectation, we use a simple trick that splits by cases:

$$\begin{aligned}\mathbb{E}_\tau[\mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S] &= \sum_{i \in V} \mathbf{d}(i) \cdot \mathbb{P}[i \in S] \\ &= \sum_{i \in V, \mathbf{z}_{sc}(i) < 0} \mathbf{d}(i) \cdot \mathbb{P}[i \in S \wedge S = S_\tau] + \sum_{i \in V, \mathbf{z}_{sc}(i) \geq 0} \mathbf{d}(i) \cdot \mathbb{P}[i \in S \wedge S \neq S_\tau] \\ &= \sum_{i \in V, \mathbf{z}_{sc}(i) < 0} \mathbf{d}(i) \cdot \mathbb{P}[\mathbf{z}_{sc}(i) < \tau \wedge \tau < 0] + \sum_{i \in V, \mathbf{z}_{sc}(i) \geq 0} \mathbf{d}(i) \cdot \mathbb{P}[\mathbf{z}_{sc}(i) \geq \tau \wedge \tau \geq 0]\end{aligned}$$

where we use the centering of  $\mathbf{z}_{sc}$  the definition of  $S$  and that the event  $\{i \in S \wedge S = S_\tau\}$  can be rewritten as the event  $\{i < \tau \wedge \tau < 0\}$  (the other case is analogous).

Let  $i$  be a vertex with  $\mathbf{z}_{sc}(i) < 0$ , then the probability  $\mathbb{P}[i \in S \wedge S = S_\tau]$  is exactly  $\mathbf{z}_{sc}(i)^2$  by choice of the density function of  $\tau$  (again the case for  $i$  with  $\mathbf{z}_{sc}(i)$  non-negative is analogous). Thus, summing over all vertices, we obtain

$$\begin{aligned}\mathbb{E}_\tau[\mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S] &= \sum_{i \in V, \mathbf{z}_{sc}(i) < 0} \mathbf{d}(i) \cdot \mathbb{P}[\mathbf{z}_{sc}(i) < \tau \wedge \tau < 0] + \sum_{i \in V, \mathbf{z}_{sc}(i) \geq 0} \mathbb{P}[\mathbf{z}_{sc}(i) \geq \tau \wedge \tau \geq 0] \\ &= \sum_{i \in V} \mathbf{d}(i) \cdot \mathbf{z}_{sc}(i)^2 = \mathbf{z}_{sc}^\top \mathbf{D} \mathbf{z}_{sc}.\end{aligned}$$

Therefore, we can plug in our result directly into Equation 7.6 and the proof is completed by dividing both sides by  $\mathbb{E}_\tau[\mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S]$ .  $\square$

While Theorem 7.3.2 only ensures our claim in expectation, this is already sufficient to conclude that there exists some set  $S$  that satisfies the same guarantees deterministically, as you will prove in Problem Set 4. This is often called the *probabilistic method of expectation* and can be seen from the definition of expectation. We have thus proven the upper bound of Cheeger's inequality.  $\square$

## 7.4 Conclusion

Today, we have introduced the concepts of conductance and formalized expanders and expander decompositions. These are crucial concepts that you will encounter often in literature and also again in this course. They are a key tool in many recent breakthroughs in Theoretical Computer Science.

In the second part of the lecture (the main part), we discussed Cheeger's inequality which allows to relate the second eigenvalue of the normalized Laplacian to a graphs conductance. We summarize the full statement here.

**Theorem 7.4.1** (Cheeger's Inequality). *We have  $\frac{\lambda_2(\mathcal{N})}{2} \leq \phi(G) \leq \sqrt{2 \cdot \lambda_2(\mathcal{N})}$ .*

We point out that this Theorem is tight as you will show in the exercises. The proof for Cheeger's inequality is probably the most advanced proof, we have seen so far in the course. The many tricks that make the proof work might sometimes seem a bit magical but it is important to remember that they are a result of many people polishing this proof over and over. The proof techniques used are extremely useful and can be re-used in various contexts. We therefore strongly encourage you to really understand the proof yourself!

# Chapter 8

## Random Walks

Today, we talk about random walks on graphs and how the spectrum of the Laplacian guides convergence of random walks. We start by giving the definition of a random walk on a weighted graph  $G = (V, E, w)$ .

### 8.1 A Primer on Random Walks

**Random Walk Basics.** We call a random sequence of vertices  $v_0, v_1, \dots$  a *random walk* on  $G$ , if  $v_0$  is a vertex in  $G$  chosen according to some probability distribution  $\mathbf{p}_0 \in \mathbb{R}^V$ ; and for any  $t \geq 0$ , we have

$$\mathbb{P}[v_{t+1} = v | v_t = u] = \begin{cases} \mathbf{w}(u, v) / \mathbf{d}(u) & \text{if } \{u, v\} \in E, \\ 0 & \text{otherwise.} \end{cases}$$

To gain some intuition for the definition, assume first that the graph  $G$  is undirected. Consider a **particle** that is placed at a random vertex  $v_0$  initially. Then at each step the particle is moved to a neighbor of the current vertex it is resting at, where the neighbor is chosen uniformly at random.

If the graph is weighted, then instead of choosing a neighbor  $v_{t+1}$  of a vertex  $v_t$  at each step uniformly at random, one chooses a neighbor  $v$  of  $v_t$  with probability  $\mathbf{w}(v, v_t)$  divided by the degree  $\mathbf{d}(v_t)$ .

**The Random Walk Matrix.** We now define the random walk matrix  $\mathbf{W}$  by

$$\mathbf{W} = \mathbf{A}\mathbf{D}^{-1}$$

and observe that for all vertices  $u, v \in V$  (and any  $t$ ), we have that

$$\mathbf{W}_{vu} = \begin{cases} \mathbf{w}(u, v) / \mathbf{d}(u) & \text{if } \{u, v\} \in E, \\ 0 & \text{otherwise.} \end{cases}$$

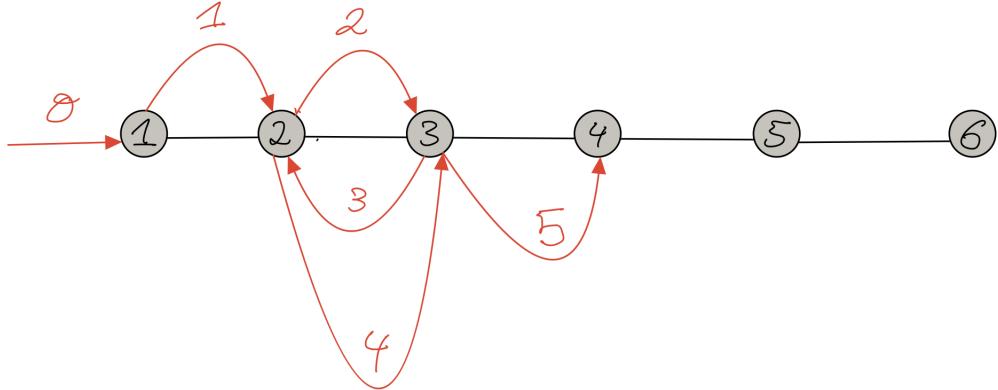


Figure 8.1: A (possibly random) walk where the red edges indicate the edges that the particle moves along. Here the walk visits the vertices  $v_0 = 1, v_1 = 2, v_2 = 3, v_3 = 2, v_4 = 3, v_5 = 4$ .

Thus,  $\mathbf{W}_{vu} = \mathbb{P}[v_{t+1} = v | v_t = u]$  (for any  $t$ ).

Therefore,  $\mathbf{W}\mathbf{1}_u$  is the distribution over the vertices that the random walk visits them at the next time step, given that it currently is at  $u$ . More generally, we can now compute the distribution  $\mathbf{p}_1$  over the vertices that they are visited at time 1 by  $\mathbf{W}\mathbf{p}_0$ , the distribution  $\mathbf{p}_2$  by  $\mathbf{W}\mathbf{p}_1 = \mathbf{W}(\mathbf{W}\mathbf{p}_0)$  and so on. Another way of writing this is  $\mathbf{p}_t = \mathbf{W}^t\mathbf{p}_0$ .

## 8.2 Convergence Results for Random Walks

In this first part of the chapter, we are interested mostly in convergence of random walks that is the two questions:

- How does a random walk behave after a large number of steps are taken?
- How many steps does it take asymptotically until the random walk behaves as if an infinite number of steps were taken?

To start shedding some light on these questions, we introduce stationary distributions.

**Stationary Distribution.** We call a distribution  $\boldsymbol{\pi} \in \mathbb{R}^V$ , a *stationary distribution* if  $\mathbf{W}\boldsymbol{\pi} = \boldsymbol{\pi}$ . That is  $\boldsymbol{\pi}$  is an eigenvector of  $\mathbf{W}$  associated with eigenvalue 1. It turns out such a stationary distribution always exists.

**Lemma 8.2.1.** *Every graph  $G$  has a stationary distribution.*

*Proof.* Let  $\boldsymbol{\pi} = \frac{\mathbf{d}}{\mathbf{1}^\top \mathbf{d}}$ . Clearly, we have that  $\|\boldsymbol{\pi}\|_1 = \sum_{v \in V} \mathbf{d}(v)/\mathbf{1}^\top \mathbf{d} = \frac{1}{\mathbf{1}^\top \mathbf{d}} \sum_{v \in V} \mathbf{d}(v) = 1$ ,

so  $\pi$  is indeed a distribution. Further note that

$$W\pi = AD^{-1} \cdot \frac{d}{1^\top d} = \frac{A1}{1^\top d} = \frac{d}{1^\top d} = \pi.$$

□

For many graphs one can show that for  $t \rightarrow \infty$ , we have that  $p_t \rightarrow \pi$ , i.e. that independent of the starting distribution  $p_0$ , the random walk always converges to distribution  $\pi$ .

Unfortunately, this is not true for all graphs: take the graph of two vertices connected by a single edge with  $p_0$  being 1 at one vertex and 0 at the other.

### 8.2.1 Making Random Walks Lazy

**Lazy Random Walks.** Luckily, we can overcome this issue by using a *lazy random walk*. A lazy random walk behaves just like a random walk, however, at each time step, with probability  $\frac{1}{2}$  instead of transitioning to a neighbor, it simply stays put. We give the lazy random walk matrix by

$$\tilde{W} = \frac{1}{2}I + \frac{1}{2}W = \frac{1}{2}(I + AD^{-1}).$$

It is not hard to see that the stationary distribution  $\pi$  for  $W$ , is also a stationary distribution for  $\tilde{W}$ .

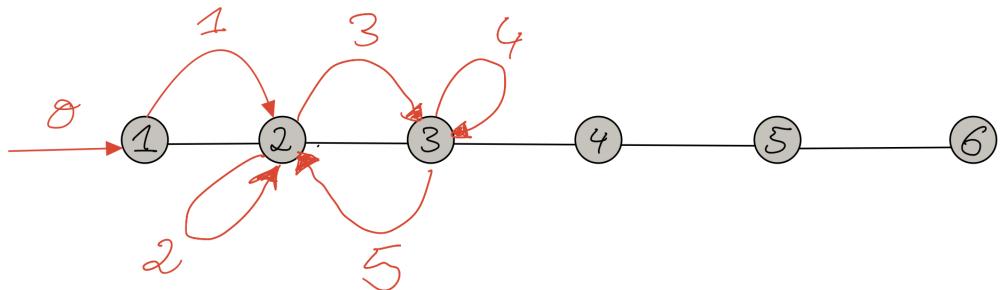


Figure 8.2: A lazy random walk where the red edges indicate the edges that the particle moves along. Here the lazy walk visits the vertices  $v_0 = 1, v_1 = 2, v_2 = 2, v_3 = 3, v_4 = 3, v_5 = 2$ .

**Lazy Random Walks and the Normalized Laplacian.** Recall that we defined  $\mathbf{N} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2} = \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} \iff \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} = \mathbf{I} - \mathbf{N}$ . We can therefore derive

$$\begin{aligned}\tilde{\mathbf{W}} &= \frac{1}{2} \mathbf{I} + \frac{1}{2} \mathbf{A} \mathbf{D}^{-1} \\ &= \frac{1}{2} \mathbf{I} + \frac{1}{2} \mathbf{D}^{1/2} \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} \mathbf{D}^{-1/2} \\ &= \frac{1}{2} \mathbf{I} + \frac{1}{2} \mathbf{D}^{1/2} (\mathbf{I} - \mathbf{N}) \mathbf{D}^{-1/2} \\ &= \frac{1}{2} \mathbf{I} + \frac{1}{2} \mathbf{D}^{1/2} \mathbf{I} \mathbf{D}^{-1/2} - \frac{1}{2} \mathbf{D}^{1/2} \mathbf{N} \mathbf{D}^{-1/2} \\ &= \mathbf{I} - \frac{1}{2} \mathbf{D}^{1/2} \mathbf{N} \mathbf{D}^{-1/2}\end{aligned}$$

We will now start to reason about the eigenvalues and eigenvectors of  $\tilde{\mathbf{W}}$  in terms of the normalized laplacian  $\mathbf{N}$  that we are already familiar with.

For the rest of the lecture, we let  $\nu_1 \leq \nu_2 \leq \dots \leq \nu_n$  be the eigenvalues of  $\mathbf{N}$  associated with the orthogonal eigenvectors  $\psi_1, \psi_2, \dots, \psi_n$  where we know that such eigenvectors exist by the Spectral Theorem. We note in particular that from the last lecture, we have that  $\psi_1 = \frac{\mathbf{d}^{1/2}}{(\mathbf{1}^\top \mathbf{d})^{1/2}}$  (see Equation 7.2 where we added a normalization such that  $\psi_1^\top \psi_1 = 1$ ).

**Lemma 8.2.2.** *For the  $i^{\text{th}}$  eigenvalue  $\nu_i$  of  $\mathbf{N}$  associated with eigenvector  $\psi_i$ , we have that  $\tilde{\mathbf{W}}$  has an eigenvalue of  $(1 - \frac{1}{2}\nu_i)$  associated with eigenvector  $\mathbf{D}^{1/2}\psi_i$ .*

*Proof.* The proof is by straight-forward calculations

$$\begin{aligned}\tilde{\mathbf{W}} \mathbf{D}^{1/2} \psi_i &= (\mathbf{I} - \frac{1}{2} \mathbf{D}^{1/2} \mathbf{N} \mathbf{D}^{-1/2}) \mathbf{D}^{1/2} \psi_i \\ &= \mathbf{D}^{1/2} \psi_i - \frac{1}{2} \mathbf{D}^{1/2} \mathbf{N} \psi_i \\ &= \mathbf{D}^{1/2} \psi_i - \frac{1}{2} \mathbf{D}^{1/2} \psi_i \nu_i = \mathbf{D}^{1/2} \psi_i (1 - \frac{1}{2} \nu_i).\end{aligned}$$

□

**Corollary 8.2.3.** *Every eigenvalue of  $\tilde{\mathbf{W}}$  is in  $[0, 1]$ .*

*Proof.* Recall that  $\mathbf{L} \preccurlyeq 2\mathbf{D}$  which implies that  $\mathbf{N} \preccurlyeq 2\mathbf{I}$ . But this implies that every eigenvalue of  $\mathbf{N}$  is in  $[0, 2]$ . Thus, using Lemma 8.2.2, the corollary follows. □

### 8.2.2 Convergence of Lazy Random Walks

We have now done enough work to obtain an interesting result. We can derive an alternative characterization of  $\mathbf{p}_t$  by expanding  $\mathbf{p}_0$  along an orthogonal eigenvectors basis and then we can repeatedly apply  $\tilde{\mathbf{W}}$  by taking powers of the eigenvalues.

Unfortunately,  $\tilde{\mathbf{W}}$  is not symmetric so its eigenvectors are not necessarily orthogonal. Instead, we use a simple trick that allows to expand along the eigenvectors of  $\mathbf{N}$

$$\forall i, \boldsymbol{\psi}_i^\top \mathbf{D}^{-1/2} \mathbf{p}_0 = \alpha_i \iff \mathbf{D}^{-1/2} \mathbf{p}_0 = \sum_{i=1}^n \alpha_i \boldsymbol{\psi}_i \iff \mathbf{p}_0 = \sum_{i=1}^n \alpha_i \mathbf{D}^{1/2} \boldsymbol{\psi}_i. \quad (8.1)$$

The above equivalences are best understood if you start from the middle. To get to the left side, you need to observe that multiplying both sides by  $\boldsymbol{\psi}_i^\top$  cancels all terms  $\boldsymbol{\psi}_j$  with  $j \neq i$  in the sum by orthogonality. To get the right hand side expression, one can simply left-multiply by  $\mathbf{D}^{1/2}$ . Technically, we have to show that  $\mathbf{D}^{-1/2} \mathbf{p}_0$  lives in the eigenspace of  $\mathbf{N}$  but we leave this as an exercise.

This allows us to express a right multiplication by  $\tilde{\mathbf{W}}$  as

$$\mathbf{p}_1 = \tilde{\mathbf{W}} \mathbf{p}_0 = \sum_{i=1}^n \alpha_i \tilde{\mathbf{W}} \mathbf{D}^{1/2} \boldsymbol{\psi}_i = \sum_{i=1}^n \alpha_i \left(1 - \frac{\nu_i}{2}\right) \mathbf{D}^{1/2} \boldsymbol{\psi}_i.$$

And as promised, if we apply  $\tilde{\mathbf{W}}$ , the lazy random walk operator,  $t$  times, we now obtain

$$\mathbf{p}_t = \sum_{i=1}^n \alpha_i \left(1 - \frac{\nu_i}{2}\right)^t \mathbf{D}^{1/2} \boldsymbol{\psi}_i = \alpha_1 \mathbf{D}^{1/2} \boldsymbol{\psi}_1 + \sum_{i=2}^n \alpha_i \left(1 - \frac{\nu_i}{2}\right)^t \mathbf{D}^{1/2} \boldsymbol{\psi}_i. \quad (8.2)$$

where we use in the last equality that  $\nu_1 = 0$ . Using this simple characterization, we immediately get that  $\mathbf{p}_t \rightarrow \boldsymbol{\pi}$  if  $\nu_i > 0$  for all  $i \geq 2$  (which is exactly when the graph is connected as you will prove in an exercise). To see this, observe that as  $t$  grows sum vanishes. We have that

$$\lim_{t \rightarrow \infty} \mathbf{p}_t = \alpha_1 \mathbf{D}^{1/2} \boldsymbol{\psi}_1 = \boldsymbol{\pi}.$$

where we used in the equality that  $\mathbf{D}^{1/2} \boldsymbol{\psi}_1 = \frac{\mathbf{d}}{(\mathbf{1}^\top \mathbf{d})^{1/2}}$  and the value of  $\alpha_1$  (from 8.1).

**Theorem 8.2.4.** *For any connected graph  $G$ , we have that the lazy random walk converges to the stationary distribution of  $G$ .*

### 8.2.3 The Rate of Convergence

Let us now come to the main result that we want to prove in this lecture.

**Theorem 8.2.5.** *In any unweighted (a.k.a. unit weight) connected graph  $G$ , for any  $\mathbf{p}_0$ , at any time step  $t$ , we have for  $\mathbf{p}_t = \tilde{\mathbf{W}}^t \mathbf{p}_0$  that*

$$\|\mathbf{p}_t - \boldsymbol{\pi}\|_\infty \leq e^{-\nu_2 \cdot t/2} \sqrt{n}$$

Instead of proving the theorem above, we prove the lemma below which gives point-wise convergence. This makes it more convenient to derive a proof and it is not hard to deduce the theorem above as a corollary.

**Lemma 8.2.6.** In any weighted connected graph  $G$ , for all  $a, b \in V$ , and any time step  $t$ , we have for  $\mathbf{p}_0 = \mathbf{1}_a$  and  $\mathbf{p}_t = \tilde{\mathbf{W}}^t \mathbf{p}_0$  that

$$|\mathbf{p}_t(b) - \boldsymbol{\pi}(b)| \leq e^{-\nu_2 \cdot t/2} \sqrt{\mathbf{d}_b / \mathbf{d}_a}$$

From Equation 8.2, we obtain that

$$\mathbf{p}_t(b) - \boldsymbol{\pi}(b) = \mathbf{1}_b^\top (\mathbf{p}_t - \boldsymbol{\pi}) = \mathbf{1}_b^\top \left( \sum_{i=2}^n \alpha_i \left(1 - \frac{\nu_i}{2}\right)^t \mathbf{D}^{1/2} \boldsymbol{\psi}_i \right) \quad (8.3)$$

$$= \sum_{i=2}^n \alpha_i \left(1 - \frac{\nu_i}{2}\right)^t \mathbf{1}_b^\top \mathbf{D}^{1/2} \boldsymbol{\psi}_i \leq \left(1 - \frac{\nu_2}{2}\right)^t \cdot \sum_{i=2}^n \alpha_i \mathbf{1}_b^\top \mathbf{D}^{1/2} \boldsymbol{\psi}_i \quad (8.4)$$

Taking the absolute value on both sides, we obtain that

$$|\mathbf{p}_t(b) - \boldsymbol{\pi}(b)| \leq \left(1 - \frac{\nu_2}{2}\right)^t \sum_{i=2}^n \left| \alpha_i \mathbf{1}_b^\top \mathbf{D}^{1/2} \boldsymbol{\psi}_i \right| \leq \left(1 - \frac{\nu_2}{2}\right)^t \sqrt{\left( \sum_{i=2}^n \alpha_i^2 \right) \left( \sum_{i=2}^n \left( \mathbf{1}_b^\top \mathbf{D}^{1/2} \boldsymbol{\psi}_i \right)^2 \right)}$$

where we use Cauchy-Schwarz in the last inequality, i.e.  $|\langle \mathbf{u}, \mathbf{v} \rangle|^2 \leq \langle \mathbf{u}, \mathbf{u} \rangle \cdot \langle \mathbf{v}, \mathbf{v} \rangle$ . Let us finally bound the two sums:

- By 8.1,  $\sum_{i=2}^n \alpha_i^2 = \sum_{i=2}^n \left( \boldsymbol{\psi}_i^\top \mathbf{D}^{-1/2} \mathbf{p}_0 \right)^2 \leq \|\mathbf{D}^{-1/2} \mathbf{p}_0\|_2^2 = \|\mathbf{D}^{-1/2} \mathbf{1}_a\|_2^2 = 1/\mathbf{d}_a$ .
- Finally, we show that  $\sum_{i=2}^n \left( \mathbf{1}_b^\top \mathbf{D}^{1/2} \boldsymbol{\psi}_i \right)^2 \leq \sum_{i=1}^n \left( \mathbf{1}_b^\top \mathbf{D}^{1/2} \boldsymbol{\psi}_i \right)^2 = \|\mathbf{D}^{1/2} \mathbf{1}_b\|_2^2 = \mathbf{d}_b$  (we only show the first equality, the other inequalities are straight-forward). We first expand the vector  $\mathbf{D}^{1/2} \mathbf{1}_b$  along the eigenvectors using some values  $\beta_i$  defined

$$\mathbf{D}^{1/2} \mathbf{1}_b = \sum_{i=1}^n \beta_i \boldsymbol{\psi}_i \iff \boldsymbol{\psi}_i^\top \mathbf{D}^{1/2} \mathbf{1}_b = \beta_i \iff \mathbf{1}_b^\top \mathbf{D}^{1/2} \boldsymbol{\psi}_i = \beta_i$$

We used orthogonality to get the first equivalence, and then just take the transpose to get the second. We can now write

$$\|\mathbf{D}^{1/2} \mathbf{1}_b\|_2^2 = (\mathbf{D}^{1/2} \mathbf{1}_b)^\top (\mathbf{D}^{1/2} \mathbf{1}_b) = \left( \sum_{i=1}^n \beta_i \boldsymbol{\psi}_i^\top \right) \left( \sum_{i=1}^n \beta_i \boldsymbol{\psi}_i \right) = \sum_{i=2}^n \beta_i^2$$

where we again used orthogonality of  $\boldsymbol{\psi}_i$ . The equality then follows by definition of  $\beta_i$ .

Putting everything together (and using  $1 + x \leq e^x, \forall x \in \mathbb{R}$ ), we obtain

$$|\mathbf{p}_t(b) - \boldsymbol{\pi}(b)| \leq \left(1 - \frac{\nu_2}{2}\right)^t \sqrt{\mathbf{d}_b / \mathbf{d}_a} \leq e^{-\nu_2 \cdot t/2} \sqrt{\mathbf{d}_b / \mathbf{d}_a}$$

## 8.3 Properties of Random Walks

We now shift our focus away from convergence of random walks and consider some interesting properties of random walks<sup>1</sup>. Here, we are no longer interested in lazy random walks, although all proofs can be straight-forwardly adapted. While in the previous section, we relied on computing the second eigenvalue of the Normalized Laplacian efficiently, here, we will discover that solving Laplacian systems, that is finding an  $\mathbf{x}$  such that  $\mathbf{Lx} = \mathbf{b}$  can solve a host of problems in random walks.

### 8.3.1 Hitting Times

One of the most natural questions one can ask about a random walk starting in a vertex  $a$  (i.e.  $\mathbf{p}_0 = \mathbf{1}_a$ ) is how many steps it takes to get to a special vertex  $s$ . This quantity is called the *hitting time* from  $a$  to  $s$  and we denote it by  $H_{a,s} = \inf\{t \mid \mathbf{v}_t = s\}$ . For the rest of this section, we are concerned with computing the expected hitting time, i.e.  $\mathbb{E}[H_{a,s}]$ .

It turns out, that it is more convenient to compute *all* expected hitting times  $H_{a,s}$  for vertices  $a \in V$  to a fixed  $s$ . We denote by  $\mathbf{h} \in \mathbf{R}^V$ , the vector with  $\mathbf{h}(a) = \mathbb{E}[H_{a,s}]$ . We now show that we can compute  $\mathbf{h}$  by solving a Laplacian system  $\mathbf{Lh} = \mathbf{b}$ . We will see later in the course that such systems (spoiler alert!) can be solved in time  $\tilde{O}(m)$ , so this will imply a near-linear time algorithm to compute the hitting times.

**Hitting Time and the Random Walk Matrix.** Let us first observe that if  $s = a$ , then the answer becomes trivially 0, i.e.  $\mathbf{h}(s) = 0$ .

We compute the rest of the vector by writing down a system of equations that recursively characterizes  $\mathbf{h}$ . Observe therefore first that for any  $a \neq s$ , we have that the random walks starting at  $a$  will next visit a neighbor  $b$  of  $a$ . If the selected neighbor  $b = s$ , the random walks stops; otherwise, the random walks needs in expectation  $\mathbb{E}[H_{b,s}]$  time to move to  $s$ .

We can express this algebraically by

$$\mathbf{h}_a = 1 + \sum_{a \sim b} \mathbb{P}[v_{t+1} = b \mid v_t = a] \cdot \mathbf{h}(b) = 1 + \sum_{a \sim b} \frac{\mathbf{w}(a,b)}{\mathbf{d}(a)} \cdot \mathbf{h}(b) = 1 + (\mathbf{W}\mathbf{1}_a)^\top \mathbf{h} = 1 + \mathbf{1}_a^\top \mathbf{W}^\top \mathbf{h}.$$

Using that  $\mathbf{h}(a) = \mathbf{1}_a^\top \mathbf{h} = \mathbf{1}_a^\top \mathbf{Ih}$ , we can rewrite this as

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

This gives a system of (linear) equations, that can be neatly summarized by

$$\mathbf{1} - \alpha \cdot \mathbf{1}_s = (\mathbf{I} - \mathbf{W}^\top) \mathbf{h}$$

---

<sup>1</sup>Note that this part of the script is in large part inspired by Aaron Sidford's lecture notes on the same subject, in his really interesting course Discrete Mathematics and Algorithms.

where we have an extra degree of freedom in choosing  $\alpha$  in formulating a constraint  $1 - \alpha = \mathbf{1}_s^\top (\mathbf{I} - \mathbf{W}^\top) \mathbf{h}$ . This extra degree of freedom stems from the fact that  $n - 1$  equations suffice for us to enforce that the returned vector  $\mathbf{x}$  from the system is indeed the hitting times (possibly shifted by the value assigned to coordinate  $t$ ).

**Finding Hitting Times via Laplacian System Solve.** Since we assume  $G$  connected, we have that multiplying with  $\mathbf{D} = \mathbf{D}^\top$  preserves equality. Further since  $\mathbf{W} = \mathbf{A}\mathbf{D}^{-1}$ , we obtain

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

Defining  $\mathbf{b} = \mathbf{d} - \alpha \cdot \mathbf{d}(s) \cdot \mathbf{1}_s$ , and observing  $\mathbf{L} = \mathbf{D} - \mathbf{A}$ , we have  $\mathbf{L}\mathbf{h} = \mathbf{b}$ .

Finally, we observe that we only have a solution to the above system if and only if  $\mathbf{b} \in \ker(\mathbf{L})^\perp = \text{span}(\mathbf{1})^\perp$ . We thus have to set  $\alpha$  such that

$$\mathbf{1}^\top (\mathbf{d} - \alpha \cdot \mathbf{d}(s) \cdot \mathbf{1}_s) = \|\mathbf{d}\|_1 - \alpha \cdot \mathbf{d}(s) \iff \alpha = \|\mathbf{d}\|_1 / \mathbf{d}(s).$$

We have now formalized  $\mathbf{L}$  and  $\mathbf{b}$  completely. A last detail that we should not forget about is that any solution  $\mathbf{x}$  to such system  $\mathbf{L}\mathbf{x} = \mathbf{b}$  is not necessarily equal  $\mathbf{h}$  but has the property that it is shifted from  $\mathbf{h}$  by the all-ones vector. Since we require  $\mathbf{h}(s) = 0$ , we can reconstruct  $\mathbf{h}$  from  $\mathbf{x}$  straight-forwardly by subtracting  $\mathbf{x}(s)\mathbf{1}$ .

**Theorem 8.3.1.** *Given a connected graph  $G$ , a special vertex  $s \in V$ . Then, we can formalize a Laplacian system  $\mathbf{L}\mathbf{x} = \mathbf{b}$  (where  $\mathbf{L}$  is the Laplacian of  $G$ ) such that the expected hitting times to  $s$  are given by  $\mathbf{h} = \mathbf{x} - \mathbf{x}(s)\mathbf{1}$ . We can reconstruct  $\mathbf{h}$  from  $\mathbf{x}$  in time  $O(n)$ .*

**Hitting Times and Electrical Networks.** Seeing that hitting times can be computed by formulating a Laplacian system  $\mathbf{L}\mathbf{x} = \mathbf{b}$ . You might remember that in the first lecture, we argued that a system  $\mathbf{L}\mathbf{x} = \mathbf{b}$  also solves the problem of routing a demand  $\mathbf{b}$  via an electrical flow with voltages  $\mathbf{x}$ .

Indeed, we can interpret computing expected hitting times  $\mathbf{h}$  to a special vertex  $s$  as the problem of computing the electrical voltages  $\mathbf{x}$  where we insert (or more technically correct apply)  $\mathbf{d}(a)$  units of current at every vertex  $a \neq s$  and where we remove  $\mathbf{1}^\top \mathbf{d} - \mathbf{d}(s)$  units of current at the vertex  $s$ . Then, we can express expected hitting time to some vertex  $a$  as the voltage difference to  $s$ :  $\mathbb{E}[H_{a,s}] = \mathbf{h}(a) = \mathbf{x}(a) - \mathbf{x}(s)$ .

### 8.3.2 Commute Time

A topic very related to hitting times are *commute times*. That is for two vertices  $a, b$ , the commute time is the time in a random walk starting in  $a$  to visit  $b$  and then to return to  $a$  again. Thus, it can be defined  $C_{a,b} = H_{a,b} + H_{b,a}$ .

**Commute Times via Electric Flows.** Recall that expected hitting times have an electric flow interpretation.

Now, let us denote by  $\mathbf{x}$  a solution to the Laplacian system  $\mathbf{Lx} = \mathbf{b}_b$  where the demand is  $\mathbf{b}_b = \mathbf{d} - \mathbf{d}^\top \mathbf{1} \cdot \mathbf{1}_b \in \ker(\mathbf{1})^\perp$ . Recall that we have  $\mathbb{E}[H_{z,b}] = \mathbf{x}(z) - \mathbf{x}(b)$  for all  $z$ . Similarly, we can compute voltages  $\mathbf{y}$  to the Laplacian system  $\mathbf{Ly} = \mathbf{b}_a$  where  $\mathbf{b}_a = \mathbf{d} - \mathbf{d}^\top \mathbf{1} \cdot \mathbf{1}_a \in \ker(\mathbf{1})^\perp$ . Again,  $\mathbb{E}[H_{z,a}] = \mathbf{y}(z) - \mathbf{y}(a)$  for all  $z$ .

But observe that this allows us to argue by linearity that  $\mathbb{E}[C_{a,b}] = \mathbb{E}[H_{a,b} + H_{b,a}] = \mathbb{E}[H_{a,b}] + \mathbb{E}[H_{b,a}] = \mathbf{x}(a) - \mathbf{x}(b) + \mathbf{y}(b) - \mathbf{y}(a) = (\mathbf{1}_a - \mathbf{1}_b)^\top (\mathbf{x} - \mathbf{y})$ . But using linearity again, we can also argue at the same time that we obtain the vector  $\mathbf{z} = (\mathbf{x} - \mathbf{y})$  as a solution to the problem  $\mathbf{Lz} = \mathbf{b}_b - \mathbf{b}_a = \mathbf{d}^\top \mathbf{1} (\mathbf{1}_a - \mathbf{1}_b)$ . That is the flow that routes  $\|\mathbf{d}\|_1$  units of flow from  $b$  to  $a$ .

**Theorem 8.3.2.** *Given a graph  $G = (V, E)$ , for any two fixed vertices  $a, b \in V$ , the expected commute time  $C_{a,b}$  is given by the voltage difference between  $a$  and  $b$  for any solution  $\mathbf{z}$  to the Laplacian system  $\mathbf{Lz} = \|\mathbf{d}\|_1 \cdot (\mathbf{1}_b - \mathbf{1}_a)$ .*

We note that the voltage difference between  $a$  and  $b$  in an electrical flow routing demand  $\mathbf{1}_b - \mathbf{1}_a$  is also called the *effective resistance*  $R_{\text{eff}}(a, b)$ . This quantity will play a crucial role in the next roles. In the next lecture, we introduce  $R_{\text{eff}}(a, b)$  slightly differently as the energy required by the electrical flow that routes  $\mathbf{1}_b - \mathbf{1}_a$ , however, it is not hard to show that these two definitions are equivalent.

Our theorem can now be restated as saying that the expected commute time  $\mathbb{E}[C_{a,b}] = \|\mathbf{d}\|_1 \cdot R_{\text{eff}}(a, b)$  which is a classic result.

# Chapter 9

## Random Matrix Concentration and Spectral Graph Sparsification

### 9.1 Matrix Sampling and Approximation

We want to begin understanding how sums of random matrices behave, in particular, whether they exhibit a tendency to concentrate in the same way that sum of scalar random variables do under various conditions.

First, let's recall a scalar Chernoff bound, which shows that a sum of bounded, non-negative random variables tend to concentrate around their mean.

**Theorem 9.1.1** (A Chernoff Concentration Bound). *Suppose  $X_1, \dots, X_k \in \mathbb{R}$  are independent, non-negative, random variables with  $X_i \leq R$  always. Let  $X = \sum_i X_i$ , and  $\mu = \mathbb{E}[X]$ , then for  $0 < \epsilon \leq 1$*

$$\Pr[X \geq (1 + \epsilon)\mu] \leq \exp\left(\frac{-\epsilon^2\mu}{4R}\right) \text{ and } \Pr[X \leq (1 - \epsilon)\mu] \leq \exp\left(\frac{-\epsilon^2\mu}{4R}\right).$$

The Chernoff bound should be familiar to most of you, but you may not have seen the following very similar bound. The Bernstein bound, which we will state in terms of zero-mean variables, is much like the Chernoff bound. It also requires bounded variables. But, when the variables have small variance, the Bernstein bound is sometimes stronger.

**Theorem 9.1.2** (A Bernstein Concentration Bound). *Suppose  $X_1, \dots, X_k \in \mathbb{R}$  are independent, zero-mean, random variables with  $|X_i| \leq R$  always. Let  $X = \sum_i X_i$ , and  $\sigma^2 = \text{Var}[X] = \sum_i \mathbb{E}[X_i^2]$ , then for  $t > 0$*

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

We will now prove the Bernstein concentration bound for scalar random variables, as a warm-up to the next section, where we will prove a version of it for matrix-valued random variables.

To help us prove Bernstein's bound, first let's recall Markov's inequality. This is a very weak concentration inequality, but also very versatile, because it requires few assumptions.

**Lemma 9.1.3** (Markov's Inequality). *Suppose  $X \in \mathbb{R}$  is a non-negative random variable, with a finite expectation. Then for any  $t > 0$ ,*

$$\Pr[X \geq t] \leq \frac{\mathbb{E}[X]}{t}.$$

*Proof.*

$$\begin{aligned}\mathbb{E}[X] &= \Pr[X \geq t] \mathbb{E}[X | X \geq t] + \Pr[X < t] \mathbb{E}[X | X < t] \\ &\geq \Pr[X \geq t] \mathbb{E}[X | X \geq t] \\ &\geq \Pr[X \geq t] \cdot t.\end{aligned}$$

We can rearrange this to get the desired statement.  $\square$

Now, we are ready to prove Bernstein's bound.

*Proof of Theorem 9.1.2.* We will focus on bounding the probability that  $\Pr[X \geq t]$ . The proof that  $\Pr[-X \geq t]$  is small proceeds in the same way.

First we observe that

$$\begin{aligned}\Pr[X \geq t] &= \Pr[\exp(\theta X) \geq \exp(\theta t)] \\ &\quad \text{for any } \theta > 0, \text{ because } x \mapsto \exp(\theta x) \text{ is strictly increasing.} \\ &\leq \exp(-\theta t) \mathbb{E}[\exp(\theta X)] \quad \text{by Lemma 9.1.3 (Markov's Inequality).}\end{aligned}$$

Now, let's require that  $\theta \leq 1/R$ . This will allow us to use the following bound: For all  $|z| \leq 1$ ,

$$\exp(z) \leq 1 + z + z^2. \tag{9.1}$$

We omit a proof of this, but the plots in Figure 9.1 suggest that this upper bound holds. The reader should consider how to prove this. With this in mind, we see that

$$\begin{aligned}\mathbb{E}[\exp(\theta X)] &= \mathbb{E}\left[\exp\left(\theta \sum_i X_i\right)\right] \\ &= \mathbb{E}[\prod_i \exp(\theta X_i)] \\ &= \prod_i \mathbb{E}[\exp(\theta X_i)] \quad \text{because } \mathbb{E}[YZ] = \mathbb{E}[Y]\mathbb{E}[Z] \text{ for independent } Y \text{ and } Z. \\ &\leq \prod_i \mathbb{E}[1 + \theta X_i + (\theta X_i)^2] \\ &= \prod_i (1 + \theta^2 \mathbb{E}[X_i^2]) \quad \text{because } X_i \text{ are zero-mean.} \\ &\leq \prod_i \exp(\theta^2 \mathbb{E}[X_i^2]) \quad \text{because } 1 + z \leq \exp(z) \text{ for all } z \in R. \\ &= \exp\left(\sum_i \theta^2 \mathbb{E}[X_i^2]\right) = \exp(\theta^2 \sigma^2).\end{aligned}$$

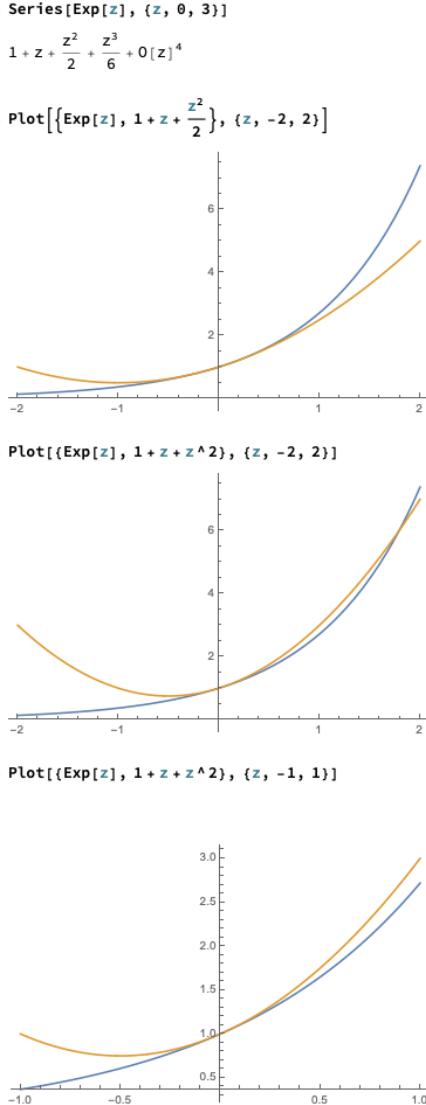


Figure 9.1: Plotting  $\exp(z)$  compared to  $1 + z + z^2$ .

Thus  $\Pr[X \geq t] \leq \exp(-\theta t) \mathbb{E}[\exp(\theta X)] \leq \exp(-\theta t + \theta^2 \sigma^2)$ . Now, to get the best possible bound, we'd like to minimize  $-\theta t + \theta^2 \sigma^2$  subject to the constraint  $0 < \theta \leq 1/R$ . Setting

$$\frac{\partial}{\partial \theta} (-\theta t + \theta^2 \sigma^2) = -t + 2\theta \sigma^2.$$

Setting this derivative to zero gives  $\theta = \frac{t}{2\sigma^2}$ , and plugging that in gives

$$-\theta t + \theta^2 \sigma^2 = -\frac{t^2}{4\sigma^2}$$

This choice only satisfies our constraints on  $\theta$  if  $\frac{t}{2\sigma^2} \leq 1/R$ . Otherwise, we let  $\theta = 1/R$  and note that in this case

$$-\theta t + \theta^2 \sigma^2 = -\frac{t}{R} + \frac{\sigma^2}{R^2} \leq -\frac{t}{R} + \frac{t}{2R} = -\frac{t}{2R}$$

where we got the inequality from  $t > 2\sigma^2/R$ . Altogether, we can conclude that there always is a choice of  $\theta$  s.t.

$$-\theta t + \theta^2 \sigma^2 \leq -\min\left(\frac{t}{2R}, \frac{t^2}{4\sigma^2}\right) \leq -\frac{t^2}{2Rt + 4\sigma^2}.$$

In fact, with the benefit of hindsight, and a little algebra, we arrive at the same conclusion in another way: One can check that the following choice of  $\theta$  is always valid and achieves the same bound:  $\theta = \frac{1}{2\sigma^2} \left( t - \frac{\sqrt{R} \cdot t^{3/2}}{\sqrt{2\sigma^2 + Rt}} \right)$ .  $\square$

We use  $\|\cdot\|$  to denote the spectral norm on matrices. Let's take a look at a version of Bernstein's bound that applies to sums of random matrices.

**Theorem 9.1.4** (A Bernstein Matrix Concentration Bound (Tropp 2011)). *Suppose  $\mathbf{X}_1, \dots, \mathbf{X}_k \in \mathbb{R}^{n \times n}$  are independent, symmetric matrix-valued random variables. Assume each  $\mathbf{X}_i$  is zero-mean, i.e.  $\mathbb{E}[\mathbf{X}_i] = \mathbf{0}_{n \times n}$ , and that  $\|\mathbf{X}_i\| \leq R$  always. Let  $\mathbf{X} = \sum_i \mathbf{X}_i$ , and  $\sigma^2 = \text{Var}[\mathbf{X}] = \sum_i \mathbb{E}[\mathbf{X}_i^2]$ , then for  $t > 0$*

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

This basically says that probability of  $\mathbf{X}$  being large in spectral norm behaves like the scalar case, except the bound is larger by a factor  $n$ , where the matrices are  $n \times n$ . We can get a feeling for why this might be a reasonable bound by considering the case of random diagonal matrices. Now  $\|\mathbf{X}\| = \max_j |\mathbf{X}(j, j)| = \max_j |\sum_i \mathbf{X}_i(j, j)|$ . In this case, we need to bound the largest of the  $n$  diagonal entries: We can do this by a union bound over  $n$  instances of the scalar problem – and this also turns out to be essentially tight in some cases, meaning we can't expect a better bound in general.

## 9.2 Matrix Concentration

In this section we will prove the Bernstein matrix concentration bound (Tropp 2011) that we saw in the previous section (Theorem 9.1.4).

**Theorem 9.2.1.** *Suppose  $\mathbf{X}_1, \dots, \mathbf{X}_k \in \mathbb{R}^{n \times n}$  are independent, symmetric matrix-valued random variables. Assume each  $\mathbf{X}_i$  is zero-mean, i.e.  $\mathbb{E}[\mathbf{X}_i] = \mathbf{0}_{n \times n}$ , and that  $\|\mathbf{X}_i\| \leq R$  always. Let  $\mathbf{X} = \sum_i \mathbf{X}_i$ , and  $\sigma^2 = \|\text{Var}[\mathbf{X}]\| = \|\sum_i \mathbb{E}[\mathbf{X}_i^2]\|$ , then for  $t > 0$*

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

But let's collect some useful tools for the proof first.

**Definition 9.2.2** (trace). The trace of a square matrix  $\mathbf{A}$  is defined as

$$\text{Tr}(\mathbf{A}) := \sum_i \mathbf{A}(i, i)$$

**Claim 9.2.3** (cyclic property of trace).  $\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA})$

Let  $S^n$  denote the set of all  $n \times n$  real symmetric matrices,  $S_+^n$  the set of all  $n \times n$  positive semidefinite matrices, and  $S_{++}^n$  the set of all  $n \times n$  positive definite matrices. Their relation is clear,  $S_{++}^n \subset S_+^n \subset S^n$ . For any  $\mathbf{A} \in S^n$  with eigenvalues  $\lambda_1(\mathbf{A}) \leq \dots \leq \lambda_n(\mathbf{A})$ , by spectral decomposition theorem,  $\mathbf{A} = \mathbf{V}\Lambda\mathbf{V}^\top$  where  $\Lambda = \text{diag}_i\{\lambda_i(\mathbf{A})\}$  and  $\mathbf{V}^\top\mathbf{V} = \mathbf{V}\mathbf{V}^\top = \mathbf{I}$ , we'll use this property without specifying in the sequel.

**Claim 9.2.4.** Given a symmetric and real matrix  $\mathbf{A}$ ,  $\text{Tr}(\mathbf{A}) = \sum_i \lambda_i$ , where  $\{\lambda_i\}$  are eigenvalues of  $\mathbf{A}$ .

*Proof.*

$$\text{Tr}(\mathbf{A}) = \text{Tr}(\mathbf{V}\Lambda\mathbf{V}^\top) = \text{Tr}\left(\Lambda \underbrace{\mathbf{V}^\top\mathbf{V}}_{\mathbf{I}}\right) = \text{Tr}(\Lambda) = \sum_i \lambda_i.$$

□

### 9.2.1 Matrix Functions

**Definition 9.2.5** (Matrix function). Given a real-valued function  $f : \mathbb{R} \rightarrow \mathbb{R}$ , we extend it to a matrix function  $f : S^n \rightarrow S^n$ . For  $\mathbf{A} \in S^n$  with spectral decomposition  $\mathbf{A} = \mathbf{V}\Lambda\mathbf{V}^\top$ , let

$$f(\mathbf{A}) = \mathbf{V} \underset{i}{\text{diag}} \{f(\lambda_i)\} \mathbf{V}^\top.$$

**Example.** Recall that every PSD matrix  $\mathbf{A}$  has a square root  $\mathbf{A}^{1/2}$ . If  $f(x) = x^{1/2}$  for  $x \in \mathbb{R}_+$ , then  $f(\mathbf{A}) = \mathbf{A}^{1/2}$  for  $\mathbf{A} \in S_+^n$ .

**Example.** If  $f(x) = \exp(x)$  for  $x \in \mathbb{R}$ , then  $f(\mathbf{A}) = \exp(\mathbf{A}) = \mathbf{V} \exp(\Lambda) \mathbf{V}^\top$  for  $\mathbf{A} \in S^n$ . Note that  $\exp(\mathbf{A})$  is positive definite for any  $\mathbf{A} \in S^n$ .

### 9.2.2 Monotonicity and Operator Monotonicity

Consider a function  $f : \mathcal{D} \rightarrow \mathcal{C}$ . If we have a partial order  $\leq_{\mathcal{D}}$  defined on  $\mathcal{D}$  and a partial order  $\leq_{\mathcal{C}}$  defined on  $\mathcal{C}$ , then we say that the function is monotone increasing (resp. decreasing) w.r.t. this pair of orderings if for all  $d_1, d_2 \in \mathcal{D}$  s.t.  $d_1 \leq_{\mathcal{D}} d_2$  we have  $f(d_1) \leq_{\mathcal{C}} f(d_2)$  (resp. decreasing if  $f(d_2) \leq_{\mathcal{C}} f(d_1)$ ).

Let's introduce some terminology for important special cases of this idea. We say that a function  $f : \mathcal{S} \rightarrow \mathbb{R}$ , where  $\mathcal{S} \subseteq S^n$ , is monotone increasing if  $\mathbf{A} \preceq \mathbf{B}$  implies  $f(\mathbf{A}) \leq f(\mathbf{B})$ .

Meanwhile, a function  $f : \mathcal{S} \rightarrow \mathcal{T}$  where  $\mathcal{S}, \mathcal{T} \subseteq S^n$  is said to be operator monotone increasing if  $\mathbf{A} \preceq \mathbf{B}$  implies  $f(\mathbf{A}) \preceq f(\mathbf{B})$ .

**Lemma 9.2.6.** *Let  $T \subseteq \mathbb{R}$ . If the scalar function  $f : T \rightarrow \mathbb{R}$  is monotone increasing, the matrix function  $\mathbf{X} \mapsto \text{Tr}(f(\mathbf{X}))$  is monotone increasing.*

*Proof.* From previous chapters, we know if  $\mathbf{A} \preceq \mathbf{B}$  then  $\lambda_i(\mathbf{A}) \leq \lambda_i(\mathbf{B})$  for all  $i$ . As  $x \mapsto f(x)$  is monotone, then  $\lambda_i(f(\mathbf{A})) \leq \lambda_i(f(\mathbf{B}))$  for all  $i$ . By Claim 9.2.4,  $\text{Tr}(f(\mathbf{A})) \leq \text{Tr}(f(\mathbf{B}))$ .  $\square$

From this, and the fact that  $x \mapsto \exp(x)$  is a monotone function on the reals, we get the following corollary.

**Corollary 9.2.7.** *If  $\mathbf{A} \preceq \mathbf{B}$ , then  $\text{Tr}(\exp(\mathbf{A})) \leq \text{Tr}(\exp(\mathbf{B}))$ , i.e.  $\mathbf{X} \mapsto \text{Tr}(\exp(\mathbf{X}))$  is monotone increasing.*

**Lemma 9.2.8.** *If  $\mathbf{0} \prec \mathbf{A} \preceq \mathbf{B}$ , then  $\mathbf{B}^{-1} \preceq \mathbf{A}^{-1}$ , i.e.  $\mathbf{X} \mapsto \mathbf{X}^{-1}$  is operator monotone decreasing on  $S_{++}^n$ .*

You will prove the above lemma in this week's exercises.

**Lemma 9.2.9.** *If  $\mathbf{0} \prec \mathbf{A} \preceq \mathbf{B}$ , then  $\log(\mathbf{A}) \preceq \log(\mathbf{B})$ .*

To prove this lemma, we first recall an integral representation of the logarithm.

**Lemma 9.2.10.**

$$\log a = \int_0^\infty \left( \frac{1}{1+t} - \frac{1}{a+t} \right) dt$$

*Proof.*

$$\begin{aligned} \int_0^\infty \left( \frac{1}{1+t} - \frac{1}{a+t} \right) dt &= \lim_{T \rightarrow \infty} \int_0^T \left( \frac{1}{1+t} - \frac{1}{a+t} \right) dt \\ &= \lim_{T \rightarrow \infty} [\log(1+t) - \log(a+t)]_0^T \\ &= \log(a) + \lim_{T \rightarrow \infty} \log \left( \frac{1+T}{a+T} \right) \\ &= \log(a) \end{aligned}$$

$\square$

*Proof sketch of Lemma 9.2.9.* Because all the matrices involved are diagonalized by the same orthogonal transformation, we can conclude from Lemma 9.2.10 that for a matrix  $\mathbf{A} \succ \mathbf{0}$ ,

$$\log(\mathbf{A}) = \int_0^\infty \left( \frac{1}{1+t} \mathbf{I} - (t\mathbf{I} + \mathbf{A})^{-1} \right) dt$$

This integration can be expressed as the limit of a sum with positive coefficients, and from this we can show that the integrand (the term inside the integration symbol) is operator monotone increasing in  $\mathbf{A}$  by Lemma 9.2.8, the result of the integral, i.e.  $\log(\mathbf{A})$  must also be operator monotone increasing.  $\square$

The following is a more general version of Lemma 1.6.

**Lemma 9.2.11.** *Let  $T \subset \mathbb{R}$ . If the scalar function  $f : T \rightarrow \mathbb{R}$  is monotone, the matrix function  $\mathbf{X} \mapsto \text{Tr}(f(\mathbf{X}))$  is monotone.*

**Remark 9.2.12.** It is not always true that when  $f : \mathbb{R} \rightarrow \mathbb{R}$  is monotone,  $f : S^n \rightarrow S^n$  is operator monotone. For example,  $\mathbf{X} \mapsto \mathbf{X}^2$  and  $\mathbf{X} \mapsto \exp(\mathbf{X})$  are *not* operator monotone.

### 9.2.3 Some Useful Facts

**Lemma 9.2.13.**  $\exp(\mathbf{A}) \preceq \mathbf{I} + \mathbf{A} + \mathbf{A}^2$  for  $\|\mathbf{A}\| \leq 1$ .

*Proof.*

$$\begin{aligned} \mathbf{I} + \mathbf{A} + \mathbf{A}^2 - \exp(\mathbf{A}) &= \mathbf{V} \mathbf{I} \mathbf{V}^\top + \mathbf{V} \mathbf{A} \mathbf{V}^\top + \mathbf{V} \mathbf{A}^2 \mathbf{V}^\top - \mathbf{V} \exp(\mathbf{A}) \mathbf{V}^\top \\ &= \mathbf{V} (\mathbf{I} + \mathbf{A} + \mathbf{A}^2 - \exp(\mathbf{A})) \mathbf{V}^\top \\ &= \mathbf{V} \underset{i}{\text{diag}} \{1 + \lambda_i + \lambda_i^2 - \exp(\lambda_i)\} \mathbf{V}^\top \end{aligned}$$

Recall  $\exp(x) \leq 1 + x + x^2$  for all  $|x| \leq 1$ . Since  $\|\mathbf{A}\| \leq 1$  i.e.  $|\lambda_i| \leq 1$  for all  $i$ , thus  $1 + \lambda_i + \lambda_i^2 - \exp(\lambda_i) \geq 0$  for all  $i$ , meaning  $\mathbf{I} + \mathbf{A} + \mathbf{A}^2 - \exp(\mathbf{A}) \succeq 0$ .  $\square$

**Lemma 9.2.14.**  $\log(\mathbf{I} + \mathbf{A}) \preceq \mathbf{A}$  for  $\mathbf{A} \succ -\mathbf{I}$ .

*Proof.*

$$\begin{aligned} \mathbf{A} - \log(\mathbf{I} + \mathbf{A}) &= \mathbf{V} \mathbf{A} \mathbf{V}^\top - \mathbf{V} \log(\mathbf{A} + \mathbf{I}) \mathbf{V}^\top \\ &= \mathbf{V} (\mathbf{A} - \log(\mathbf{A} + \mathbf{I})) \mathbf{V}^\top \\ &= \mathbf{V} \underset{i}{\text{diag}} \{\lambda_i - \log(1 + \lambda_i)\} \mathbf{V}^\top \end{aligned}$$

Recall  $x \geq \log(1 + x)$  for all  $x > -1$ . Since  $\|\mathbf{A}\| \succ -\mathbf{I}$  i.e.  $\lambda_i > -1$  for all  $i$ , thus  $\lambda_i - \log(1 + \lambda_i) \geq 0$  for all  $i$ , meaning  $\mathbf{A} - \log(\mathbf{I} + \mathbf{A}) \succeq 0$ .  $\square$

**Theorem 9.2.15** (Lieb). *Let  $f : S_{++}^n \rightarrow \mathbb{R}$  be a matrix function given by*

$$f(\mathbf{A}) = \text{Tr}(\exp(\mathbf{H} + \log(\mathbf{A})))$$

*for some  $\mathbf{H} \in S^n$ . Then  $-f$  is convex (i.e.  $f$  is concave).*

The Lieb's theorem will be crucial in our proof of Theorem 9.1.4, but it is also highly non-trivial and we will omit its proof here. The interested reader can find a proof in Chapter 8 of [T<sup>+</sup>15].

**Lemma 9.2.16** (Jensen's inequality).  $\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$  when  $f$  is convex;  $\mathbb{E}[f(X)] \leq f(\mathbb{E}[X])$  when  $f$  is concave.

### 9.2.4 Proof of Matrix Bernstein Concentration Bound

Now, we are ready to prove the Bernstein matrix concentration bound.

*Proof of Theorem 9.1.4.* For any  $\mathbf{A} \in S^n$ , its spectral norm  $\|\mathbf{A}\| = \max\{|\lambda_n(\mathbf{A})|, |\lambda_1(\mathbf{A})|\} = \max\{\lambda_n(\mathbf{A}), -\lambda_1(\mathbf{A})\}$ . Let  $\lambda_1 \leq \dots \leq \lambda_n$  be the eigenvalues of  $\mathbf{X}$ . Then,

$$\Pr[\|\mathbf{X}\| \geq t] = \Pr[(\lambda_n \geq t) \vee (-\lambda_1 \geq t)] \leq \Pr[\lambda_n \geq t] + \Pr[-\lambda_1 \geq t].$$

Let  $\mathbf{Y} := \sum_i -\mathbf{X}_i$ . It's easy to see that  $-\lambda_n \leq \dots \leq -\lambda_1$  are eigenvalues of  $\mathbf{Y}$ , implying  $\lambda_n(\mathbf{Y}) = -\lambda_1(\mathbf{X})$ . Since  $\mathbb{E}[-\mathbf{X}_i] = \mathbb{E}[\mathbf{X}_i] = 0$  and  $\|-\mathbf{X}_i\| = \|\mathbf{X}_i\| \leq R$  for all  $i$ , if we can bound  $\Pr[\lambda_n(\mathbf{X}) \geq t]$ , then applying to  $\mathbf{Y}$ , we can bound  $\Pr[\lambda_n(\mathbf{Y}) \geq t]$ . As

$$\Pr[-\lambda_1(\mathbf{X}) \geq t] = \Pr[\lambda_n(\mathbf{Y}) \geq t],$$

it suffices to bound  $\Pr[\lambda_n \geq t]$ .

For any  $\theta > 0$ ,  $\lambda_n \geq t \iff \exp(\theta\lambda_n) \geq \exp(\theta t)$  and  $\text{Tr}(\exp(\theta\mathbf{X})) = \sum_i \exp(\theta\lambda_i)$  by Claim 9.2.4, thus  $\lambda_n \geq t \Rightarrow \text{Tr}(\exp(\theta\mathbf{X})) \geq \exp(\theta t)$ . Then, using Markov's inequality,

$$\begin{aligned} \Pr[\lambda_n \geq t] &\leq \Pr[\text{Tr}(\exp(\theta\mathbf{X})) \geq \exp(\theta t)] \\ &\leq \exp(-\theta t) \mathbb{E}[\text{Tr}(\exp(\theta\mathbf{X}))] \end{aligned}$$

For two independent random variables  $\mathbf{U}$  and  $\mathbf{V}$ , we have

$$\mathbb{E}_{\mathbf{U}, \mathbf{V}} f(\mathbf{U}, \mathbf{V}) = \mathbb{E}_{\mathbf{U}} \mathbb{E}_{\mathbf{V}} [f(\mathbf{U}, \mathbf{V}) | \mathbf{U}] = \mathbb{E}_{\mathbf{U}} \mathbb{E}_{\mathbf{V}} [f(\mathbf{U}, \mathbf{V})].$$

Define  $\mathbf{X}_{< i} = \sum_{j < i} \mathbf{X}_j$ . Let  $0 < \theta \leq 1/R$ ,

$$\begin{aligned}
\mathbb{E} \operatorname{Tr}(\exp(\theta \mathbf{X})) &= \underset{\mathbf{X}_1, \dots, \mathbf{X}_{k-1}}{\mathbb{E}} \underset{\mathbf{X}_k}{\mathbb{E}} \operatorname{Tr} \exp \left( \underbrace{\theta \mathbf{X}_{<k}}_{\mathbf{H}} + \underbrace{\theta \mathbf{X}_k}_{=\log \exp(\theta \mathbf{X}_k)} \right), \quad \{\mathbf{X}_i\} \text{ are independent} \\
&\leq \underset{\mathbf{X}_1, \dots, \mathbf{X}_{k-1}}{\mathbb{E}} \operatorname{Tr} \exp \left( \theta \mathbf{X}_{<k} + \log \mathbb{E} \exp(\theta \mathbf{X}_k) \right), \quad \text{by 9.2.15 and 9.2.16} \\
&\leq \underset{\mathbf{X}_1, \dots, \mathbf{X}_{k-1}}{\mathbb{E}} \operatorname{Tr} \exp \left( \theta \mathbf{X}_{<k} + \log \mathbb{E} [\mathbf{I} + \theta \mathbf{X}_k + \theta^2 \mathbf{X}_k^2] \right), \quad \text{by 9.2.13, 9.2.7, and 9.2.9} \\
&\leq \underset{\mathbf{X}_1, \dots, \mathbf{X}_{k-1}}{\mathbb{E}} \operatorname{Tr} \exp \left( \theta \mathbf{X}_{<k} + \theta^2 \mathbb{E} \mathbf{X}_k^2 \right), \quad \text{by 9.2.14 and 9.2.7} \\
&= \underset{\mathbf{X}_1, \dots, \mathbf{X}_{k-2}}{\mathbb{E}} \underset{\mathbf{X}_{k-1}}{\mathbb{E}} \operatorname{Tr} \exp \left( \underbrace{\theta^2 \mathbb{E} \mathbf{X}_k^2 + \theta \mathbf{X}_{<k-1} + \theta \mathbf{X}_{k-1}}_{\mathbf{H}} \right), \\
&\vdots \\
&\leq \operatorname{Tr} \exp \left( \theta^2 \sum_i \mathbb{E} [\mathbf{X}_i^2] \right), \\
&\leq \operatorname{Tr} \exp (\theta^2 \sigma^2 \mathbf{I}), \quad \text{by 9.2.7 and } \sum_i \mathbb{E} [\mathbf{X}_i^2] \preceq \sigma^2 \mathbf{I} \\
&= n \cdot \exp(\theta^2 \sigma^2).
\end{aligned}$$

Then,

$$\Pr[\lambda_n \geq t] \leq n \cdot \exp(-\theta t + \theta^2 \sigma^2),$$

and

$$\Pr[\|\mathbf{X}\| \geq t] \leq 2n \cdot \exp(-\theta t + \theta^2 \sigma^2).$$

Similar to the proof of Bernstein concentration bound for one-dimension random variable, minimize the RHS over  $0 < \theta \leq 1/R$  yields

$$\Pr[\|\mathbf{X}\| \geq t] \leq 2n \cdot \exp \left( \frac{-t^2}{2Rt + 4\sigma^2} \right).$$

□

### 9.3 Spectral Graph Sparsification

In this section, we will see that for any dense graph, we can find another sparser graph whose graph Laplacian is approximately the same as measured by their quadratic forms. This turns out to be a very useful tool for designing algorithms.

**Definition 9.3.1.** Given  $\mathbf{A}, \mathbf{B} \in S_+^n$  and  $\epsilon > 0$ , we say

$$\mathbf{A} \approx_\epsilon \mathbf{B} \text{ if and only if } \frac{1}{1+\epsilon} \mathbf{A} \preceq \mathbf{B} \preceq (1+\epsilon) \mathbf{A}.$$

Suppose we start with a connected graph  $G = (V, E, \mathbf{w})$ , where as usual we say that  $|V| = n$  and  $|E| = m$ . We want to produce another graph  $\tilde{G} = (V, \tilde{E}, \tilde{\mathbf{w}})$  s.t  $|\tilde{E}| \ll |E|$  and at the same time  $\mathbf{L}_G \approx_\epsilon \mathbf{L}_{\tilde{G}}$ . We call  $\tilde{G}$  a *spectral sparsifier* of  $G$ . Our construction will also ensure that  $\tilde{E} \subseteq E$ , although this is not important in most applications. Figure 9.2 shows an example of a graph  $G$  and spectral sparsifier  $\tilde{G}$ .

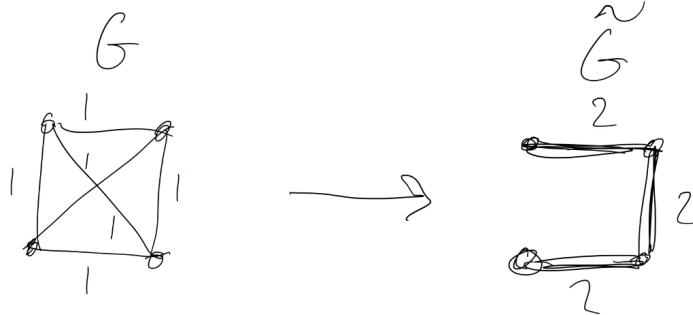


Figure 9.2: A graph  $G$  and a spectral sparsifier  $\tilde{G}$ , satisfying  $\mathbf{L}_G \approx_\epsilon \mathbf{L}_{\tilde{G}}$  for  $\epsilon = 2.42$ .

We are going to construct  $\tilde{G}$  by sampling some of the edges of  $G$  according to a suitable probability distribution and scaling up their weight to make up for the fact that we pick fewer of them.

To get a better understanding for the notion of approximation given in 9.3.1 means, let's observe a simple consequence of it.

Given a vertex subset  $T \subseteq V$ , we say that  $(T, V \setminus T)$  is a *cut* in  $G$  and that the value of the cut is

$$c_G(T) = \sum_{e \in E \cap (T \times V \setminus T)} \mathbf{w}(e).$$

Figure 9.3 shows the  $c_G(T)$  in a graph  $G$ .

**Theorem 9.3.2.** If  $\mathbf{L}_G \approx_\epsilon \mathbf{L}_{\tilde{G}}$ , then for all  $T \subseteq V$ ,

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

*Proof.* Let  $\mathbf{1}_T \in \mathbb{R}^V$  be the indicator of the set  $T$ , i.e.  $\mathbf{1}_T(u) = 1$  for  $u \in T$  and  $\mathbf{1}_T(u) = 0$  otherwise. We can see that  $\mathbf{1}_T^\top \mathbf{L}_G \mathbf{1}_T = c_G(T)$ , and hence the theorem follows by comparing the quadratic forms.  $\square$

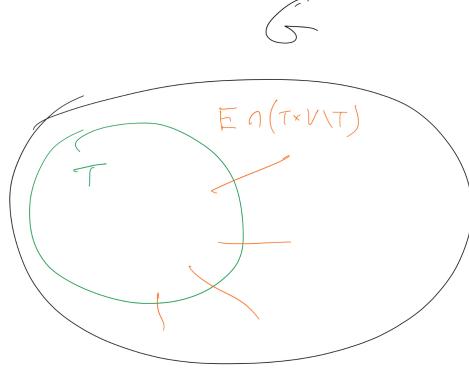


Figure 9.3: The cut  $c_G(T)$  in  $G$ .

But how well can we spectrally approximate a graph with a sparse graph? The next theorem gives us a nearly optimal answer to this question.

**Theorem 9.3.3** (Spectral Graph Approximation by Sampling, (Spielman-Srivastava 2008)). *Consider a connected graph  $G = (V, E, \mathbf{w})$ , with  $n = |V|$ . For any  $0 < \epsilon < 1$  and  $0 < \delta < 1$ , there exist sampling probabilities  $p_e$  for each edge  $e \in E$  s.t. if we include each edge  $e$  in  $\tilde{E}$  independently with probability  $p_e$  and set its weight  $\tilde{\mathbf{w}}(e) = \frac{1}{p_e} \mathbf{w}(e)$ , then with probability at least  $1 - \delta$  the graph  $\tilde{G} = (V, \tilde{E}, \tilde{\mathbf{w}})$  satisfies*

$$\mathbf{L}_G \approx_{\epsilon} \mathbf{L}_{\tilde{G}} \text{ and } |\tilde{E}| \leq O(n\epsilon^{-2} \log(n/\delta)).$$

The original proof can be found in [SS11].

**Remark 9.3.4.** For convenience, we will abbreviate  $\mathbf{L}_G$  as  $\mathbf{L}$  and  $\mathbf{L}_{\tilde{G}}$  as  $\tilde{\mathbf{L}}$  in the rest of this section.

We are going to analyze a sampling procedure by turning our goal into a problem of matrix concentration. Recall that

**Fact 9.3.5.**  $\mathbf{A} \preceq \mathbf{B}$  implies  $\mathbf{C}\mathbf{A}\mathbf{C}^\top \preceq \mathbf{C}\mathbf{B}\mathbf{C}^\top$  for any  $\mathbf{C} \in \mathbb{R}^{n \times n}$ .

By letting  $\mathbf{C} = \mathbf{L}^{+/2}$ , we can see that

$$\mathbf{L} \approx_{\epsilon} \tilde{\mathbf{L}} \text{ implies } \mathbf{\Pi}_L \approx_{\epsilon} \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2}, \quad (9.2)$$

where  $\mathbf{\Pi}_L = \mathbf{L}^{+/2} \mathbf{L} \mathbf{L}^{+/2}$  is the orthogonal projection to the complement of the kernel of  $\mathbf{L}$ .

**Definition 9.3.6.** Given a matrix  $\mathbf{A}$ , we define  $\mathbf{\Pi}_A$  to be the orthogonal projection to the complement of the kernel of  $\mathbf{A}$ , i.e.  $\mathbf{\Pi}_A \mathbf{v} = \mathbf{0}$  for  $\mathbf{v} \in \ker(\mathbf{A})$  and  $\mathbf{\Pi}_A \mathbf{v} = \mathbf{v}$  for  $\mathbf{v} \in \ker(\mathbf{A})^\perp$ . Recall that  $\ker(\mathbf{A})^\perp = \text{im}(\mathbf{A}^\top)$ .

**Claim 9.3.7.** For a matrix  $\mathbf{A} \in S^n$  with spectral decomposition  $\mathbf{A} = \mathbf{V} \Lambda \mathbf{V}^\top = \sum_i \lambda_i \mathbf{v}_i \mathbf{v}_i^\top$  s.t.  $\mathbf{V}^\top \mathbf{V} = \mathbf{I}$ , we have  $\mathbf{\Pi}_A = \sum_{i: \lambda_i \neq 0} \mathbf{v}_i \mathbf{v}_i^\top$ , and  $\mathbf{\Pi}_A = \mathbf{A}^{+/2} \mathbf{A} \mathbf{A}^{+/2} = \mathbf{A} \mathbf{A}^+ = \mathbf{A}^+ \mathbf{A}$ .

From the definition, we can see that  $\Pi_L = \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$ .

Now that we understand the projection  $\Pi_L$ , it is not hard to show the following claim.

**Claim 9.3.8.**

1.  $\Pi_L \approx_\epsilon \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2}$  implies  $\mathbf{L} \approx_\epsilon \tilde{\mathbf{L}}$ .
2. For  $\epsilon \leq 1$ , we have that  $\|\Pi_L - \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2}\| \leq \epsilon/2$  implies  $\Pi_L \approx_\epsilon \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2}$ .

Really, the only idea needed here is that when comparing quadratic forms in matrices with the same kernel, we necessarily can't have the quadratic forms disagree on vectors in the kernel. Simple! But we are going to write it out carefully, since we're still getting used to these types of calculations.

*Proof of Claim 9.3.8.* To prove Part 1, we assume  $\Pi_L \approx_\epsilon \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2}$ . Recall that  $G$  is a connected graph, so  $\ker(\mathbf{L}) = \text{span}\{\mathbf{1}\}$ , while  $\tilde{\mathbf{L}}$  is the Laplacian of a graph which may or may not be connected, so  $\ker(\tilde{\mathbf{L}}) \supseteq \ker(\mathbf{L})$ , and equivalently  $\text{im}(\tilde{\mathbf{L}}) \subseteq \text{im}(\mathbf{L})$ . Now, for any  $\mathbf{v} \in \ker(\mathbf{L})$  we have  $\mathbf{v}^\top \tilde{\mathbf{L}} \mathbf{v} = 0 = \mathbf{v}^\top \mathbf{L} \mathbf{v}$ . For any  $\mathbf{v} \in \ker(\mathbf{L})^\perp$  we have  $\mathbf{v} = \mathbf{L}^{+/2} \mathbf{z}$  for some  $\mathbf{z}$ , as  $\ker(\mathbf{L})^\perp = \text{im}(\mathbf{L}) = \text{im}(\mathbf{L}^{+/2})$ . Hence

$$\mathbf{v}^\top \tilde{\mathbf{L}} \mathbf{v} = \mathbf{z}^\top \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2} \mathbf{z} \geq \frac{1}{1+\epsilon} \mathbf{z}^\top \mathbf{L}^{+/2} \mathbf{L} \mathbf{L}^{+/2} \mathbf{z} = \frac{1}{1+\epsilon} \mathbf{v}^\top \mathbf{L} \mathbf{v}$$

and similarly

$$\mathbf{v}^\top \tilde{\mathbf{L}} \mathbf{v} = \mathbf{z}^\top \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2} \mathbf{z} \leq (1+\epsilon) \mathbf{z}^\top \mathbf{L}^{+/2} \mathbf{L} \mathbf{L}^{+/2} \mathbf{z} = (1+\epsilon) \mathbf{v}^\top \mathbf{L} \mathbf{v}.$$

Thus we have established  $\mathbf{L} \approx_\epsilon \tilde{\mathbf{L}}$ .

To prove Part 2, we assume  $\|\Pi_L - \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2}\| \leq \epsilon/2$ . This is equivalent to

$$-\frac{\epsilon}{2} \mathbf{I} \preceq \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2} - \Pi_L \preceq \frac{\epsilon}{2} \mathbf{I}$$

But since

$$\mathbf{1}^\top (\mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2} - \Pi_L) \mathbf{1} = 0,$$

we can in fact sharpen this to

$$-\frac{\epsilon}{2} \Pi_L \preceq \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2} - \Pi_L \preceq \frac{\epsilon}{2} \Pi_L.$$

Rearranging, we then conclude

$$(1 - \frac{\epsilon}{2}) \Pi_L \preceq \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2} \preceq (1 + \frac{\epsilon}{2}) \Pi_L.$$

Finally, we note that  $1/(1+\epsilon) \leq (1 - \frac{\epsilon}{2})$  to reach our conclusion,  $\Pi_L \approx_\epsilon \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2}$ .  $\square$

We now have most of the tools to prove Theorem 9.3.3, but to help us, we are going to establish one small piece of helpful notation: We define a matrix function  $\Phi : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$  by

$$\Phi(\mathbf{A}) = \mathbf{L}^{+/-2} \mathbf{A} \mathbf{L}^{+/-2}.$$

We sometimes call this a “normalizing map”, because it transforms a matrix to the space where spectral norm bounds can be translated into relative error guarantees compare to the  $\mathbf{L}$  quadratic form.

*Proof of Theorem 9.3.3.* By Claim 9.3.8, it suffices to show

$$\left\| \mathbf{\Pi}_L - L^{+/-2} \tilde{L} L^{+/-2} \right\| \leq \epsilon/2. \quad (9.3)$$

We introduce a set of independent random variables, one for each edge  $e$ , with a probability  $p_e$  associated with the edge which we will fix later. We then let

$$\mathbf{Y}_e = \begin{cases} \frac{w(e)}{p_e} \mathbf{b}_e \mathbf{b}_e^\top & \text{with probability } p_e \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

This way,  $\tilde{L} = \sum_e \mathbf{Y}_e$ . Note that  $\mathbb{E}[\mathbf{Y}_e] = p_e \frac{w(e)}{p_e} \mathbf{b}_e \mathbf{b}_e^\top = w(e) \mathbf{b}_e \mathbf{b}_e^\top$ , and so

$$\mathbb{E}[\tilde{L}] = \sum_e \mathbb{E}[\mathbf{Y}_e] = L.$$

By linearity of  $\Phi$ ,

$$\mathbb{E}[\Phi(\tilde{L})] = \Phi(\mathbb{E}[\tilde{L}]) = \mathbf{\Pi}_L.$$

Let us also define

$$\mathbf{X}_e = \Phi(\mathbf{Y}_e) - \mathbb{E}[\Phi(\mathbf{Y}_e)] \text{ and } \mathbf{X} = \sum_e \mathbf{X}_e$$

Note that this ensures  $\mathbb{E}[\mathbf{X}_e] = \mathbf{0}$ . We are now going to fix the edge sampling probabilities, in a way that depends on some overall scaling parameter  $\alpha > 0$ . We let

$$p_e = \min(\alpha \|\Phi(w(e) \mathbf{b}_e \mathbf{b}_e^\top)\|, 1)$$

then we see from the definition of  $\mathbf{Y}_e$  that whenever  $p_e < 1$

$$\|\Phi(\mathbf{Y}_e)\| \leq \frac{1}{\alpha}$$

from this, we can conclude, with a bit of work, that for all  $e$

$$\|\mathbf{X}_e\| \leq \frac{1}{\alpha}. \quad (9.4)$$

We can also show that

$$\left\| \sum_e \mathbb{E}[\mathbf{X}_e^2] \right\| \leq \frac{1}{\alpha}. \quad (9.5)$$

In the exercises for this chapter, we will ask you to show that Equations (9.4) and (9.5) hold.

This means that we can apply Theorem 9.1.4 to our  $\mathbf{X} = \sum_e \mathbf{X}_e$ , with  $R = \frac{1}{\alpha}$  and  $\sigma^2 = \frac{1}{\alpha}$ , to get

$$\Pr \left[ \left\| \mathbf{\Pi}_L - \mathbf{L}^{+/-} \tilde{\mathbf{L}} \mathbf{L}^{+/-} \right\| \geq \epsilon/2 \right] \leq 2n \exp \left( \frac{-0.25\epsilon^2}{(\epsilon+4)/\alpha} \right)$$

Since  $0 < \epsilon < 1$ , this means that if  $\alpha = 40\epsilon^{-2} \log(n/\delta)$ , then

$$\Pr \left[ \left\| \mathbf{\Pi}_L - \mathbf{L}^{+/-} \tilde{\mathbf{L}} \mathbf{L}^{+/-} \right\| \geq \epsilon/2 \right] \leq \frac{2n\delta^2}{n^2} \leq \delta/2.$$

In the last step, we assumed  $n \geq 4$ .

Lastly, we'd like to know that the graph  $\tilde{G}$  is sparse. The number of edges in  $\tilde{G}$  is equal to the number of  $\mathbf{Y}_e$  that come out nonzero. Thus, the expected value of  $|\tilde{E}|$  is

$$\mathbb{E} \left[ |\tilde{E}| \right] = \sum_e p_e \leq \alpha \sum_e \mathbf{w}(e) \left\| \mathbf{L}^{+/-} \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/-} \right\|$$

We can bound the sum of the norms with a neat trick relating it to the trace of  $\mathbf{\Pi}_L$ . Note that in general for a vector  $\mathbf{a} \in \mathbb{R}^n$ , we have  $\|\mathbf{a}\mathbf{a}^\top\| = \mathbf{a}^\top \mathbf{a} = \text{Tr}(\mathbf{a}\mathbf{a}^\top)$ . And hence

$$\begin{aligned} \sum_e \mathbf{w}(e) \left\| \mathbf{L}^{+/-} \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/-} \right\| &= \sum_e \mathbf{w}(e) \text{Tr} \left( \mathbf{L}^{+/-} \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/-} \right) \\ &= \text{Tr} \left( \mathbf{L}^{+/-} \left( \sum_e \mathbf{w}(e) \mathbf{b}_e \mathbf{b}_e^\top \right) \mathbf{L}^{+/-} \right) \\ &= \text{Tr}(\mathbf{\Pi}_L) = n - 1. \end{aligned}$$

Thus with our choice of  $\alpha$ ,

$$\mathbb{E} \left[ |\tilde{E}| \right] \leq 40\epsilon^{-2} \log(n/\delta)n.$$

With a scalar Chernoff bound, can show that  $|\tilde{E}| \leq O(\epsilon^{-2} \log(n/\delta)n)$  with probability at least  $1 - \delta/2$ . Thus by a union bound, the this condition and Equation (9.3) are both satisfied with probability at least  $1 - \delta$ .  $\square$

**Remark 9.3.9.** Note that

$$\|\Phi(\mathbf{w}(e)\mathbf{b}_e\mathbf{b}_e^\top)\| = \mathbf{w}(e) \left\| \mathbf{L}^{+/-} \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/-} \right\| \leq \mathbf{w}(e) \left\| \mathbf{L}^{+/-} \mathbf{b}_e \right\|_2^2.$$

Recall that in Chapter 5, we saw that the effective resistance between vertex  $v$  and vertex  $u$  is given by  $\left\| \mathbf{L}^{+/-}(\mathbf{e}_u - \mathbf{e}_v) \right\|_2^2$ , and for an edge  $e$  connecting vertex  $u$  and  $v$ , we have  $\mathbf{b}_e = \mathbf{e}_u - \mathbf{e}_v$ . That means the norm of the “baby Laplacian”  $\mathbf{w}(e)\mathbf{b}_e\mathbf{b}_e^\top$  of a single edge with weight  $\mathbf{w}(e)$  is exactly  $\mathbf{w}(e)$  times the effective resistance between the two endpoints of the edge.

We haven't shown how to compute the sampling probabilities efficiently, so right now, it isn't clear whether we can efficiently find  $\tilde{G}$ . It turns out that if we have access to a fast algorithm for solving Laplacian linear equations, then we can find sufficiently good approximations to the effective resistances quickly, and use these to compute  $\tilde{G}$ . An algorithm for this is described in [SS11].

# Chapter 10

## Solving Laplacian Linear Equations

### 10.1 Solving Linear Equations Approximately

Given a Laplacian  $\mathbf{L}$  of a connected graph and a demand vector  $\mathbf{d} \perp \mathbf{1}$ , we want to find  $\mathbf{x}^*$  solving the linear equation  $\mathbf{L}\mathbf{x}^* = \mathbf{d}$ . We are going to focus on fast algorithms for finding approximate (but highly accurate) solutions.

This means we need a notion of an approximate solution. Since our definition is not special to Laplacians, we state it more generally for positive semi-definite matrices.

**Definition 10.1.1.** Given PSD matrix  $\mathbf{M}$  and  $\mathbf{d} \in \ker(\mathbf{M})^\perp$ , let  $\mathbf{M}\mathbf{x}^* = \mathbf{d}$ . We say that  $\tilde{\mathbf{x}}$  is an  $\epsilon$ -approximate solution to the linear equation  $\mathbf{M}\mathbf{x} = \mathbf{d}$  if

$$\|\tilde{\mathbf{x}} - \mathbf{x}^*\|_{\mathbf{M}}^2 \leq \epsilon \|\mathbf{x}^*\|_{\mathbf{M}}^2.$$

**Remark 10.1.2.** The requirement  $\mathbf{d} \in \ker(\mathbf{M})^\perp$  can be removed, but this is not important for us.

**Theorem 10.1.3** (Spielman and Teng (2004) [ST04]). *Given a Laplacian  $\mathbf{L}$  of a weighted undirected graph  $G = (V, E, \mathbf{w})$  with  $|E| = m$  and  $|V| = n$  and a demand vector  $\mathbf{d} \in \mathbb{R}^V$ , we can find  $\tilde{\mathbf{x}}$  that is an  $\epsilon$ -approximate solution to  $\mathbf{L}\mathbf{x} = \mathbf{d}$ , using an algorithm that takes time  $O(m \log^c n \log(1/\epsilon))$  for some fixed constant  $c$  and succeeds with probability  $1 - 1/n^{10}$ .*

In the original algorithm of Spielman and Teng, the exponent on the log in the running time was  $c \approx 70$ .

Today, we are going to see a simpler algorithm. But first, we'll look at one of the key tools behind all algorithms for solving Laplacian linear equations quickly.

## 10.2 Preconditioning and Approximate Gaussian Elimination

Recall our definition of two positive semi-definite matrices being approximately equal.

**Definition 10.2.1** (Spectral approximation). Given  $\mathbf{A}, \mathbf{B} \in S_+^n$ , we say that

$$\mathbf{A} \approx_K \mathbf{B} \text{ if and only if } \frac{1}{1+K}\mathbf{A} \preceq \mathbf{B} \preceq (1+K)\mathbf{A}.$$

Suppose we have a positive definite matrix  $\mathbf{M} \in S_{++}^n$  and want to solve a linear equation  $\mathbf{M}\mathbf{x} = \mathbf{d}$ . We can do this using gradient descent or accelerated gradient descent, as we covered in Graded Homework 1. But if we have access to an easy-to-invert matrix that happens to also be a good spectral approximation of  $\mathbf{M}$ , then we can use this to speed up the (accelerated) gradient descent algorithm. An example of this would be that we have a factorization  $\mathcal{L}\mathcal{L}^\top \approx_K \mathbf{M}$ , where  $\mathcal{L}$  is lower triangular and sparse, which means we can invert it quickly.

The following lemma, which you will prove in Problem Set 6, makes this preconditioning precise.

**Lemma 10.2.2.** *Given a matrix  $\mathbf{M} \in S_{++}^n$ , a vector  $\mathbf{d}$  and a decomposition  $\mathbf{M} \approx_K \mathcal{L}\mathcal{L}^\top$ , we can find  $\tilde{\mathbf{x}}$  that  $\epsilon$ -approximately solves  $\mathbf{M}\mathbf{x} = \mathbf{d}$ , using  $O((1+K)\log(K/\epsilon)(T_{\text{matvec}} + T_{\text{sol}} + n))$  time.*

- $T_{\text{matvec}}$  denotes the time required to compute  $\mathbf{M}\mathbf{z}$  given a vector  $\mathbf{z}$ , i.e. a “matrix-vector multiplication”.
- $T_{\text{sol}}$  denotes the time required to compute  $\mathcal{L}^{-1}\mathbf{z}$  or  $(\mathcal{L}^\top)^{-1}\mathbf{z}$  given a vector  $\mathbf{z}$ .

**Dealing with pseudo-inverses.** When our matrices have a null space, preconditioning becomes slightly more complicated, but as long as it is easy to project to the complement of the null space, there’s no real issue. The following describes precisely what we need (but you can ignore the null-space issue when first reading these notes without losing anything significant).

**Lemma 10.2.3.** *Given a matrix  $\mathbf{M} \in S_+^n$ , a vector  $\mathbf{d} \in \ker(\mathbf{M})^\perp$  and a decomposition  $\mathbf{M} \approx_K \mathcal{L}\mathcal{D}\mathcal{L}^\top$ , where  $\mathcal{L}$  is invertible, we can find  $\tilde{\mathbf{x}}$  that  $\epsilon$ -approximately solves  $\mathbf{M}\mathbf{x} = \mathbf{d}$ , using  $O((1+K)\log(K/\epsilon)(T_{\text{matvec}} + T_{\text{sol}} + T_{\text{proj}} + n))$  time.*

- $T_{\text{matvec}}$  denotes the time required to compute  $\mathbf{M}\mathbf{z}$  given a vector  $\mathbf{z}$ , i.e. a “matrix-vector multiplication”.
- $T_{\text{sol}}$  denotes the time required to compute  $\mathcal{L}^{-1}\mathbf{z}$  and  $(\mathcal{L}^\top)^{-1}\mathbf{z}$  and  $\mathcal{D}^+\mathbf{z}$  given a vector  $\mathbf{z}$ .

- $T_{\text{proj}}$  denotes the time required to compute  $\Pi_M \mathbf{z}$  given a vector  $\mathbf{z}$ .

**Theorem 10.2.4** (Kyng and Sachdeva (2015) [KS16]). *Given a Laplacian  $\mathbf{L}$  of a weighted undirected graph  $G = (V, E, \mathbf{w})$  with  $|E| = M$  and  $|V| = n$ , we can find a decomposition  $\mathcal{L}\mathcal{L}^\top \approx_{0.5} \mathbf{L}$ , such that  $\mathcal{L}$  has number of non-zeroes  $\text{nnz}(\mathcal{L}) = O(m \log^3 n)$ , with probability at least  $1 - 3/n^5$ . in time  $O(m \log^3 n)$ .*

We can combine Theorem 10.2.4 with Lemma 10.2.3 to get a fast algorithm for solving Laplacian linear equations.

**Corollary 10.2.5.** *Given a Laplacian  $\mathbf{L}$  of a weighted undirected graph  $G = (V, E, \mathbf{w})$  with  $|E| = m$  and  $|V| = n$  and a demand vector  $\mathbf{d} \in \mathbb{R}^V$ , we can find  $\tilde{\mathbf{x}}$  that is an  $\epsilon$ -approximate solution to  $\mathbf{L}\mathbf{x} = \mathbf{d}$ , using an algorithm that takes time  $O(m \log^3 n \log(1/\epsilon))$  and succeeds with probability  $1 - 1/n^{10}$ .*

*Proof sketch.* First we need to get a factorization that confirms to Lemma 10.2.3. The decomposition  $\mathcal{L}\mathcal{L}^\top$  provided by Theorem 10.2.4 can be rewritten as  $\mathcal{L}\mathcal{L}^\top = \tilde{\mathcal{L}}\mathcal{D}(\tilde{\mathcal{L}})^\top$  where  $\tilde{\mathcal{L}}$  is equal to  $\mathcal{L}$  except  $\mathcal{L}(n, n) = 1$  and we let  $\mathcal{D}$  be the identity matrix, except  $\mathcal{D}(n, n) = 0$ . This ensures  $\mathcal{D}^+ = \mathcal{D}$  and that  $\tilde{\mathcal{L}}$  is invertible and lower triangular with  $O(m \log^3 n)$  non-zeros. We note that the inverse of an invertible lower or upper triangular matrix with  $N$  non-zeros can be applied in time  $O(N)$  given an adjacency list representation of the matrix. Finally, as  $\ker(\mathcal{L}\mathcal{L}^\top) = \text{span}\{\mathbf{1}\}$ , we have  $\Pi_{\tilde{\mathcal{L}}\mathcal{D}(\tilde{\mathcal{L}})^\top} = \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$ , and this projection matrix can be applied in  $O(n)$  time. Altogether, this means that  $T_{\text{matvec}} + T_{\text{sol}} + T_{\text{proj}} = O(n)$ , which suffices to complete the proof.  $\square$

## 10.3 Approximate Gaussian Elimination Algorithm

Recall *Gaussian Elimination / Cholesky decomposition* of a graph Laplacian  $\mathbf{L}$ . We will use  $\mathbf{A}(:, i)$  to denote the the  $i$ th column of a matrix  $\mathbf{A}$ . We can write the algorithm as

---

**Algorithm 1:** Gaussian Elimination / Cholesky Decomposition

---

**Input:** Graph Laplacian  $\mathbf{L}$

**Output:** Lower triangular  $\mathcal{L}$  s.t.  $\mathcal{L}\mathcal{L}^\top = \mathbf{L}$

```

1 Let  $\mathbf{S}_0 = \mathbf{L}$ ;
2 for  $i = 1$  to  $i = n - 1$  do
3    $\mathbf{l}_i = \frac{1}{\sqrt{s_{i-1}(i,i)}} \mathbf{S}_{i-1}(:, i);$ 
4    $\mathbf{S}_i = \mathbf{S}_{i-1} - \mathbf{l}_i \mathbf{l}_i^\top.$ 
5    $\mathbf{l}_n = \mathbf{0}_{n \times 1};$ 
6 return  $\mathcal{L} = [\mathbf{l}_1 \cdots \mathbf{l}_n];$ 

```

---

We want to introduce some notation that will help us describe and analyze a faster version of Gaussian elimination – one that uses sampling to create a sparse approximation of the decomposition.

Consider a Laplacian  $\mathbf{S}$  of a graph  $H$  and a vertex  $v$  of  $H$ . We define  $\text{STAR}(v, \mathbf{S})$  to be the Laplacian of the subgraph of  $H$  consisting of edges incident on  $v$ . We define

$$\text{CLIQUE}(v, \mathbf{S}) = \text{STAR}(v, \mathbf{S}) - \frac{1}{\mathbf{S}_{(v,v)}} \mathbf{S}(:,v) \mathbf{S}(:,v)^\top$$

For example, suppose

$$\mathbf{L} = \begin{pmatrix} W & -\mathbf{a}^\top \\ -\mathbf{a} & \text{diag}(\mathbf{a}) + \mathbf{L}_{-1} \end{pmatrix}$$

Then

$$\text{STAR}(1, \mathbf{L}) = \begin{pmatrix} W & -\mathbf{a}^\top \\ -\mathbf{a} & \text{diag}(\mathbf{a}) \end{pmatrix} \text{ and } \text{CLIQUE}(1, \mathbf{L}) = \begin{pmatrix} 0 & \mathbf{0} \\ \mathbf{0} & \text{diag}(\mathbf{a}) - \frac{1}{W} \mathbf{a} \mathbf{a}^\top \end{pmatrix}$$

which is illustrated in Figure 10.1.

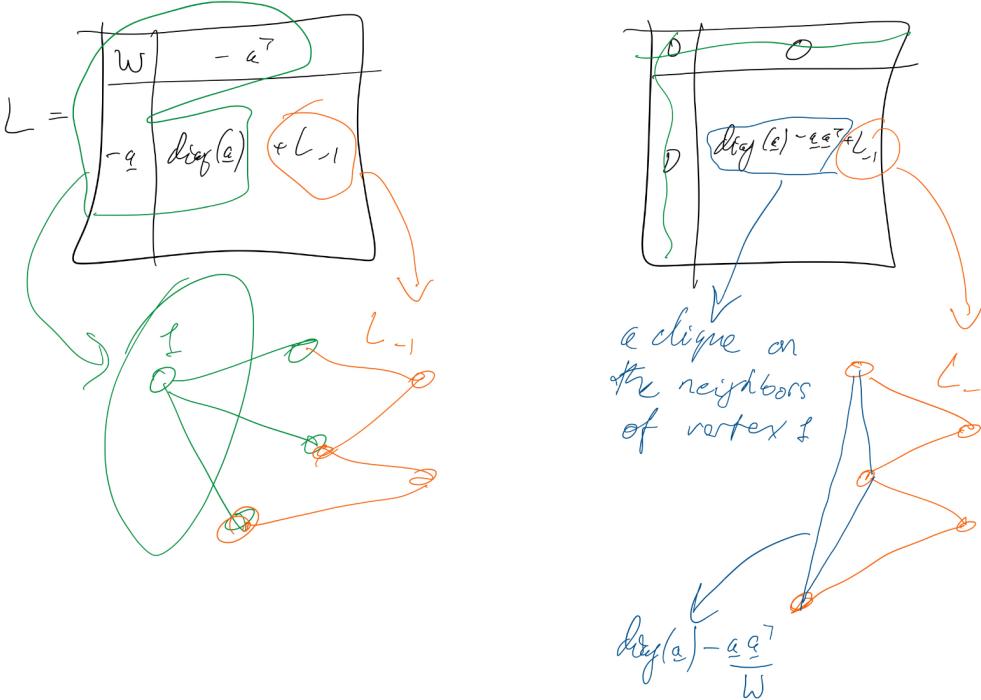


Figure 10.1: Gaussian Elimination:  $\text{CLIQUE}(1, \mathbf{L}) = \text{STAR}(1, \mathbf{L}) - \frac{1}{\mathbf{L}(1,1)} \mathbf{L}(:,1) \mathbf{L}(:,1)^\top$ .

In Chapter 6, we proved that  $\text{CLIQUE}(v, \mathbf{S})$  is a graph Laplacian – it follows from the proof of Claim 6.1.1 in that chapter. Thus we have that following.

**Claim 10.3.1.** *If  $\mathbf{S}$  is the Laplacian of a connected graph, then  $\text{CLIQUE}(v, \mathbf{S})$  is a graph Laplacian.*

Note that in Algorithm 1, we have  $\mathbf{l}_i \mathbf{l}_i^\top = \text{STAR}(v_i, \mathbf{S}_{i-1}) - \text{CLIQUE}(v_i, \mathbf{S}_{i-1})$ . The update rule can be rewritten as

$$\mathbf{S}_i = \mathbf{S}_{i-1} - \text{STAR}(v_i, \mathbf{S}_{i-1}) + \text{CLIQUE}(v_i, \mathbf{S}_{i-1}),$$

This also provides way to understand why Gaussian Elimination is slow in some cases. At each step, one vertex is eliminated, but a clique is added to the subgraph on the remaining vertices, making the graph denser. And at the  $i$ th step, computing  $\text{STAR}(v_i, \mathbf{S}_{i-1})$  takes around  $\deg(v_i)$  time, but computing  $\text{CLIQUE}(v_i, \mathbf{S}_{i-1})$  requires around  $\deg(v_i)^2$  time. In order to speed up Gaussian Elimination, the algorithmic idea of [KS16] is to plug in a sparser approximate of the intended clique instead of the entire one.

The following procedure  $\text{CLIQUESAMPLE}(v, \mathbf{S})$  produces a sparse approximation of  $\text{CLIQUE}(v, \mathbf{S})$ . Let  $V$  be the vertex set of the graph associated with  $\mathbf{S}$  and  $E$  the edge set. We define  $\mathbf{b}_{i,j} \in \mathbb{R}^V$  to be the vector with

$$\mathbf{b}_{i,j}(i) = 1 \text{ and } \mathbf{b}_{i,j}(j) = -1 \text{ and } \mathbf{b}_{i,j}(k) = 0 \text{ for } k \neq i, j.$$

Given weights  $\mathbf{w} \in \mathbb{R}^E$  and a vertex  $v \in V$ , we let

$$\mathbf{w}_v = \sum_{(u,v) \in E} \mathbf{w}(u, v).$$

---

**Algorithm 2:**  $\text{CLIQUESAMPLE}(v, \mathbf{S})$ 


---

**Input:** Graph Laplacian  $\mathbf{S} \in \mathbb{R}^{V \times V}$ , of a graph with edge weights  $\mathbf{w}$ , and vertex  $v \in V$

**Output:**  $\mathbf{Y}_v \in \mathbb{R}^{V \times V}$  sparse approximation of  $\text{CLIQUE}(v, \mathbf{S})$

- 1  $\mathbf{Y}_v \leftarrow \mathbf{0}_{n \times n};$
  - 2 **foreach** Multiedge  $e = (v, i)$  from  $v$  to a neighbor  $i$  **do**
  - 3     Randomly pick a neighbor  $j$  of  $v$  with probability  $\frac{\mathbf{w}(j, v)}{\mathbf{w}_v};$
  - 4     If  $i \neq j$ , let  $\mathbf{Y}_v \leftarrow \mathbf{Y}_v + \frac{\mathbf{w}(i, v)\mathbf{w}(j, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} \mathbf{b}_{i,j} \mathbf{b}_{i,j}^\top;$
  - 5 **return**  $\mathbf{Y}_v;$
- 

**Remark 10.3.2.** We can implement each sampling of a neighbor  $j$  in  $O(1)$  time using a classical algorithm known as Walker's method (also known as the Alias method or Vose's method). This algorithm requires an additional  $O(\deg_S(v))$  time to initialize a data structure used for sampling. Overall, this means the total time for  $O(\deg_S(v))$  samples is still  $O(\deg_S(v))$ .

**Lemma 10.3.3.**  $\mathbb{E}[\mathbf{Y}_v] = \text{CLIQUE}(v, \mathbf{S}).$

*Proof.* Let  $\mathbf{C} = \text{CLIQUE}(v, \mathbf{S})$ . Observe that both  $\mathbb{E}[\mathbf{Y}_v]$  and  $\mathbf{C}$  are Laplacians. Thus it suffices to verify  $\mathbb{E}_{\mathbf{Y}_v(i,j)} = \mathbf{C}(i, j)$  for  $i \neq j$ .

$$\mathbf{C}(i, j) = -\frac{\mathbf{w}(i, v)\mathbf{w}(j, v)}{\mathbf{w}_v},$$

$$\mathbb{E}_{\mathbf{Y}_v(i,j)} = -\frac{\mathbf{w}(i, v)\mathbf{w}(j, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} \left( \frac{\mathbf{w}(j, v)}{\mathbf{w}_v} + \frac{\mathbf{w}(i, v)}{\mathbf{w}_v} \right) = -\frac{\mathbf{w}(i, v)\mathbf{w}(j, v)}{\mathbf{w}_v} = \mathbf{C}(i, j).$$

□

**Remark 10.3.4.** Lemma 10.3.3 shows that  $\text{CLIQUESAMPLE}(v, \mathbf{L})$  produces the original  $\text{CLIQUE}(v, \mathbf{L})$  in expectation.

Now, we define *Approximate Gaussian Elimination*.

---

**Algorithm 3:** Approximate Gaussian Elimination / Cholesky Decomposition

---

**Input:** Graph Laplacian  $\mathbf{L}$

**Output:** Lower triangular<sup>a</sup>  $\mathcal{L}$  as given in Theorem 10.2.4

- 1 Let  $\mathbf{S}_0 = \mathbf{L}$ ;
  - 2 Generate a random permutation  $\pi$  on  $[n]$ ;
  - 3 **for**  $i = 1$  to  $i = n - 1$  **do**
  - 4      $\mathbf{l}_i = \frac{1}{\sqrt{\mathbf{S}_{i-1}(\pi(i), \pi(i))}} \mathbf{S}_{i-1}(:, \pi(i))$ ;
  - 5      $\mathbf{S}_i = \mathbf{S}_{i-1} - \text{STAR}(\pi(i), \mathbf{S}_{i-1}) + \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1})$
  - 6      $\mathbf{l}_n = \mathbf{0}_{n \times 1}$ ;
  - 7 **return**  $\mathcal{L} = [\mathbf{l}_1 \cdots \mathbf{l}_n]$  and  $\pi$ ;
- 

<sup>a</sup> $\mathcal{L}$  is not actually lower triangular. However, if we let  $\mathbf{P}_\pi$  be the permutation matrix corresponding to  $\pi$ , then  $\mathbf{P}_\pi \mathcal{L}$  is lower triangular. Knowing the ordering that achieves this is enough to let us implement forward and backward substitution for solving linear equations in  $\mathcal{L}$  and  $\mathcal{L}^\top$ .

Note that if we replace  $\text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1})$  by  $\text{CLIQUE}(\pi(i), \mathbf{S}_{i-1})$  at each step, then we can recover Gaussian Elimination, but with a random elimination order.

## 10.4 Analyzing Approximate Gaussian Elimination

In this Section, we're going to analyze Approximate Gaussian Elimination, and see why it works.

Ultimately, the main challenge in proving Theorem 10.2.4 will be to prove for the output  $\mathcal{L}$  of Algorithm 3 that with high probability

$$0.5\mathbf{L} \preceq \mathcal{L}\mathcal{L}^\top \preceq 1.5\mathbf{L}. \quad (10.1)$$

We can reduce this to proving that with high probability

$$\left\| \mathbf{L}^{+/-} (\mathcal{L}\mathcal{L}^\top - \mathbf{L}) \mathbf{L}^{+/-} \right\| \leq 0.5 \quad (10.2)$$

Ultimately, the proof is going to have a lot in common with our proof of Matrix Bernstein in Chapter 9. Overall, the lesson there was that when we have a sum of independent, zero-mean random matrices, we can show that the sum is likely to have small spectral norm if the spectral norm of each random matrix is small, and the matrix-valued variance is also small.

Thus, to replicate the proof, we need control over

1. The *sample norms*.

## 2. The *sample variance*.

But, there is seemingly another major obstacle: We are trying to analyze a process where the samples are far from independent. Each time we sample edges, we add new edges to the remaining graph, which we will later sample again. This creates a lot of dependencies between the samples, which we have to handle.

However, it turns out that independence is more than what is needed to prove concentration. Instead, it suffices to have a sequence of random variables such that each is mean-zero in expectation, conditional on the previous ones. This is called a martingale difference sequence. We'll now learn about those.

### 10.4.1 Normalization, a.k.a. Isotropic Position

Since our analysis requires frequently measuring matrices after right and left-multiplication by  $\mathbf{L}^{+/-}$ , we reintroduce the “normalizing map”  $\Phi : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$  defined by

$$\Phi(\mathbf{A}) = \mathbf{L}^{+/-} \mathbf{A} \mathbf{L}^{+/-}.$$

We previously saw this in Chapter 9.

### 10.4.2 Martingales

A scalar martingale is a sequence of random variables  $Z_0, \dots, Z_k$ , such that

$$\mathbb{E}[Z_i | Z_0, \dots, Z_{i-1}] = Z_{i-1}. \quad (10.3)$$

That is, conditional on the outcome of all the previous random variables, the expectation of  $Z_i$  equals  $Z_{i-1}$ . If we unravel the sequence of conditional expectations, we get that *without conditioning*  $\mathbb{E}[Z_k] = \mathbb{E}[Z_0]$ .

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

We can also think of a martingale in terms of the sequence of changes in the  $Z_i$  variables. Let  $X_i = Z_i - Z_{i-1}$ . The sequence of  $X_i$ s is called a martingale difference sequence. We can now state the martingale condition as

$$\mathbb{E}[X_i | Z_0, \dots, Z_{i-1}] = 0.$$

And because  $Z_0$  and  $X_1, \dots, X_{i-1}$  completely determine  $Z_1, \dots, Z_{i-1}$ , we could also write the martingale condition equivalently as

$$\mathbb{E}[X_i | Z_0, X_1, \dots, X_{i-1}] = 0.$$

Crucially, we can write

$$Z_k = Z_0 + \sum_{i=1}^k Z_i - Z_{i-1} = Z_0 + \sum_{i=1}^k X_i$$

and when we are trying to prove concentration, the martingale difference property of the  $X_i$ 's is often “as good as” independence, meaning that  $\sum_{i=1}^k X_i$  concentrates similarly to a sum of independent random variables.

**Matrix-valued martingales.** We can also define matrix-valued martingales. In this case, we replace the martingalue condition of Equation (10.3), with the condition that the whole matrix stays the same in expectation. For example, we could have a sequence of random matrices  $\mathbf{Z}_0, \dots, \mathbf{Z}_k \in \mathbb{R}^{n \times n}$ , such that

$$\mathbb{E} [\mathbf{Z}_i | \mathbf{Z}_0, \dots, \mathbf{Z}_{i-1}] = \mathbf{Z}_{i-1}. \quad (10.4)$$

**Lemma 10.4.1.** Let  $\mathbf{L}_i = \mathbf{S}_i + \sum_{j=1}^i \mathbf{l}_j \mathbf{l}_j^\top$  for  $i = 1, \dots, n$  and  $\mathbf{L}_0 = \mathbf{S}_0 = \mathbf{L}$ . Then

$$\mathbb{E} [\mathbf{L}_i | \text{all random variables before CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1})] = \mathbf{L}_{i-1}.$$

*Proof.* Let's only consider  $i = 1$  here as other cases are similar.

$$\mathbf{L}_0 = \mathbf{L} = \mathbf{l}_1 \mathbf{l}_1^\top + \text{CLIQUE}(v, \mathbf{L}) + \mathbf{L}_{-1}$$

$$\begin{aligned} \mathbf{L}_1 &= \mathbf{l}_1 \mathbf{l}_1^\top + \text{CLIQUESAMPLE}(v, \mathbf{L}) + \mathbf{L}_{-1} \\ \mathbb{E} [\mathbf{L}_1 | \pi(1)] &= \mathbf{l}_1 \mathbf{l}_1^\top + \mathbb{E} [\text{CLIQUESAMPLE}(v, \mathbf{L}) | \pi(1)] + \mathbf{L}_{-1} \\ &= \mathbf{l}_1 \mathbf{l}_1^\top + \text{CLIQUE}(v, \mathbf{L}) + \mathbf{L}_{-1} \\ &= \mathbf{L}_0 \end{aligned}$$

where we used Lemma 10.3.3 to get  $\mathbb{E} [\text{CLIQUESAMPLE}(v, \mathbf{L}) | \pi(1)] = \text{CLIQUE}(v, \mathbf{L})$ .  $\square$

**Remark 10.4.2.**  $\sum_{j=1}^i \mathbf{l}_j \mathbf{l}_j^\top$  can be treated as what has already been eliminated by (Approximate) Gaussian Elimination, while  $\mathbf{S}_i$  is what still left or going to be eliminated. In Approximate Gaussian Elimination,  $\mathbf{L}_n = \sum_{i=1}^n \mathbf{l}_i \mathbf{l}_i^\top$  and our goal is to show that  $\mathbf{L}_n \approx_K \mathbf{L}$ . Note that  $\mathbf{L}_i$  is always equal to the original Laplacian  $\mathbf{L}$  for all  $i$  in Gaussian Elimination. Lemma 10.4.1 demonstrates that  $\mathbf{L}_0, \mathbf{L}_1, \dots, \mathbf{L}_n$  forms a matrix martingale.

Ultimately, our plan is to use this matrix martingale structure to show that “ $\mathbf{L}_n$  is concentrated around  $\mathbf{L}$ ” in some appropriate sense. More precisely, the spectral approximation we would like to show can be established by showing that “ $\Phi(\mathbf{L}_n)$  is concentrated around  $\Phi(\mathbf{L})$ ”

### 10.4.3 Martingale Difference Sequence as Edge-Samples

We start by taking a slightly different view of the observations we used to prove Lemma 10.4.1. Recall that  $\mathbf{L}_i = \mathbf{S}_i + \sum_{j=1}^i \mathbf{l}_j \mathbf{l}_j^\top$ , and  $\mathbf{L}_{i-1} = \mathbf{S}_{i-1} + \sum_{j=1}^{i-1} \mathbf{l}_j \mathbf{l}_j^\top$  and

$$\mathbf{S}_i = \mathbf{S}_{i-1} - \text{STAR}(\pi(i), \mathbf{S}_{i-1}) + \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}).$$

Putting these together, we get

$$\begin{aligned} \mathbf{L}_i - \mathbf{L}_{i-1} &= \mathbf{l}_i \mathbf{l}_i^\top + \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) - \text{STAR}(\pi(i), \mathbf{S}_{i-1}) \\ &= \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) - \text{CLIQUE}(\pi(i), \mathbf{S}_{i-1}) \\ &= \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) - \mathbb{E} [\text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) \mid \text{preceding samples}] \end{aligned} \tag{10.5}$$

by Lemma 10.3.3.

In particular, recall that by Lemma 10.3.3, conditional on the randomness before the call to  $\text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1})$ , we have

$$\mathbb{E} [\text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) \mid \text{preceding samples}] = \text{CLIQUE}(\pi(i), \mathbf{S}_{i-1})$$

Adopting the notation of Lemma 10.3.3 we write

$$\mathbf{Y}_{\pi(i)} = \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1})$$

and we further introduce notation each multi-edge sample for  $e \in \text{STAR}(\pi(i), \mathbf{S}_{i-1})$ , as  $\mathbf{Y}_{\pi(i),e}$ , denoting the random edge Laplacian sampled when the algorithm is processing multi-edge  $e$ . Thus, conditional on preceding samples, we have

$$\mathbf{Y}_{\pi(i)} = \sum_{e \in \text{STAR}(\pi(i), \mathbf{S}_{i-1})} \mathbf{Y}_{\pi(i),e} \tag{10.6}$$

Note that even the number of multi-edges in  $\text{STAR}(\pi(i), \mathbf{S}_{i-1})$  depends on the preceding samples. We also want to associate zero-mean variables with each edge. Conditional on preceding samples, we also define

$$\mathbf{X}_{i,e} = \Phi(\mathbf{Y}_{\pi(i),e} - \mathbb{E}[\mathbf{Y}_{\pi(i),e}]) \text{ and } \mathbf{X}_i = \sum_{e \in \text{STAR}(\pi(i), \mathbf{S}_{i-1})} \mathbf{X}_{i,e}$$

and combining this with Equations (10.5) and (10.6)

$$\mathbf{X}_i = \Phi(\mathbf{Y}_{\pi(i)} - \mathbb{E}[\mathbf{Y}_{\pi(i)}]) = \Phi(\mathbf{L}_i - \mathbf{L}_{i-1})$$

Altogether, we can write

$$\Phi(\mathbf{L}_n - \mathbf{L}) = \sum_{i=1}^n \Phi(\mathbf{L}_i - \mathbf{L}_{i-1}) = \sum_{i=1}^n \mathbf{X}_i = \sum_{i=1}^n \sum_{e \in \text{STAR}(\pi(i), \mathbf{S}_{i-1})} \mathbf{X}_{i,e}$$

Note that the  $\mathbf{X}_{i,e}$  variables form a martingale difference sequence, because the linearity of  $\Phi$  ensures they are zero-mean conditional on preceding randomness.

#### 10.4.4 Stopped Martingales

Unfortunately, directly analyzing the concentration properties of the  $\mathbf{L}_i$  martingale that we just introduced turns out to be difficult. The reason is that we're trying to prove some very delicate multiplicative error guarantees. And, if we analyze  $\mathbf{L}_i$ , we find that the multiplicative error is not easy to control, *after it's already gotten big*. But that's not really what we care about anyway: We want to say it never gets big in the first place, with high probability. So we need to introduce another martingale, that lets us ignore the bad case when the error has already gotten too big. At the same time, we also need to make sure that statements about our new martingale can help us prove guarantees about  $\mathbf{L}_i$ . Fortunately, we can achieve both at once. The technique we use is related to the much broader topic of martingale *stopping times*, which we only scratch the surface of here. We're also going to be quite informal about it, in the interest of brevity. Lecture notes by Tropp [Tro19] give a more formal introduction for those who are interested.

We define the stopped martingale sequence  $\tilde{\mathbf{L}}_i$  by

$$\tilde{\mathbf{L}}_i = \begin{cases} \mathbf{L}_i & \text{if for all } j < i \text{ we have } \mathbf{L}_j \preceq 1.5\mathbf{L} \\ \mathbf{L}_{j^*} & \text{for } j^* \text{ being the least } j \text{ such that } \mathbf{L}_j \not\preceq 1.5\mathbf{L} \end{cases} \quad (10.7)$$

Figure 10.2 shows the  $\tilde{\mathbf{L}}_i$  martingale getting stuck at the first time  $\mathbf{L}_{j^*} \not\preceq 1.5\mathbf{L}$ .

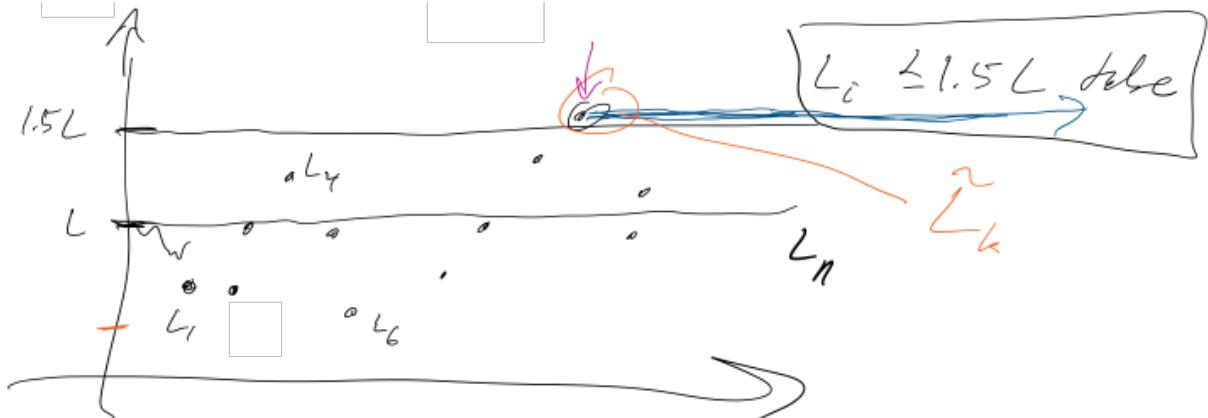


Figure 10.2: Gaussian Elimination : CLIQUE(1,  $\mathbf{L}$ ) = STAR(1,  $\mathbf{L}$ ) -  $\frac{1}{\mathbf{L}(1,1)}\mathbf{L}(:,1)\mathbf{L}(:,1)^\top$ .

We state the following without proof:

**Claim 10.4.3.**

1. The sequence  $\{\tilde{\mathbf{L}}_i\}$  for  $i = 0, \dots, n$  is a martingale.
2.  $\|\mathbf{L}^{+1/2}(\tilde{\mathbf{L}}_i - \mathbf{L})\mathbf{L}^{+1/2}\| \leq 0.5$  implies  $\|\mathbf{L}^{+1/2}(\mathbf{L}_i - \mathbf{L})\mathbf{L}^{+1/2}\| \leq 0.5$

The martingale property also implies that the unconditional expectation satisfies  $\mathbb{E} [\tilde{\mathbf{L}}_n] = \mathbf{L}$ . The proof of the claim is easy to sketch: For Part 1, each difference is zero-mean if the condition has not been violated, and is identically zero (and hence zero-mean) if it has been violated. For Part 2, if the martingale  $\{\tilde{\mathbf{L}}_i\}$  has stopped, then  $\|\mathbf{L}^{+/-}(\tilde{\mathbf{L}}_i - \mathbf{L})\mathbf{L}^{+/-}\| \leq 0.5$  is false, and the implication is vacuously true. If, on the other hand, the martingale has not stopped, the quantities are equal, because  $\tilde{\mathbf{L}}_i = \mathbf{L}_i$ , and again it's easy to see the implication holds.

Thus, ultimately, our strategy is going to be to show that  $\|\mathbf{L}^{+/-}(\tilde{\mathbf{L}}_i - \mathbf{L})\mathbf{L}^{+/-}\| \leq 0.5$  with high probability. Expressed using the normalizing map  $\Phi(\cdot)$ , our goal is to show that with high probability

$$\|\Phi(\tilde{\mathbf{L}}_n - \mathbf{L})\| \leq 0.5.$$

**Stopped martingale difference sequence.** In order to prove the spectral norm bound, we want to express the  $\{\tilde{\mathbf{L}}_i\}$  martingale in terms of a sequence of martingale differences. To this end, we define  $\tilde{\mathbf{X}}_i = \Phi(\tilde{\mathbf{L}}_i - \tilde{\mathbf{L}}_{i-1})$ . This ensures that

$$\tilde{\mathbf{X}}_i = \begin{cases} \mathbf{X}_i & \text{if for all } j < i \text{ we have } \mathbf{L}_i \preceq 1.5\mathbf{L} \\ \mathbf{0} & \text{otherwise} \end{cases} \quad (10.8)$$

Whenever the modified martingale  $\tilde{\mathbf{X}}_i$  has not yet stopped, we also introduce individual modified edge samples  $\tilde{\mathbf{X}}_{i,e} = \mathbf{X}_{i,e}$ . If the martingale *has* stopped, i.e.  $\tilde{\mathbf{X}}_i = \mathbf{0}$ , then we can take these edge samples  $\tilde{\mathbf{X}}_{i,e}$  to be zero. We can now write

$$\Phi(\tilde{\mathbf{L}}_n - \mathbf{L}) = \sum_{i=1}^n \Phi(\tilde{\mathbf{L}}_i - \tilde{\mathbf{L}}_{i-1}) = \sum_{i=1}^n \tilde{\mathbf{X}}_i = \sum_{i=1}^n \sum_{e \in \text{STAR}(\pi(i), \mathbf{S}_{i-1})} \tilde{\mathbf{X}}_{i,e}.$$

Thus, we can see that Equation (10.2) is implied by

$$\left\| \sum_{i=1}^n \tilde{\mathbf{X}}_i \right\| \leq 0.5. \quad (10.9)$$

#### 10.4.5 Sample Norm Control

In this Subsection, we're going to see that the norms of each multi-edge sample is controlled throughout the algorithm.

**Lemma 10.4.4.** *Given two Laplacians  $\mathbf{L}$  and  $\mathbf{S}$  on the same vertex set.<sup>1</sup> If each multiedge  $e$  of  $\text{STAR}(v, \mathbf{S})$  has bounded norm in the following sense,*

$$\|\mathbf{L}^{+/-} \mathbf{w}_{\mathbf{S}(e)} \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/-}\| \leq R,$$

---

<sup>1</sup>  $\mathbf{L}$  can be regarded as the original Laplacian we care about, while  $\mathbf{S}$  can be regarded as some intermediate Laplacian appearing during Approximate Gaussian Elimination.

then each possible sampled multiedge  $e'$  of CLIQUESAMPLE( $v, \mathbf{S}$ ) also satisfies

$$\left\| \mathbf{L}^{+/-2} \mathbf{w}_{\text{new}}(e') \mathbf{b}_{e'} \mathbf{b}_{e'}^\top \mathbf{L}^{+/-2} \right\| \leq R.$$

*Proof.* Let  $\mathbf{w} = \mathbf{w}_S$  for simplicity. Consider a sampled edge between  $i$  and  $j$  with weight  $\mathbf{w}_{\text{new}}(i, j) = \mathbf{w}(i, v)\mathbf{w}(j, v)/(\mathbf{w}(i, v) + \mathbf{w}(j, v))$ .

$$\begin{aligned} \left\| \mathbf{L}^{+/-2} \mathbf{w}_{\text{new}}(i, j) \mathbf{b}_{ij} \mathbf{b}_{ij}^\top \mathbf{L}^{+/-2} \right\| &= \mathbf{w}_{\text{new}}(i, j) \left\| \mathbf{L}^{+/-2} \mathbf{b}_{ij} \mathbf{b}_{ij}^\top \mathbf{L}^{+/-2} \right\| \\ &= \mathbf{w}_{\text{new}}(i, j) \left\| \mathbf{L}^{+/-2} \mathbf{b}_{ij} \right\|^2 \\ &\leq \mathbf{w}_{\text{new}}(i, j) \left( \left\| \mathbf{L}^{+/-2} \mathbf{b}_{iv} \right\|^2 + \left\| \mathbf{L}^{+/-2} \mathbf{b}_{jv} \right\|^2 \right) \\ &= \frac{\mathbf{w}(j, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} \left\| \mathbf{L}^{+/-2} \mathbf{w}(i, v) \mathbf{b}_{iv} \mathbf{b}_{iv}^\top \mathbf{L}^{+/-2} \right\| + \\ &\quad \frac{\mathbf{w}(i, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} \left\| \mathbf{L}^{+/-2} \mathbf{w}(j, v) \mathbf{b}_{jv} \mathbf{b}_{jv}^\top \mathbf{L}^{+/-2} \right\| \\ &\leq \frac{\mathbf{w}(j, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} R + \frac{\mathbf{w}(i, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} R \\ &= R \end{aligned}$$

The first inequality uses the triangle inequality of effective resistance in  $\mathbf{L}$ , in that effective resistance is a distance as we proved in Chapter 5. The second inequality just uses the conditions of this lemma.  $\square$

**Remark 10.4.5.** Lemma 10.4.4 only requires that each single multiedge has small norm instead of that the sum of all edges between a pair of vertices have small norm. And this lemma tells us, after sampling, each multiedge in the new graph still satisfies the bounded norm condition.

From the Lemma, we can conclude that each edge sample  $\mathbf{Y}_{\pi(i), e}$  satisfies  $\|\Phi(\mathbf{Y}_{\pi(i), e})\| \leq R$  provided the assumptions of the Lemma hold. Let's record this observation as a Lemma.

**Lemma 10.4.6.** *If for all  $e \in \text{STAR}(v, \mathbf{S}_i)$ ,*

$$\|\Phi(\mathbf{w}_{\mathbf{S}_i}(e) \mathbf{b}_e \mathbf{b}_e^\top)\| \leq R.$$

*then all  $e \in \text{STAR}(\pi(i), \mathbf{S}_i)$ ,*

$$\|\Phi(\mathbf{Y}_{\pi(i), e})\| \leq R.$$

**Preprocessing by multi-edge splitting.** In the original graph of Laplacian  $\mathbf{L}$  of graph  $G = (V, E, \mathbf{w})$ , we have for each edge  $\hat{e}$  that

$$\mathbf{w}(\hat{e}) \mathbf{b}_{\hat{e}} \mathbf{b}_{\hat{e}}^\top \preceq \sum_e \mathbf{w}(e) \mathbf{b}_e \mathbf{b}_e^\top = \mathbf{L}$$

This also implies that

$$\left\| \mathbf{L}^{+/-2} \mathbf{w}(\hat{e}) \mathbf{b}_{\hat{e}} \mathbf{b}_{\hat{e}}^\top \mathbf{L}^{+/-2} \right\| \leq 1.$$

Now, that means that if we split every original edge  $e$  of the graph into  $K$  multi-edges  $e_1, \dots, e_K$ , with a fraction  $1/K$  of the weight, we get a new graph  $G' = (V, E', \mathbf{w}')$  such that

**Claim 10.4.7.**

1.  $G'$  and  $G$  have the same graph Laplacian.
2.  $|E'| = K |E|$
3. For every multi-edge in  $G'$

$$\left\| \mathbf{L}^{+/-2} \mathbf{w}'(e) \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/-2} \right\| \leq 1/K.$$

Before we run Approximate Gaussian Elimination, we are going to do this multi-edge splitting to ensure we have control over multi-edge sample norms. Combined with Lemma 10.4.4 immediately establishes the next lemma, because we start off with all multi-edges having bounded norm and only produce multi-edges with bounded norm.

**Lemma 10.4.8.** *When Algorithm 3 is run on the (multi-edge) Laplacian of  $G'$ , arising from splitting edges of  $G$  into  $K$  multi-edges, the every edge sample  $\mathbf{Y}_{\pi(i),e}$  satisfies*

$$\|\Phi(\mathbf{Y}_{\pi(i),e})\| \leq 1/K.$$

As we will see later  $K = 200 \log^2 n$  suffices.

#### 10.4.6 Random Matrix Concentration from Trace Exponentials

Let us recall how matrix-valued variances come into the picture when proving concentration following the strategy from Matrix Bernstein in Chapter 9.

For some matrix-valued random variable  $\mathbf{X} \in S^n$ , we'd like to show  $\Pr[\|\mathbf{X}\| \leq 0.5]$ . Using Markov's inequality, and some observations about matrix exponentials and traces, we saw that for all  $\theta > 0$ ,

$$\Pr[\|\mathbf{X}\| \geq 0.5] \leq \exp(-0.5\theta) (\mathbb{E}[\text{Tr}(\exp(\theta\mathbf{X}))] + \mathbb{E}[\text{Tr}(\exp(-\theta\mathbf{X}))]). \quad (10.10)$$

We then want to bound  $\mathbb{E}[\text{Tr}(\exp(\theta\mathbf{X}))]$  using Lieb's theorem. We can handle  $\mathbb{E}[\text{Tr}(\exp(-\theta\mathbf{X}))]$  similarly.

**Theorem 10.4.9** (Lieb). *Let  $f : S_{++}^n \rightarrow \mathbb{R}$  be a matrix function given by*

$$f(\mathbf{A}) = \text{Tr}(\exp(\mathbf{H} + \log(\mathbf{A})))$$

*for some  $\mathbf{H} \in S^n$ . Then  $-f$  is convex (i.e.  $f$  is concave).*

As observed by Tropp, this is useful for proving matrix concentration statements. Combined with Jensen's inequality, it gives that for a random matrix  $\mathbf{X} \in S^n$  and a fixed  $\mathbf{H} \in S^n$

$$\mathbb{E} [\text{Tr} (\exp (\mathbf{H} + \mathbf{X}))] \leq \text{Tr} (\exp (\mathbf{H} + \log (\mathbb{E} [\exp (\mathbf{X})]))) .$$

The next crucial step was to show that it suffices to obtain an upper bound on the matrix  $\mathbb{E} [\exp (\mathbf{X})]$  w.r.t the Loewner order. Using the following three lemmas, this conclusion is an immediate corollary.

**Lemma 10.4.10.** *If  $\mathbf{A} \preceq \mathbf{B}$ , then  $\text{Tr} (\exp (\mathbf{A})) \leq \text{Tr} (\exp (\mathbf{B}))$ .*

**Lemma 10.4.11.** *If  $0 \prec \mathbf{A} \preceq \mathbf{B}$ , then  $\log (\mathbf{A}) \preceq \log (\mathbf{B})$ .*

**Lemma 10.4.12.**  $\log (\mathbf{I} + \mathbf{A}) \preceq \mathbf{A}$  for  $\mathbf{A} \succ -\mathbf{I}$ .

**Corollary 10.4.13.** *For a random matrix  $\mathbf{X} \in S^n$  and a fixed  $\mathbf{H} \in S^n$ , if  $\mathbb{E} [\exp (\mathbf{X})] \preceq \mathbf{I} + \mathbf{U}$  where  $\mathbf{U} \succ -\mathbf{I}$ , then*

$$\mathbb{E} [\text{Tr} (\exp (\mathbf{H} + \mathbf{X}))] \leq \text{Tr} (\exp (\mathbf{H} + \mathbf{U})) .$$

#### 10.4.7 Mean-Exponential Bounds from Variance Bounds

To use Corollary 10.4.13, we need to construct useful upper bounds on  $\mathbb{E} [\exp (\mathbf{X})]$ . This can be done, starting from the following lemma.

**Lemma 10.4.14.**  $\exp (\mathbf{A}) \preceq \mathbf{I} + \mathbf{A} + \mathbf{A}^2$  for  $\|\mathbf{A}\| \leq 1$ .

If  $\mathbf{X}$  is zero-mean and  $\|\mathbf{X}\| \leq 1$ , this means that  $\mathbb{E} [\exp (\mathbf{X})] \preceq \mathbf{I} + \mathbb{E} [\mathbf{X}^2]$ , which is how we end up wanting to bound the matrix-valued variance  $\mathbb{E} [\mathbf{X}^2]$ . In the rest of this Sub-section, we're going to see the matrix-valued variance of the stopped martingale is bounded throughout the algorithm.

Firstly, we note that for a single edge sample  $\tilde{\mathbf{X}}_{i,e}$ , by Lemma 10.4.8, we have that

$$\left\| \tilde{\mathbf{X}}_{i,e} \right\| \leq \left\| \Phi (\mathbf{Y}_{\pi(i),e} - \mathbb{E} [\mathbf{Y}_{\pi(i),e}]) \right\| \leq 1/K,$$

using that  $\|\mathbf{A} - \mathbf{B}\| \leq \max(\|\mathbf{A}\|, \|\mathbf{B}\|)$ , for  $\mathbf{A}, \mathbf{B} \succeq \mathbf{0}$ , and  $\|\mathbb{E} [\mathbf{A}]\| \leq \mathbb{E} [\|\mathbf{A}\|]$  by Jensen's inequality.

Thus, if  $0 < \theta \leq K$ , we have that

$$\begin{aligned} \mathbb{E} [\exp (\theta \tilde{\mathbf{X}}_{i,e}) \mid \text{preceding samples}] &\preceq \mathbf{I} + \mathbb{E} [(\theta \tilde{\mathbf{X}}_{i,e})^2 \mid \text{preceding samples}] \quad (10.11) \\ &\preceq \mathbf{I} + \frac{1}{K} \theta^2 \cdot \mathbb{E} [\Phi (\mathbf{Y}_{\pi(i),e}) \mid \text{preceding samples}] \end{aligned}$$

### 10.4.8 The Overall Mean-Trace-Exponential Bound

We will use  $\mathbb{E}_{(<i)}$  to denote expectation over variables preceding the  $i$ th elimination step. We are going to refrain from explicitly writing out conditioning in our expectations, but any *inner* expectation that appears inside another *outer* expectation should be taken as conditional on the outer expectation. We are going to use  $d_i$  to denote the multi-edge degree of vertex  $\pi(i)$  in  $\mathbf{S}_{i-1}$ . This is exactly the number of edge samples in the  $i$ th elimination. Note that there is no elimination at step  $n$  (the algorithm is already finished). As a notational convenience, let's write  $\hat{n} = n - 1$ . With all that in mind, we bound the mean-trace-exponential for some parameter  $0 < \theta \leq 0.5/\sqrt{K}$

$$\begin{aligned}
& \mathbb{E} \text{Tr} \left( \exp(\theta \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i) \right) \\
&= \mathbb{E}_{(<\hat{n})} \mathbb{E}_{\pi(\hat{n})} \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},1}} \cdots \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}-1}} \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}}} \text{Tr} \exp \left( \underbrace{\sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \sum_{e=1}^{d_{\hat{n}}-1} \theta \tilde{\mathbf{X}}_{\hat{n},e}}_{\mathbf{H}} + \theta \tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}} \right) \\
& \quad \tilde{\mathbf{X}}_{\hat{n},1} \dots, \tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}} \text{ are independent conditional on } (<\hat{n}), \pi(\hat{n}) \\
&\leq \mathbb{E}_{(<\hat{n})} \mathbb{E}_{\pi(\hat{n})} \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},1}} \cdots \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}-1}} \text{Tr} \exp \left( \sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \sum_{e=1}^{d_{\hat{n}}-1} \theta \tilde{\mathbf{X}}_{\hat{n},e} + \frac{1}{K} \theta^2 \cdot \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}}} \Phi(\mathbf{Y}_{\pi(\hat{n}), d_{\hat{n}}}) \right) \\
& \quad \text{By Equation (10.11) and Corollary 10.4.13 .} \\
& \vdots \quad \text{Repeat for each multi-edge sample } \tilde{\mathbf{X}}_{\hat{n},1} \dots, \tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}-1} \\
&\leq \mathbb{E}_{(<\hat{n})} \mathbb{E}_{\pi(\hat{n})} \text{Tr} \exp \left( \sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \sum_{e=1}^{d_{\hat{n}}} \frac{1}{K} \theta^2 \cdot \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},e}} \Phi(\mathbf{Y}_{\pi(\hat{n}), e}) \right) \\
&= \mathbb{E}_{(<\hat{n})} \mathbb{E}_{\pi(\hat{n})} \text{Tr} \exp \left( \sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \frac{1}{K} \theta^2 \Phi(\text{CLIQUE}(\pi(\hat{n}), \mathbf{S}_{\hat{n}-1})) \right)
\end{aligned} \tag{10.12}$$

To further bound the this quantity, we now need to deal with the random choice of  $\pi(\hat{n})$ . We'll be able to use this to bound the trace-exponential in a very strong way. From a random matrix perspective, it's the following few steps that give the analysis it's surprising strength.

We can treat  $\frac{1}{K} \theta^2 \Phi(\text{CLIQUE}(\pi(\hat{n}), \mathbf{S}_{\hat{n}-1}))$  as a random matrix. It is not zero-mean, but we can still bound the trace-exponential using Corollary 10.4.13.

We can also bound the expected matrix exponential in that case, using a simple corollary of Lemma 10.4.14.

**Corollary 10.4.15.**  $\exp(\mathbf{A}) \preceq \mathbf{I} + (1 + R)\mathbf{A}$  for  $\mathbf{0} \preceq \mathbf{A}$  with  $\|\mathbf{A}\| \leq R \leq 1$ .

*Proof.* The conclusion follows after observing that for  $\mathbf{0} \preceq \mathbf{A}$  with  $\|\mathbf{A}\| \leq R$ , we have  $\mathbf{A}^2 \preceq R\mathbf{A}$ . We can see this by considering the spectral decomposition of  $\mathbf{A}$  and dealing with each eigenvalue separately.  $\square$

Next, we need a simple structural observation about the cliques created by elimination:

**Claim 10.4.16.**

$$\text{CLIQUE}(\pi(i), \mathbf{S}_i) \preceq \text{STAR}(\pi(i), \mathbf{S}_i) \preceq \mathbf{S}_i$$

*Proof.* The first inequality is immediate from  $\text{CLIQUE}(\pi(i), \mathbf{S}_i) \preceq \text{CLIQUE}(\pi(i), \mathbf{S}_i) + \mathbf{l}_i \mathbf{l}_i^\top = \text{STAR}(\pi(i), \mathbf{S}_i)$ . The latter inequality  $\text{STAR}(\pi(i), \mathbf{S}_i) \preceq \mathbf{S}_i$  follows from the star being a subgraph of the whole Laplacian  $\mathbf{S}_i$ .  $\square$

Next we make use of the fact that  $\tilde{\mathbf{X}}_i$  is from the difference sequence of the *stopped* martingale. This means we can assume

$$\mathbf{S}_i \preceq 1.5\mathbf{L},$$

since otherwise  $\tilde{\mathbf{X}}_i = \mathbf{0}$  and we get an even better bound on the trace-exponential. To make this formal, in Equation (10.12), we ought to do a case analysis that also includes the case  $\tilde{\mathbf{X}}_i = \mathbf{0}$  when the martingale has stopped, but we omit this.

Thus we can conclude by Claim 10.4.16 that

$$\|\Phi(\text{CLIQUE}(\pi(i), \mathbf{S}_i))\| \leq 1.5.$$

By our assumption  $0 < \theta \leq 0.5/\sqrt{K}$ , we have  $\left\| \frac{1}{K}\theta^2\Phi(\text{CLIQUE}(\pi(i), \mathbf{S}_{i-1})) \right\| \leq 1$ , so that by Corollary 10.4.15,

$$\begin{aligned} \mathbb{E}_{\pi(i)} \exp \left( \frac{1}{K}\theta^2\Phi(\text{CLIQUE}(\pi(i), \mathbf{S}_{i-1})) \right) &\preceq \mathbf{I} + \frac{2}{K}\theta^2 \mathbb{E}_{\pi(i)} \Phi(\text{CLIQUE}(\pi(i), \mathbf{S}_{i-1})) \\ &\preceq \mathbf{I} + \frac{2}{K}\theta^2 \mathbb{E}_{\pi(i)} \Phi(\text{STAR}(\pi(i), \mathbf{S}_{i-1})) \end{aligned} \quad (10.13)$$

by Claim 10.4.16.

Next we observe that, because every multi-edge appears in exactly two stars, and  $\pi(i)$  is chosen uniformly at random among the  $n+1-i$  vertices that  $\mathbf{S}_{i-1}$  is supported on, we have

$$\mathbb{E}_{\pi(i)} \text{STAR}(\pi(i), \mathbf{S}_{i-1}) = 2 \frac{1}{n+1-i} \mathbf{S}_{i-1}.$$

And, since we assume  $\mathbf{S}_i \preceq 1.5\mathbf{L}$ , we further get

$$\mathbb{E}_{\pi(i)} \exp \left( \frac{1}{K}\theta^2\Phi(\text{CLIQUE}(\pi(i), \mathbf{S}_{i-1})) \right) \preceq \mathbf{I} + \frac{6\theta^2}{K(n+1-i)} \mathbf{I}.$$

We can combine this with Equation (10.12) and Corollary 10.4.13 to get

$$\begin{aligned} &\mathbb{E} \text{Tr} \left( \exp \left( \theta \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i \right) \right) \\ &\leq \mathbb{E}_{(<\hat{n})} \mathbb{E}_{\pi(\hat{n})} \text{Tr} \exp \left( \sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \frac{1}{K}\theta^2\Phi(\text{CLIQUE}(\pi(\hat{n}), \mathbf{S}_{\hat{n}-1})) \right) \\ &\leq \mathbb{E}_{(<\hat{n})} \text{Tr} \exp \left( \sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \frac{6\theta^2}{K(n+1-i)} \mathbf{I} \right) \end{aligned}$$

And by repeating this analysis for each term  $\tilde{\mathbf{X}}_i$ , we get

$$\begin{aligned}\mathbb{E} \operatorname{Tr} \left( \exp(\theta \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i) \right) &\leq \operatorname{Tr} \exp \left( \sum_{i=1}^{\hat{n}} \frac{6\theta^2}{K(n+1-i)} \mathbf{I} \right) \\ &\leq \operatorname{Tr} \exp \left( \frac{7\theta^2 \log(n)}{K} \mathbf{I} \right) \\ &= n \exp \left( \frac{7\theta^2 \log(n)}{K} \right)\end{aligned}$$

Then, by choosing  $K = 200 \log^2 n$  and  $\theta = 0.5\sqrt{K}$ , we get

$$\exp(-0.5\theta) \mathbb{E} \operatorname{Tr} \left( \exp(\theta \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i) \right) \leq \exp(-0.5\theta) n \exp \left( \frac{7\theta^2 \log(n)}{K} \right) \leq 1/n^5.$$

$\mathbb{E} \operatorname{Tr} \left( \exp(-\theta \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i) \right)$  can be bounded by an identical argument, so that Equation (10.10) gives

$$\Pr \left[ \left\| \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i \right\| \geq 0.5 \right] \leq 2/n^5.$$

Thus we have established  $\left\| \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i \right\| \leq 0.5$  with high probability (Equation (10.9)), and this in turn implies Equation (10.2), and finally Equation (10.1):

$$0.5\mathbf{L} \preceq \mathcal{L}\mathcal{L}^\top \preceq 1.5\mathbf{L}.$$

Now, all that's left to note is that the running time is linear in the multi-edge degree of the vertex being eliminated in each iteration (and this also bounds the number of non-zero entries being created in  $\mathcal{L}$ ). The total number of multi-edges left in the remaining graph stays constant at  $Km = O(m \log^2 n)$ . Thus the expected degree in the  $i$ th elimination is  $Km/(n+i-1)$ , because the remaining number of vertices is  $n+i-1$ . Hence the total running time and total number of non-zero entries created can both be bounded as

$$Km \sum_i 1/(n+i-1) = O(m \log^3 n).$$

We can further prove that the bound  $O(m \log^3 n)$  on running time and number of non-zeros in  $\mathbf{L}$  holds with high probability (e.g.  $1 - 1/n^5$ ). To show this, we essentially need a scalar Chernoff bound, in except the degrees are in fact not independent, and so we need a scalar martingale concentration result, e.g. Azuma's Inequality. This way, we complete the proof of Theorem 10.2.4.

## **Part III**

# **Combinatorial Graph Algorithms**

# Chapter 11

## Classical Algorithms for Maximum Flow I

### 11.1 Maximum Flow

In this chapter, we will study the *Maximum Flow Problem* and discuss some classical algorithms for finding solutions to it.

**Setup.** Consider a directed graph  $G = (V, E, \mathbf{c})$ , where  $V$  denotes vertices,  $E$  denotes edges and  $\mathbf{c} \in \mathbb{R}^E$ ,  $\mathbf{c} \geq 0$  denotes edge capacities. In contrast to earlier chapters, the direction of each edge will be important, and not just as a book keeping tool (which is how we previously used it in undirected graphs). We consider an edge  $(u, v) \in E$  to be *from*  $u$  and *to*  $v$ . Edge capacities are associated with directed edges, and we allow both edge  $(u, v)$  and  $(v, u)$  to exist, and they may have different capacities.

A *flow* is any vector  $\mathbf{f} \in \mathbb{R}^E$ .

We say that a flow is feasible when  $\mathbf{0} \leq \mathbf{f} \leq \mathbf{c}$ . The constraint  $\mathbf{0} \leq \mathbf{f}$  ensures that the flow respects edge directions, while the constraint  $\mathbf{f} \leq \mathbf{c}$  ensures that the flow does not exceed capacity on any edge.

We still define the edge-vertex incidence matrix  $\mathbf{B} \in \mathbb{R}^{V \times E}$  of  $G$  by

$$\mathbf{B}(v, e) = \begin{cases} 1 & \text{if } e = (u, v) \\ -1 & \text{if } e = (v, u) \\ 0 & \text{o.w.} \end{cases}$$

As in the undirected case, a *demand vector* is a vector  $\mathbf{d} \in \mathbb{R}^V$ . And as in the undirected case, we can express the net flow constraint that  $\mathbf{f}$  routes the demand  $\mathbf{d}$  by

$$\mathbf{B}\mathbf{f} = \mathbf{d}.$$

We will focus on the case:

$$\mathbf{B}\mathbf{f} = F(-\mathbf{e}_s + \mathbf{e}_t) = F\mathbf{b}_{st}.$$

Flows that satisfy the above equation for some scalar  $F$  are called  $s$ - $t$  flows where  $s, t \in V$ . The vertex  $s$  is called the source,  $t$  is called the sink and  $\mathbf{e}_s, \mathbf{e}_t$  are indicator vectors for source and sink nodes respectively. The vector  $\mathbf{b}_{st}$  has -1 at source and 1 at sink. The maximum flow can be expressed as a linear program as follows

$$\begin{aligned} & \max_{\mathbf{f} \in \mathbb{R}^E, F} F \\ \text{s.t. } & \mathbf{B}\mathbf{f} = F\mathbf{b}_{st} \\ & 0 \leq \mathbf{f} \leq \mathbf{c} \end{aligned} \tag{11.1}$$

We use  $\text{val}(\mathbf{f})$  to denote  $F$  when  $\mathbf{B}\mathbf{f} = F\mathbf{b}_{st}$ .

## 11.2 Flow Decomposition

We now look at a way of simplifying flows

**Definition 11.2.1.** A  $s$ - $t$  path flow is a flow  $\mathbf{f} \geq \mathbf{0}$  that can be written as

$$\mathbf{f} = \alpha \sum_{e \in p} \mathbf{e}_e$$

where  $p$  is a simple path from  $s$  to  $t$ .

**Definition 11.2.2.** A cycle flow is a flow  $\mathbf{f} \geq \mathbf{0}$  that can be written as a cycle i.e.

$$\mathbf{f} = \alpha \sum_{e \in c} \mathbf{e}_e$$

where  $c$  is a simple cycle.

**Lemma 11.2.3** (The path-cycle decomposition lemma). *Any  $s$ - $t$  flow  $\mathbf{f}$  can be decomposed into a sum of  $s$ - $t$  path flows and cycle flows such that the sum contains at most  $\text{nnz}(\mathbf{f})$  terms. Note  $\text{nnz}(\mathbf{f}) \leq |E|$ .*

*Proof.* We perform induction on the number of edges with non-zero flow, which we denote by  $\text{nnz}(\mathbf{f})$ . Note that by “the support of  $\mathbf{f}$ ”, we mean the set  $\{e \in E : \mathbf{f}(e) \neq 0\}$ .

**Base case:**  $\mathbf{f} = 0$ : nothing to do.

**Inductive step:** Try to find a path from  $s$  to  $t$  OR a cycle in the support of  $\mathbf{f}$ .

“Path” case. If there exists such a  $s$ - $t$  path, let  $\alpha$  be the minimum flow value along the edges of the path, i.e.

$$\alpha = \min_{(a,b) \in p} \mathbf{f}(a,b)$$

$$\mathbf{f}' = \alpha \sum_{e \in p} \mathbf{e}_e$$

Update the flow  $\mathbf{f}$  by

$$\mathbf{f} \leftarrow \mathbf{f} - \mathbf{f}'$$

The value of the flow will still be non-negative after this update as we subtracted the minimum entry along any positive edge on the path. The number of non-zeros,  $\text{nnz}(\mathbf{f})$ , went down by at least one. Note that the updated  $\mathbf{f}$  must again be an  $s$ - $t$  flow, as it is the difference of two  $s$ - $t$  flows.

“Cycle” case. Suppose we find a cycle  $c$  in the support of  $\mathbf{f}$ . Let  $\alpha$  be the minimum flow value along the edges of the cycle, i.e.

$$\alpha = \min_{(a,b) \in c} \mathbf{f}(a,b)$$

$$\mathbf{f}' = \alpha \sum_{e \in c} \mathbf{e}_e$$

Update the flow  $\mathbf{f}$  by

$$\mathbf{f} \leftarrow \mathbf{f} - \mathbf{f}'$$

As in the path case,  $\mathbf{f}$  stays non-negative, and number of non-zeros,  $\text{nnz}(\mathbf{f})$ , goes down by at least one. Note that the updated  $\mathbf{f}$  must again be an  $s$ - $t$  flow, as it is the difference of two  $s$ - $t$  flows.

“No path or cycle” case. Suppose we can find neither a path nor a cycle, and  $\mathbf{f} \neq \mathbf{0}$ . Then there must be an edge  $(u, v)$  with non-zero flow leading into a vertex  $v \neq s, t$  and with no outgoing edge from  $v$  in the support of  $\mathbf{f}$ . In that case, we must have  $(\mathbf{B}\mathbf{f})(v) > 0$ . But since  $v \neq t$ , this contradicts  $\mathbf{B}\mathbf{f} = F\mathbf{b}_{st}$ . So this case cannot occur.  $\square$

**Lemma 11.2.4.** *In any  $s$ - $t$  max flow problem instance, there is an optimal flow  $\mathbf{f}^*$  with a path-cycle decomposition that has only paths and no cycles.*

*Proof.* Let  $\tilde{\mathbf{f}}$  be an optimal flow. Let  $\mathbf{f}^*$  be the sum of path flows in the path cycle decomposition of  $\tilde{\mathbf{f}}$ . They route the same flow (as cycles contribute to no net flow from  $s$  to  $t$ ). Thus

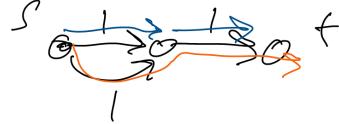
$$\mathbf{B}\mathbf{f}^* = \mathbf{B}\tilde{\mathbf{f}}$$

and hence  $\text{val}(\mathbf{f}^*) = \text{val}(\tilde{\mathbf{f}})$ . Furthermore

$$\mathbf{0} \leq \mathbf{f}^* \leq \tilde{\mathbf{f}} \leq \mathbf{c}$$

The first inequality follows from  $\mathbf{f}^*$  being a sum of positive path flow. The second inequality holds as  $\mathbf{f}^*$  is upper bounded in every single entry by  $\tilde{\mathbf{f}}$ , because we reduced it by positive entry cycles. The third inequality holds because  $\tilde{\mathbf{f}}$  is a feasible flow, so it is upper bounded by the capacities.  $\square$

An optimal flow solving the maximum flow problem may not be unique. For example, consider the graph below with source  $s$  and sink  $t$ :



There are two optimal paths in this example. Maximum flow is a convex optimization problem but not a strongly convex problem as the solutions are not unique, and this is part of what makes it hard to solve.

### 11.3 Cuts and Minimum Cuts

The decomposition shown earlier provides a way to show that the maximum flow in a graph is upper bounded by constructing graph cuts.

Given a vertex subset  $S \subseteq V$ , we say that  $(S, V \setminus S)$  is a *cut* in  $G$  and that the value of the cut is

$$c_G(S) = \sum_{e \in E \cap (S \times V \setminus S)} \mathbf{c}(e).$$

Note that in a directed graph, only edges crossing the cut going *from*  $S$  and *to*  $V \setminus S$  count toward the cut.

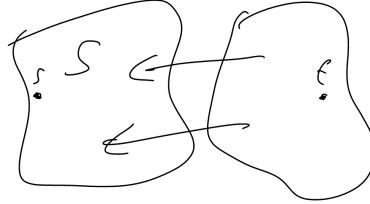


Figure 11.1: Example of a cut: No edges go from  $S$  to  $(S, V \setminus S)$ , and so the value of this cut is zero.

**Definition.** ( $s$ - $t$  cuts). We define an  $s$ - $t$  cut to be a subset  $S \subset V$ , where  $s \in S$  and  $t \in V \setminus S$ .

**A decision problem:** “Given an instance of the Maximum Flow problem, is there a flow from  $s$  to  $t$  such that  $\mathbf{f} \neq \mathbf{0}$ ? ”

**If YES:** We can decompose this flow into  $s$ - $t$  path flows and cycle flows.

**If NO:** There is no flow path from  $s$  to  $t$ . Let  $S$  be the set of vertices reachable from source  $s$ . Then  $(S, V \setminus S)$  is a cut in the graph, with no edges crossing from  $S$  to  $V \setminus S$ . Figure 11.1 gives an example.

**Upper bounding the maximum possible flow value.** How can we recognize a maximum flow? Is there a way to confirm that a flow is of maximum value?

We can now introduce the *Minimum Cut* problem.

$$\begin{aligned} & \min_{S \subseteq V} c_G(S) \\ & \text{s.t. } s \in V \text{ and } t \notin V \end{aligned} \tag{11.2}$$

The Minimum Cut problem can also be phrased as a linear program, although we won’t see that today.

We’d like to obtain a tight connection between flows and cuts in the graph. As a first step, we won’t get that, but we can at least observe that the value of any  $s$ - $t$  cut provides an upper bound to the maximum possible  $s$ - $t$  flow value.

**Theorem 11.3.1** (Max Flow  $\leq$  Min Cut). *The maximum  $s$ - $t$  flow value in a directed graph  $G$  (Program (11.1)) is upper bounded by the minimum value of any  $s$ - $t$  cut (Program (11.2)). I.e. if  $S$  is an  $s$ - $t$  cut, and  $\mathbf{f}$  a feasible  $s$ - $t$  flow then*

$$\text{val}(\mathbf{f}) \leq c_G(S)$$

*And in particular this holds for any minimum cut  $S^*$  and maximum flow  $\mathbf{f}^*$ , i.e.  $\text{val}(\mathbf{f}^*) \leq c_G(S^*)$ .*

*Proof.* Consider any feasible flow  $0 \leq \mathbf{f} \leq \mathbf{c}$  and a cut  $S, T = V \setminus S$ . Consider a path-cycle decomposition of  $\mathbf{f}$ , where each  $s$ - $t$  path must cross the cut going forward from  $S$  to  $T$  at least once. We pick a cut  $S, T$  with source on one side and sink on the other side as shown in the Figure (11.2). Every time the path flow passes through the cut, it has to use one of the edges that connect  $S$  and  $T$ . Total amount of flow crossing the cut is bounded above by total amount of capacity of the cut, otherwise the capacities would be violated, thus  $\text{val}(\mathbf{f}) \leq c_G(S, T) = \sum_{e \in E \cap S \times T} \mathbf{c}(e)$ .

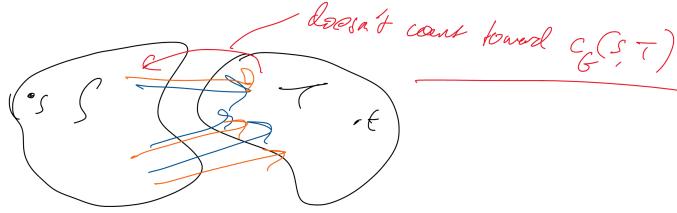


Figure 11.2:  $s$ - $t$  Cut

Note that the edges from  $T$  to  $S$  do not count towards the cut. The above equation holds for all flows with all  $s$ - $t$  cuts. This implies that **max flow  $\leq$  min cut**.  $\square$

This theorem is an instance of a general pattern, known as *weak duality*. Weak duality is a relationship between two optimization programs, a maximization problem and a minimization problem, where any solution to the former has its value upper bounded by the value of any solution to the latter.

## 11.4 Algorithms for Max flow

**How can we find a good flow?**

---

**Algorithm 4:** A first attempt – bad idea?

---

- 1  $f \leftarrow 0;$
  - 2 **repeat**
  - 3     Find an  $s$ - $t$  path flow  $\tilde{f}$  that is feasible with respect to  $c - f$ .
  - 4      $f \leftarrow f + \tilde{f}$
- 

**Does Algorithm 1 work?** Consider the graph below with directed edges with capacities 1 at every edge. If we make a single path update as shown by the orange lines in Figure 11.3, then afterwards, using the remaining capacity, there's no path flow we can route, as shown in Figure 11.3. But the max flow is 2, as shown by the flow on orange edges in Figure 11.5. So, the above algorithm does not always find a maximum flow: It can get stuck at a point where it cannot make progress despite not having reached the maximum value.

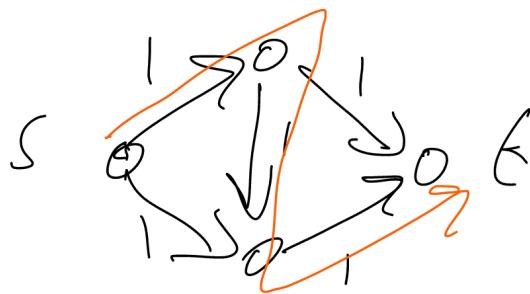


Figure 11.3: Sending a unit  $s$ - $t$  path flow through the graph.

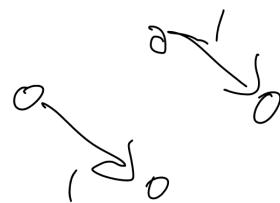


Figure 11.4: Remaining edge capacities after sending a path flow through the graph as depicted in Figure 11.3.

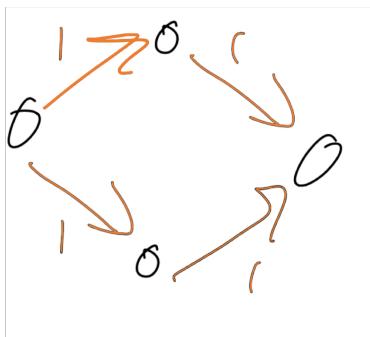


Figure 11.5: The flow depicted routes two units from  $s$  to  $t$ .

**A better approach.** It turns out we can fix the problem with the previous algorithm using a simple fix. This idea is known as *residual graphs*.

---

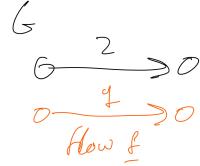
**Algorithm 5:** Better Idea (Residual Graph)

---

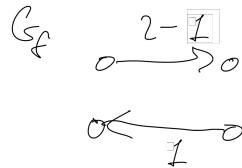
- 1  $f \leftarrow 0;$
  - 2 **repeat**
  - 3     Find an  $s$ - $t$  path flow  $\tilde{f}$  that is feasible with respect to  $-f \leq \tilde{f} \leq c - f$ .
  - 4      $f \leftarrow f + \tilde{f}$
-

The  $-\mathbf{f}(e)$  can be treated as an edge going in the other direction with capacity  $\mathbf{f}(e)$ . By convention, an edge in  $G$  with  $\mathbf{f}(e) = \mathbf{c}(e)$  is called *saturated*, and we do not include the edge in  $G_f$ . The graph defined above with such capacities is called **the residual graph of  $f$** ,  $G_f$ .  $G_f$  is only defined for feasible  $f$ , since otherwise the constraint  $\tilde{\mathbf{f}} \leq \mathbf{c} - \mathbf{f}$  gives trouble.

Suppose we start with a graph having a single edge with capacity 2 and we introduce a flow of 1 unit.



The residual graph  $G_f$  has an edge with capacity 1 going forward and  $-1$  capacity going forward, but we can treat the latter as  $+1$  capacity going backwards. So it is an edge that allows you to undo the choice made to send flow along that direction.



Let us consider the same example with its residual graph. The original graph is shown in Figure 11.6.

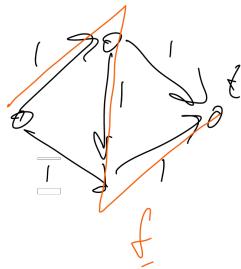


Figure 11.6: Original graph  $G$  and an  $s-t$  path flow in  $G$  shown in orange.

The residual graph for the same is shown in Figure 11.7.

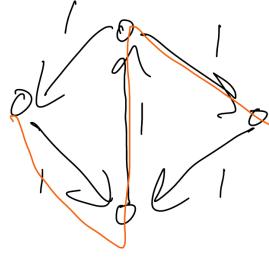


Figure 11.7: The residual graph w.r.t. the flow from Figure 11.6, and new  $s-t$  path flow which is feasible in the residual graph.

Adding both flows together, we get the paths as shown in Figure 11.8 with value 2, which is the optimum.



Figure 11.8: Maximum flow in the graph.

Let's prove some important properties of residual graphs.

**Lemma 11.4.1.** Suppose  $\mathbf{f}, \hat{\mathbf{f}}$  are feasible in  $G$ . Then this implies that  $\hat{\mathbf{f}} - \mathbf{f}$  (where negative entries count as flow in opposite direction) is feasible in  $G_f$ .

*Proof.*  $\mathbf{0} \leq \mathbf{f} \leq \mathbf{c}$  and  $\mathbf{0} \leq \tilde{\mathbf{f}} \leq \mathbf{c}$ , this implies  $-\mathbf{f} \leq \tilde{\mathbf{f}} - \mathbf{f} \leq \mathbf{c} - \mathbf{f}$ . Hence, proved.  $\square$

**Lemma 11.4.2.** Suppose that  $\mathbf{f}$  is feasible in  $G$  and  $\hat{\mathbf{f}}$  is feasible in  $G_f$ . Then,  $\mathbf{f} + \hat{\mathbf{f}}$  is feasible in  $G$ .

*Proof.*  $\mathbf{0} \leq \mathbf{f} \leq \mathbf{c}$  and  $-\mathbf{f} \leq \tilde{\mathbf{f}} \leq \mathbf{c} - \mathbf{f}$ , this implies  $\mathbf{0} \leq \tilde{\mathbf{f}} + \mathbf{f} \leq \mathbf{c}$ .  $\square$

**Lemma 11.4.3.** A feasible  $\mathbf{f}$  is optimal if and only if  $t$  is not reachable from  $s$  in  $G_f$ .

*Proof.* Let  $\mathbf{f}$  be optimal, and suppose  $t$  is reachable from  $s$  in  $G_f$  then, we can find a  $s-t$  path flow  $\tilde{\mathbf{f}}$  that is feasible in  $G_f$ , and  $\text{val}(\mathbf{f} + \tilde{\mathbf{f}}) > \text{val}(\mathbf{f})$ .  $\mathbf{f} + \tilde{\mathbf{f}}$  is feasible in  $G$  by Lemma 11.4.2. This is a contradiction, as we assumed  $\mathbf{f}$  was supposed to be optimal.

Suppose  $t$  is not reachable from  $s$  in  $G_f$ , and  $\mathbf{f}$  is feasible, but not optimal. Let  $\mathbf{f}^*$  be optimal, then by Lemma 11.4.1, the flow  $\mathbf{f}^* - \mathbf{f}$  is feasible in  $G_f$  and  $\text{val}(\mathbf{f}^* - \mathbf{f}) > 0$ . So there exists a  $s-t$  path flow from  $s$  to  $t$  in  $G_f$  (as we can do a path decomposition of  $\mathbf{f}^* - \mathbf{f}$ ). But, this is a contradiction as  $t$  is not reachable from  $s$  in  $G_f$ .  $\square$

**Theorem 11.4.4** (Max Flow = Min Cut theorem). *The maximum flow in a directed graph  $G$  equals the minimum cut.*

*Proof.* Consider the set  $S = \{\text{vertices reachable from } s \text{ in } G_f\}$ . Note that  $\mathbf{f}$  saturates the edge capacities in cut  $S, V \setminus S$  in  $G$ : Consider any edge from  $S$  to  $T := V \setminus S$  in  $G$ . Since this edge does not exist in the residual graph, we must have  $\mathbf{f}(e) = \mathbf{c}(e)$ .

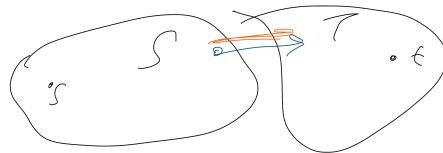


Figure 11.9: The cut between vertices reachable from  $S$  and everything else in  $G_f$  must have all outgoing edges saturated by  $\mathbf{f}$ .

This means that

$$\text{val}(\mathbf{f}) \geq c_G(S, V \setminus S).$$

Since we already know the opposite inequality by weak duality, we have shown that

$$\text{val}(\mathbf{f}) = c_G(S, V \setminus S)$$

which proves the **Max Flow = Min Cut** theorem; also called strong duality.  $\square$

**Remark 11.4.5.** Note that this proof also gives an algorithm to compute an  $s, t$ -min cut given an  $s, t$  maximum flow  $\mathbf{f}^*$ : Take one side of the cut to be the set  $S$  of nodes reachable from  $s$  in the residual graph w.r.t.  $\mathbf{f}^*$ .

### Ford-Fulkerson Algorithm.

---

#### Algorithm 6: Ford-Fulkerson Algorithm

---

- 1 **repeat**
  - 2   | Add update by arbitrary  $s-t$  path flow in  $G_f$  (augment the flow  $\mathbf{f}$  by the path flow)
- 

### Convergence properties and analysis of runtime of Ford-Fulkerson algorithm

- **Does this algorithm terminate?**

The algorithm terminates if the capacities are integers. However for irrational capacities the algorithm may not terminate.

- Does it converge towards the max flow value?

No, it does not converge to max flow value if the updates are poor and the capacities are irrational.

**Lemma 11.4.6.** *Consider Ford-Fulkerson algorithm with integer capacities. The algorithm terminates in  $\text{val}(\mathbf{f}^*)$  augmentations i.e.  $O(m \text{val}(\mathbf{f}^*))$  time.*

*Proof.* Each iteration increases the flow by at least one unit as the capacities in  $G_f$  are integral and each iteration can be computed in  $O(m)$  time.  $\square$

**Can we do better than this?** Suppose we pick the maximum bottleneck capacity (minimum capacity along path) augmenting path. This gives an algorithm that is better in some regimes.

How to pick the maximum bottleneck capacity augmenting path? We are going to perform a binary search on the capacities in  $G_f$ , to find a path with maximum bottleneck capacity. Each time our binary search has picked a threshold capacity, we then try to find an  $s-t$  path flow in  $G_f$  using only edges with capacity above that threshold. If we find a path, the binary search tries a higher capacity next. If we don't find a path, the binary search tries a lower capacity next.

Using this approach, the time to find a single path is  $O(m \log(n))$  where  $m$  is number of edges in the graph. This path must carry at least a  $\frac{1}{m}$  fraction of the total amount of flow left in  $G_f$ . For instance, if  $\hat{F}$  is the amount of flow left in  $G_f$ , then the path must carry  $\frac{\hat{F}}{m}$  flow (from the path decomposition lemma, we know that there exists a decomposition into at most  $m$  paths, and the one carrying the most flow must carry at least the average amount).

So if the flow is integral, the algorithms completes when

$$\left(1 - \frac{1}{m}\right)^T \text{val}(\mathbf{f}^*) < 1$$

where  $T$  is the number of augmentations.

This means

$$T = m \log F$$

$$\begin{aligned} \text{Total time} &= O(m \log n T) \\ &= O(m^2 \log n \log F) \\ &\leq O(m^2 \log n \log mU) \text{ as } F \leq mU \text{ where U is the maximum capacity} \end{aligned}$$

**Current state-of-the-art approaches for Max Flow.** For the interested reader, we'll briefly mention the current state-of-the-art for Maximum Flow algorithms, which is a very active area of research.

- Strongly polynomial time algorithm:
  - has to work with real valued capacities, then the best time  $O(mn)$  by Orlin.
- General integer capacities: For a long time, the state-of-the-art was a result by Goldberg and Rao achieving a runtime of  $O(m \min(\sqrt{m}, n^{2/3}) \log(n) \log(U))$  where  $U$  is the maximum capacity. Very recently, an algorithm that runs in almost-linear time  $m^{1+o(1)} \log(U)$  was given (see <https://arxiv.org/abs/2203.00671>). There are many other recent results prior to this work that we don't have time to cover here.

# Chapter 12

## Classical Algorithms for Maximum Flow II

### 12.1 Overview

In this chapter we continue looking at classical max flow algorithms. We derive Dinic's algorithm which (unlike Ford-Fulkerson) converges in a polynomial number of iterations. We will also show faster convergence on unit capacity graphs. The setting is the same as last chapter: we have a directed graph  $G = (V, E, \mathbf{c})$  with positive capacities on the edges. Our goal is to solve the maximum flow problem on this graph for a given source  $s$  and sink  $t \neq s$ :

$$\begin{aligned} \max_{F > 0, f} \quad & F \quad \text{s.t.} \quad \mathbf{B}\mathbf{f} = F\mathbf{b}_{s,t} \\ & \mathbf{0} \leq \mathbf{f} \leq \mathbf{c} \end{aligned} \tag{12.1}$$

### 12.2 Blocking Flows

Let  $\mathbf{f}$  be any feasible flow in  $G$  and let  $G_f$  be its residual graph. Observe that we can partition the vertices of  $G_f$  into ‘layers’: first  $s$  itself, then all vertices reachable from  $s$  in one hop, then all vertices reachable in two hops, etc. For each vertex  $v \in V$  define its **level** in  $G_f$  as the length of the shortest path in  $G_f$  from  $s$  to  $v$ , denoted by  $\ell_{G_f}(v)$  (or just  $\ell(v)$  if the graph is clear from context). An edge  $(u, v) \in E$  can only take you up ‘one level’: if  $\ell(v) \geq 2 + \ell(u)$  this would imply we can find a shorter  $s$ - $v$  path by appending  $(u, v)$  to the shortest  $s$ - $u$  path. However, edges can ‘drop down’ multiple levels (or be contained in the same level).

A key strategy in Dinic's algorithm will be to focus on ‘short’ augmenting paths. We can use the levels we defined above to isolate a subgraph of  $G_f$  containing all information we

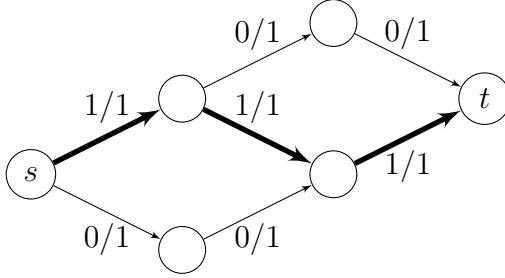


Figure 12.1: A blocking flow is not a maximum flow.

need to find shortest paths.

**Definition 12.2.1.** Call an edge  $(u, v)$  **admissible** if  $\ell(u) + 1 = \ell(v)$ . Let  $L$  be the set of admissible edges in  $G_f$  and define the **level graph** of  $G_f$  to be the subgraph induced by these edges.

**Definition 12.2.2.** Define a **blocking flow** in a residual graph  $G_f$  to be a flow  $\hat{f}$  feasible in  $G_f$ , such that  $\hat{f}$  only uses admissible edges. Furthermore we require that for any  $s-t$  path in the level graph of  $G_f$ ,  $\hat{f}$  saturates at least one edge on this path.

The last condition makes a blocking flow ‘blocking’: it blocks any shortest augmenting path in  $G_f$  by saturating one of its edges. Note however that a blocking flow is not necessarily a maximum flow in the level graph.

## 12.3 Dinic’s Algorithm

We can now formulate Dinic<sup>1</sup>’s algorithm: start with  $f = \mathbf{0}$ , and then repeatedly add to  $f$  a blocking flow  $\hat{f}$  in  $G_f$ , until no more  $s-t$  paths exist in  $G_f$ .

Note that by doing a path decomposition on  $\hat{f}$  and adding these paths to  $f$  one by one, we see that our algorithm is a ‘special case’ of Ford-Fulkerson (with a particular augmenting path strategy) and hence inherits all the behaviors/bounds we proved in the last chapter.

**Lemma 12.3.1.** Let  $f$  be a feasible flow,  $\hat{f}$  a blocking flow in  $G_f$  and define  $f' = f + \hat{f}$  (think of  $f$  and  $f'$  to be the flows at some steps  $k$  and  $k+1$  in the algorithm). Then  $\ell_{G_{f'}}(t) \geq \ell_{G_f}(t) + 1$ .

*Proof.* Let  $L, L'$  be the edge sets of the level graphs of  $G_f, G_{f'}$  respectively. We would now like to show that  $d_{L'}(s, t) \geq 1 + d_L(s, t)$ . Let’s first assume that in fact  $d_{L'}(s, t) < d_L(s, t)$ . Now take a shortest  $s-t$  path in the level graph of  $G_{f'}$ , say  $s = v_0, v_1, \dots, v_{d_{L'}(s,t)} = t$ . Let  $v_j$  be the first vertex along the path such that  $d_{L'}(s, v_j) < d_L(s, v_j)$ . As  $d_{L'}(s, t) < d_L(s, t)$ , such a  $v_j$  must exist.

---

<sup>1</sup>Sometimes also transliterated as Dinitz.

We'd like to understand the edges in the level graph  $L'$ . Let  $E(G_f)$  denote the edges of the  $G_f$ , and let  $\text{rev}(L)$  denote the set of reversed edges of  $L$ . The level graph edges of  $G_{f'}$  must satisfy  $L' \subseteq L \cup \text{rev}(L) \cup (E(G_f) \setminus L)$ .

In our  $s$ - $t$  path in  $L'$ , the vertex  $v_j$  is reached by an edge from  $v_{j-1}$ , and the level of  $v_{j-1}$  in  $G_f$  is at most the level it receives in  $G_{f'}$  so the edge  $(v_{j-1}, v_j) \in L$  skips at least one level forward in  $G_f$ . But, edges in  $L$  do not skip a level in  $G_f$ , and edges in  $\text{rev}(L)$  or  $E(G_f) \setminus L$  do not move from a lower level to higher level in  $G_f$ . So this edge cannot exist and we have reached a contradiction.

Now suppose that  $d_{L'}(s, t) = d_L(s, t)$ . This means there is an  $s$ - $t$  path in  $L'$  using edges in  $L$  – but such a path must contain an edge saturated by the blocking flow  $\hat{\mathbf{f}}$ .

(Note: the reason this path must use only edges in  $L$  is similar to the previous case: the other possible types of edges do not move to higher levels in  $G_f$  at all, making the path too long if we use any of them.)

So we must have  $d_{L'}(s, t) \geq 1 + d_L(s, t)$ . □

An immediate corollary of this lemma is the convergence of Dinic's algorithm:

**Theorem 12.3.2.** *Dinic's algorithm terminates in  $O(n)$  iterations.*

*Proof.* By the previous lemma, the level of  $t$  increases by at least one each iteration. A shortest path can only contain each vertex once (otherwise it would contain a cycle) so the level of any vertex is never more than  $n$ . □

For graphs with unit capacities ( $c = 1$ ) we can prove even better bounds on the number of iterations.

**Theorem 12.3.3.** *On unit capacity graphs, Dinic's algorithm terminates in*

$$O(\min\{m^{1/2}, n^{2/3}\})$$

*iterations.*

*Proof.* We prove the two bounds separately.

1. Suppose we run Dinic's algorithm for  $k$  iterations, obtaining a flow  $\mathbf{f}$  that is feasible but not necessarily yet optimal. Let  $\tilde{\mathbf{f}}$  be an optimal flow in  $G_f$  (i.e.  $\mathbf{f} + \tilde{\mathbf{f}}$  is a maximum flow for the whole graph), and consider a path decomposition of  $\tilde{\mathbf{f}}$ . Since after  $k$  iterations any  $s$ - $t$  path has length  $k$  or more, we use up a total capacity of at least  $\text{val}(\tilde{\mathbf{f}})k$  across all edges. But the edges in  $G_f$  are either edges from the original graph  $G$  or their reversals (but never both) meaning the total capacity of  $G_f$  is at most  $m$ , hence  $\text{val}(\tilde{\mathbf{f}}) \leq m/k$ .

Recalling our earlier observation that our algorithm is a special case of Ford-Fulkerson, this implies that our algorithm will terminate after at most another  $m/k$  iterations.

Hence the number of iterations is bounded by  $k + m/k$  for any  $k > 0$ . Substituting  $k = \sqrt{m}$  gives the first desired bound.

2. Suppose again that we run Dinic's algorithm for  $k$  iterations obtaining a flow  $f$ . The level graph of  $G_f$  partitions the vertices into sets  $D_i = \{u \mid \ell(u) = i\}$  for  $i \geq 0$ . As shown before, the sink  $t$  must be at a level of at least  $k$ , meaning we have at least this many non-empty levels starting from  $D_0 = \{s\}$ . To simplify, discard all vertices and levels beyond  $D_{\ell(t)}$ .

Now, consider choosing a level  $I$  uniformly at random from  $\{1, \dots, k\}$ . Since there are at most  $n - 1$  vertices in total across these levels,  $\mathbb{E}[|D_I|] \leq (n - 1)/k$ , and by Markov's inequality

$$\Pr[|D_I| \geq 2n/k] \leq \frac{(n - 1)/k}{2n/k} < \frac{1}{2}.$$

Thus strictly more than  $k/2$  of the levels  $i \in \{1, \dots, k\}$  have  $|D_i| < 2n/k$  and so there must be two adjacent levels  $j, j + 1$  for which this upper bound on the size holds. There can be at most  $|D_j| \cdot |D_{j+1}| \leq 4n^2/k^2$  edges between these levels, and by the min-cut max-flow theorem we saw in the previous chapter, this is an upper bound on the flow in  $G_f$ , and hence on the number of iterations still needed for our algorithm to terminate.

This means the number of iterations is bounded by  $k + 4n^2/k^2$  for any  $k > 0$ , which is  $O(n^{2/3})$  at  $k = 2n^{2/3}$ .

□

## 12.4 Finding Blocking Flows

What has been missing from our discussion so far is the process of actually finding a blocking flow. We can achieve this using repeated depth-first search. We repeatedly do a search in the level graph (so only using edges  $L$ ) for  $s$ - $t$  paths and augment these. We erase edges whose subtrees have been exhausted and do not contain any augmenting paths to  $t$ . Pseudocode for the algorithm is given below.

---

**Algorithm 7:** FINDBLOCKINGFLOW

---

```
1  $\mathbf{f} \leftarrow \mathbf{0}$ ;  
2  $H \leftarrow L$ ;  
3 repeat  
4    $P \leftarrow \text{DFS}(s, H, t)$ ;  
5   if  $P \neq \emptyset$  then  
6     Let  $\hat{\mathbf{f}}$  be a flow that saturates path  $P$ .  
7      $\mathbf{f} \leftarrow \mathbf{f} + \hat{\mathbf{f}}$ ;  
8     Remove from  $H$  all edges saturated by  $\hat{\mathbf{f}}$ .  
9   else  
10    return  $\mathbf{f}$ ;  
11 end
```

---

---

**Algorithm 8:** DFS( $u, H, t$ )

---

```
1 if  $u = t$  then  
2   return the path  $P$  on the dfs-stack.  
3 end  
4 for  $(u, v) \in H$  do  
5    $P \leftarrow \text{DFS}(v, H, t)$ ;  
6   if  $P \neq \emptyset$  then  
7     return  $P$ ;  
8   else  
9     Erase  $(u, v)$  from  $H$ .  
10  end  
11 end  
12 return  $\emptyset$ ;
```

---

**Lemma 12.4.1.** *For general graphs, FINDBLOCKINGFLOW returns a blocking flow in  $O(nm)$  time. Hence Dinic's algorithm runs in  $O(n^2m)$  time on general capacity graphs.*

*Proof.* First, consider the amount of work spent pushing edges onto the stack which eventually results in augmentation along an  $s-t$  path consisting of those edges (i.e. adding flow along that path). Since each augmenting path saturates at least one edge, we do at most  $m$  augmentations. Each path has length at most  $n$ . Thus the total amount of work pushing these edges to the stack, and removing them from the stack upon augmentation, and deleting saturated edges, can be bounded by  $O(mn)$ . Now consider the work spent pushing edges onto the stack which are later deleted because we “retreat” after not finding an  $s-t$ . An edge can only be pushed onto the stack once in this way, since it is then deleted. So the total amount spent pushing edges to the stack and deleting them this way is  $O(m)$ .  $\square$

**Lemma 12.4.2.** *For unit capacity graphs, FINDBLOCKINGFLOW returns a blocking flow in  $O(m)$  time. Hence Dinic's algorithm runs in  $O(\min\{m^{3/2}, mn^{2/3}\})$  time on unit capacity graphs.*

*Proof.* When our depth-first search traverses some edge  $(u, v)$ , one of two things will happen: either we find no augmenting path in the subtree of  $v$ , leading to the erasure of this edge, or we find an augmenting path which will necessarily saturate  $(u, v)$ , again leading to its erasure. This means each edge will be traversed at most once by the depth-first search.  $\square$

Another interesting bound on the runtime, which we will not prove here, is that Dinic's algorithm will run in  $O(|E| \sqrt{|V|})$  time on bipartite matching graphs.

## 12.5 Minimum Cut as a Linear Program

Finally we show that minimum cut may be formulated as a linear program. For a subset  $S \subseteq V$  write  $c_G(S)$  for the sum of  $\mathbf{c}(e)$  over all edges  $e = (u, v)$  that cross the cut, i.e.  $u \in S, v \notin S$ . Note that ‘reverse’ edges are not counted in this cut. The minimum cut problem asks us to find some  $S \subseteq V$  such that  $s \in S$  and  $t \notin S$  such that  $c_G(S)$  is minimal.

$$\begin{aligned} \min_{S \subseteq V} \quad & c_G(S) \\ \text{s.t.} \quad & s \in S \\ & t \notin S \end{aligned} \tag{12.2}$$

We claim this is equivalent to the following minimization problem:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^V} \quad & \sum_{e \in E} \mathbf{c}(e) \max \{ \mathbf{b}_e^\top \mathbf{x}, 0 \} \\ \text{s.t.} \quad & \mathbf{x}(s) = 0 \\ & \mathbf{x}(t) = 1 \\ & \mathbf{0} \leq \mathbf{x} \leq \mathbf{1} \end{aligned} \tag{12.3}$$

Recall  $\mathbf{b}_e^\top \mathbf{x} = \mathbf{x}(v) - \mathbf{x}(u)$  for  $e = (u, v)$ . We can rewrite this as a proper linear program by introducing extra variables for the maximum:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^V, \mathbf{u} \in \mathbb{R}^E} \quad & \mathbf{c}^\top \mathbf{u} \\ \text{s.t.} \quad & \mathbf{b}_{s,t}^\top \mathbf{x} = 1 \\ & \mathbf{0} \leq \mathbf{x} \leq \mathbf{1} \\ & \mathbf{u} \geq \mathbf{0} \\ & \mathbf{u} \geq \mathbf{B}^\top \mathbf{x} \end{aligned} \tag{12.4}$$

And, in fact, we'll also see that we don't need the constraint  $\mathbf{0} \leq \mathbf{x} \leq \mathbf{1}$ , leading to the

following simpler program:

$$\begin{aligned}
& \min_{x \in \mathbb{R}^V, u \in \mathbb{R}^E} \quad c^\top u \\
& \text{s.t.} \quad b_{s,t}^\top x = 1 \\
& \quad u \geq 0 \\
& \quad u \geq B^\top x
\end{aligned} \tag{12.5}$$

This last linear program is the dual program to the maximum flow linear program.

**Lemma 12.5.1.** *Programs (12.2), (12.4), and (12.5) have equal optimal values.*

*Proof.* We start by considering equality between optimal values of Program (12.2) and Program (12.4). Let  $S$  be an optimal solution to Program (12.2) and take  $x = \mathbf{1}_{V \setminus S}$ . Then  $x$  is a feasible solution to Program (12.4) with cost equal to  $c_G(S)$ .

Conversely, suppose  $x$  is an optimal solution to Program (12.4). Let  $\alpha$  be uniform in  $[0, 1]$  and define  $S = \{v \in V \mid x(v) \leq \alpha\}$ . This is called a threshold cut. We can verify that

$$\Pr[e \text{ is cut by } S] = \max\{b_e^\top x, 0\}.$$

and hence  $\mathbb{E}_t[c_G(S)]$  is exactly the optimization function of 12.3 (and hence 12.4). Since at least one outcome must do as well as the average, there is a subset  $S \subseteq V$  achieving this value (or less).

We can use the same threshold cut to show that we can round a solution to Program (12.5) to an equal or smaller value cut feasible for Program (12.2). The only difference is that in this case, we get

$$\Pr[e \text{ is cut by } S] \leq \max\{b_e^\top x, 0\}$$

which is still sufficient.  $\square$

# Chapter 13

## The Cut-Matching Game: Expanders via Max Flow

In this chapter, we learn about a new algorithm to compute expanders that employs max flow as a subroutine.

### 13.1 Introduction

We start with a review of expanders where make a subtle change the notion of an expander in comparison to chapter 7 to ease the exposition.

**Definitions.** We let  $G = (V, E)$  be an unweighted, connected graph in this chapter, and let  $\mathbf{d}$  be the degree vector,  $E(S, V \setminus S)$  denote the set of edges crossing the cut  $A, B$

Given set  $\emptyset \subset S \subset V$ , then we define the **sparsity**  $\psi(S)$  of  $S$  by

$$\psi(S) = \frac{|E(S, V \setminus S)|}{\min\{|S|, |V \setminus S|\}}$$

Note that **sparsity**  $\psi(S)$  differs from **conductance**  $\phi(S) = \frac{|E(S, V \setminus S)|}{\min\{\text{vol}(S), \text{vol}(V \setminus S)\}}$ , as defined in 7, in the denominator. It is straight-forward to see that in a connected graph  $\psi(S) \geq \phi(S)$  for all  $S$ .

Clearly, we again have  $\psi(S) = \psi(V \setminus S)$ . We define the *sparsity* of a graph  $G$  by  $\psi(G) = \min_{\emptyset \subset S \subset V} \psi(S)$ . For any  $\psi \in (0, n]$ , we say a graph  $G$  is a  $\psi$ -expander with regard to sparsity, if  $\psi(G) \geq \psi$ . When the context is clear, we simply say that  $G$  is a  $\psi$ -expander.

**The Main Result.** The main result of this chapter is the following theorem.

**Theorem 13.1.1.** *There is an algorithm  $\text{SPARSITYCERTIFYORCUT}(G, \psi)$  that given a graph  $G$  and a parameter  $0 < \psi \leq 1$  either:*

- Certifies that  $G$  is a  $\Omega(\psi / \log^2 n)$ -expander with regard to sparsity, or
- Presents a cut  $S$  such that  $\psi(S) \leq O(\psi)$ .

The algorithm runs in time  $O(\log^2 n) \cdot T_{\text{max\_flow}}(G) + \tilde{O}(m)$  where  $T_{\text{max\_flow}}(G)$  is the time it takes to solve a Max Flow problem on  $G$ <sup>1</sup>.

The bounds above can further be extended to compute  $\phi$ -expanders (with regard to conductance). Using current state-of-the-art Max Flow results, the above problem can be solved in  $m^{1+o(1)}$  time (see [?]).

## 13.2 Embedding Graphs into Expanders

Let us start by exploring the first key idea behind the algorithm. We therefore need a definition of what it means to embed one graph into another.

**Definition of Embedding.** Given graphs  $H$  and  $G$  that are defined over the same vertex set, then we say that a function  $\text{EMBED}_{H \rightarrow G}$  is an *embedding* if it maps each edge  $(u, v) \in H$  to a  $u$ -to- $v$  path  $P_{u,v} = \text{EMBED}_{H \rightarrow G}(u, v)$  in  $G$ .

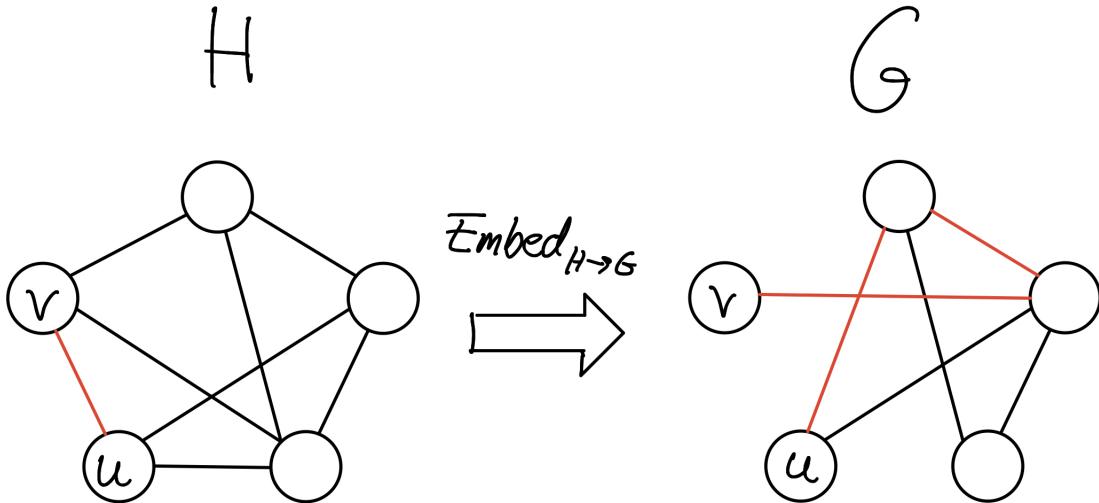


Figure 13.1: In this example the red edge  $(u, v)$  in  $H$  is mapped to the red  $u$ -to- $v$  path in  $G$ .

---

<sup>1</sup>Technically, we will solve problems with two additional vertices and  $n$  additional edges but this will not change the run-time of any known max-flow algorithm asymptotically.

We say that the *congestion* of  $\text{EMBED}_{H \rightarrow G}$  is the maximum number of times that any edge  $e \in E(G)$  appears on any embedding path:

$$\text{cong}(\text{EMBED}_{H \rightarrow G}) = \max_{e \in E(G)} |\{e' \in E(H) \mid e \in \text{EMBED}_{H \rightarrow G}(e')\}|.$$

**Certifying Expander Graphs via Embeddings.** Let us next prove the following lemma that is often consider Folklore.

**Lemma 13.2.1.** *Given a  $\frac{1}{2}$ -expander graph  $H$  and an embedding of  $H$  into  $G$  with congestion  $C$ , then  $G$  must be an  $\Omega(\frac{1}{C})$ -expander.*

*Proof.* Consider any cut  $(S, V \setminus S)$  with  $|S| \leq |V \setminus S|$ . Since  $H$  is a  $\psi$ -expander, we have that  $|E_H(S, V \setminus S)| \geq |S|/2$ . We also know by the embedding of  $H$  into  $G$ , that for each edge  $(u, v) \in E_H(S, V \setminus S)$ , we can find path a  $P_{u,v}$  in  $G$  that also has to cross the cut  $(S, V \setminus S)$  at least once. But since each edge in  $G$  is on at most  $C$  such paths, we can conclude that at least  $|E_H(S, V \setminus S)|/C \geq |S|/2C$  edges in  $G$  cross the cut  $(S, V \setminus S)$ .  $\square$

Unfortunately, the reverse of the above lemma is not true, i.e. even if there exists no embedding from  $1/2$ -expander  $H$  into  $G$  of congestion  $C$ , then  $G$  might still be an  $\Omega(\frac{1}{C})$ -expander.

### 13.3 The Cut-Matching Algorithm

Although there are still some missing pieces, let us next discuss the algorithm 9. The algorithm runs for  $T$  iterations where we will later find that the right value to set  $T$  to is in  $\Theta(\log^2 n)$ .

---

**Algorithm 9: SPARSITYCERTIFYORCUT( $G, \psi$ )**


---

```

1 for  $i = 1, 2, \dots, T$  do
2    $(S_i, \bar{S}_i) \leftarrow \text{FINDBiPARTITION}(G, \{M_1, M_2, \dots, M_i\})$ ; // Assume  $|S| = |\bar{S}| = n/2$ 
3   Solve the flow problem on  $G$  where each vertex  $e \in E$  receives capacity  $\mathbf{c}(e) = 1/\psi$ 
   and the demand at each edge  $v \in S$  is  $+1$  and for each  $v \in \bar{S}$  is  $-1$ , by introducing
   a super-source  $s$  and super-sink  $t$ ;
4   if flow procedure returns a flow  $\mathbf{f}$  with  $\text{val}(\mathbf{f}) = n/2$  then
5     Remove  $s, t$  to derive  $S-\bar{S}$  flow; then decompose flow  $\mathbf{f}$  into flow paths
       $P_1, P_2, \dots, P_{n/2}$ ;
6     Create a matching  $M_i$  where for each  $s$ -to- $\bar{s}$  flow path  $P_i$ , we add  $(s, \bar{s})$  to  $M_i$ .
7   else
8     Let  $(X \cup \{s\}, (V \setminus X) \cup \{t\})$  be the  $st$ -minimum cut in the above flow problem.
9     return  $(X, V \setminus X)$ 
10 return  $H = \bigcup_i M_i$ 

```

---

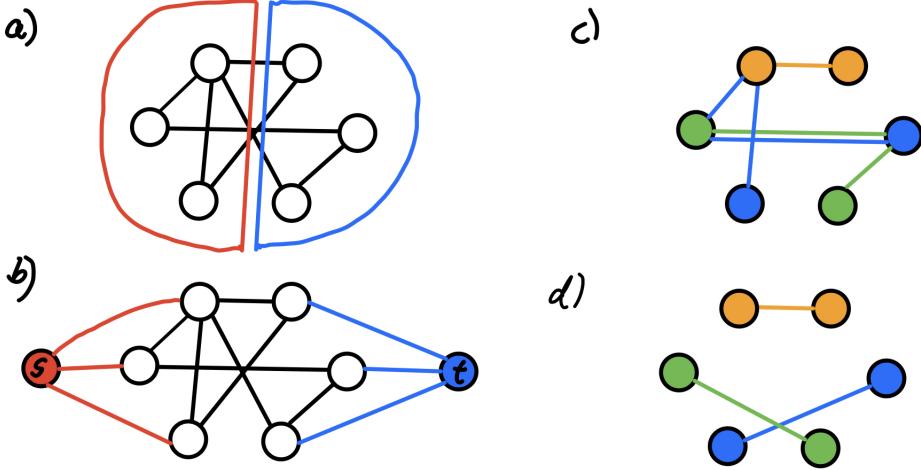


Figure 13.2: Illustration of the steps of the Algorithm. In a), a bi-partition of  $V$  is found. In b), the bi-partition is used to obtain a flow problem where we inject one unit of flow to each vertex in  $S$  via super-source  $s$  and extract one unit via super-sink  $t$ . c) A path flow decomposition. For each path, the first vertex is in  $S$  and the last vertex in  $\bar{S}$ . d) We find  $M_i$  to be the one-to-one matching between endpoints in  $S$  and  $\bar{S}$  defined by the path flows.

**Beginning of an Iteration.** In each iteration of the algorithm, we first invoke a sub-procedure  $\text{FINDBiPARTITION}(\cdot)$  that returns a cut  $(S, \bar{S})$  (where  $\bar{S} = V \setminus S$ ) that partitions the vertex set into two equal-sized sides. Here, we implicitly assume that the number of vertices  $n$  is even which is w.l.o.g.

Next, the algorithm creates a flow problem where each vertex  $s \in S$  has to send exactly one unit of flow and each vertex  $\bar{s} \in \bar{S}$  has to receive exactly one unit of flow. We therefore introduce dummy nodes  $s$  and  $t$ , add for each vertex  $v \in S$  an edge of capacity 1 between  $s$  and  $v$ ; and for each vertex  $v \in \bar{S}$  and edge between  $t$  and  $v$  of capacity 1. We set the capacity of edges in the original graph  $G$  to  $1/\psi$  (which we assume wlog to be integer).

**The If Statement.** If the flow problem can be solved exactly, then we can find a path decomposition of the flow  $\mathbf{f}$  in  $\tilde{O}(m)$  time (for example using a DFS) where each path starts in  $S$  ends in  $\bar{S}$  and carries one unit of flow<sup>2</sup>. This defines a one-to-one correspondence between the vertices in  $S$  and the vertices in  $\bar{S}$ . We capture this correspondences in the matching  $M_i$ . We will later prove the following lemma.

**Lemma 13.3.1.** *If the algorithm returns after constructing  $T$  matchings, for an appropriately chosen  $T = \Theta(\log^2 n)$ , then the graph  $H$  returned by the algorithm is a  $\frac{1}{2}$ -expander and  $H$  can be embedded into  $G$  with congestion  $O(\log^2 n / \psi)$ .*

<sup>2</sup>For simplicity, assume that the returned flow is integral.

**The Else Statement.** On the other hand, if the flow problem on  $G$  could not be solved, then we return the min-cut of the flow problem. Such a cut can be found in  $O(m)$  time by using the above reduction to an  $s$ - $t$  flow problem on which one can compute maximum flow  $f$  from which the  $s$ - $t$  min-cut can be constructed by following the construction in the proof of Theorem 11.4.4.

It turns out that this min-cut already is a sparse cut by the way our flow problem is defined.

**Lemma 13.3.2.** *If the algorithm terminates with a cut  $(X, V \setminus X)$ , then  $\psi(X) = O(\psi)$ .*

*Proof.* Let  $X_s = X \cup \{s\}$ ,  $\bar{X}_s = (V \setminus X) \cup \{t\}$ . First observe that since the flow was not routed, and  $(X_s, \bar{X}_s)$  is an  $st$ -min cut for this flow problem, we have that  $c(E_G(X_s, \bar{X}_s)) < n/2$  (otherwise we could have routed the demands).

Let  $n_s$  be the number of edges incident to the super-source  $s$  that cross the cut  $(X_s, \bar{X}_s)$ . Let  $n_t$  be the number of edges incident to  $t$  that cross the cut  $(X_s, \bar{X}_s)$ .

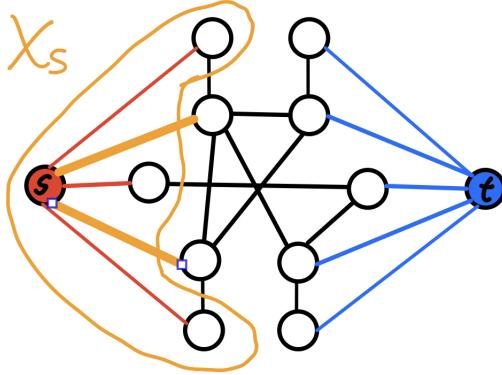


Figure 13.3: Set  $X_s$  is enclosed by the orange circle. The thick orange edges are in the cut and incident to super-source  $s$ . Thus they count towards  $n_s$ . Here  $n_s = 2, n_t = 0$ . Note that all remaining edges in the cut are black, i.e. were originally in  $G$  and therefore have capacity  $1/\psi$ .

Observe that after taking away the vertices  $s$  and  $t$ , the cut  $(X, V \setminus X)$  has less than  $n/2 - n_s - n_t$  capacity. But each remaining edge has capacity  $1/\psi$ , so the total number of edges in the cut can be at most  $\psi \cdot (n/2 - n_s - n_t)$ . Since  $X = X_s \setminus \{s\}$  is of size at least  $n/2 - n_s$ , and  $V \setminus X = \bar{X}_s \setminus \{t\}$  is of size at least  $n/2 - n_t$ , we have that the induced cut in  $G$  has

$$\psi(X) < \frac{\psi \cdot (n/2 - n_s - n_t)}{\min\{n/2 - n_s, n/2 - n_t\}} \leq \psi.$$

□

## 13.4 Constructing an Expander via Random Walks

Next, we give the implementation and analysis for the procedure  $\text{FINDBiPARTITION}(\cdot)$ . We start however by giving some more preliminaries.

**Random Walk on Matchings.** Let  $\{M_1, M_2, \dots, M_T\}$  be the set of matchings we compute (if we never find a cut). In the  $i^{th}$ -step of the lazy random walk, we let the mass at each vertex  $j$  stay put with probability  $1/2$ , and otherwise traverses the edge in matching  $M_i$  incident to  $j$  with probability  $1/2$ .

We let  $\mathbf{p}_{j \rightarrow i}^t$  denote the probability that a particle that started at vertex  $j$  is at vertex  $i$  after a  $t$ -step lazy random walk. We let  $\mathbf{p}_i^t = [\mathbf{p}_{1 \rightarrow i}^t \ \mathbf{p}_{2 \rightarrow i}^t \ \dots \ \mathbf{p}_{n \rightarrow i}^t]$ . Note that for each edge  $(i, j) \in M_{t+1}$ , we have that

$$\mathbf{p}_i^{(t+1)} = \frac{1}{2}\mathbf{p}_i^t + \frac{1}{2}\mathbf{p}_j^t = \mathbf{p}_j^{(t+1)}.$$

We define the projection matrix  $\mathbf{\Pi}^t = [\mathbf{p}_1^t, \mathbf{p}_2^t, \dots, \mathbf{p}_n^t]^\top$  that maps an initial probability distribution  $\mathbf{d}$  to the probability distribution over the vertices that the random walk visits them at step  $t$ . You will prove in the exercises that  $\mathbf{\Pi}^t$  is *doubly-stochastic*.

We say that a lazy random walk is *mixing* at step  $t$ , if for each  $i, j$ ,  $\mathbf{p}_{j \rightarrow i}^t \geq 1/(2n)$ .

**Lemma 13.4.1.** *If  $t$ -step lazy random walk is mixing, then  $H = \bigcup_{i \leq t} M_i$  is a  $\frac{1}{2}$ -expander.*

*Proof.* Consider any cut  $(S, \bar{S})$  with  $|S| \leq |\bar{S}|$ . It is convenient to think about the random walks in terms of probability mass that is moved around. Observe that each vertex  $j \in \bar{S}$  has to push at least  $1/(2n)$  units of mass from  $j$  to  $i$  (by definition of mixing).

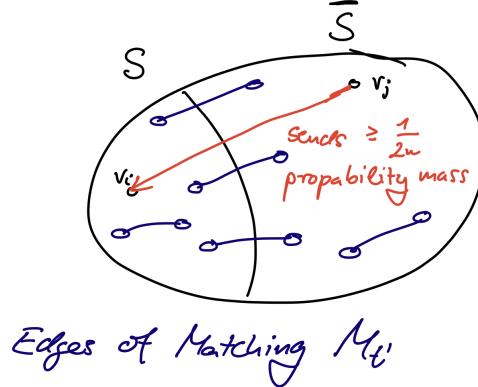


Figure 13.4: Each vertex  $j \in \bar{S}$  sends at least  $1/(2n)$  probability mass to  $i$  (red arrow). But in order to transport it, it does have to push the mass through edges in the matchings  $M_1, M_2, \dots, M_t$  that cross the cut.

Clearly, to move the mass from  $\bar{S}$  to  $S$  it has to use the matching edges that also cross the cut.

Now observe that since there are  $\geq n/2$  vertices in  $\bar{S}$ , and each of them has to push  $\geq 1/(2n)$  mass to  $i$ , the total amount of probability mass pushed through the cut for  $i$  is  $\geq 1/4$ . Since there are  $|S|$  such vertices  $i$ , the total amount of mass that has to cross the cut is  $\geq |S|/4$ . But note, that after each step of the random walk, the total probability mass at each vertex is exactly 1. Thus, at each step  $i \leq t$ , each edge in  $M_i$  crossing the cut can push at most  $1/2$  units of probability mass over the cut (and thereafter the edge is gone).

It follows that there must be at least  $|S|/2$  edges in the matchings  $M_1, M_2, \dots, M_t$ . But this implies that  $H = \cup_i M_i$  is a  $\frac{1}{2}$ -expander.  $\square$

**Implementing `FindBiPartition`( $\cdot$ ).** We can now give away the implementation of `FINDBiPARTITION`( $\cdot$ ) which you can find below.

---

**Algorithm 10:** `FINDBiPARTITION`( $G, \{M_1, M_2, \dots, M_t\}$ )

---

- 1 Choose random  $n$ -dimensional vector  $\mathbf{r}$  orthogonal to  $\mathbf{1}$ ;
  - 2 Compute vector  $\mathbf{u} = \mathbf{\Pi}^t \mathbf{r}$ , i.e. each  $\mathbf{u}(i) = \mathbf{p}_i^t \cdot \mathbf{r}$ ;
  - 3 Let  $S$  be the  $n/2$  smallest vertices w.r.t.  $\mathbf{u}$ ; and  $\bar{S}$  be the  $n/2$  largest w.r.t  $\mathbf{u}$  (ties broken arbitrarily but consistently);;
  - 4 **return**  $(S, \bar{S})$
- 

The central claim, we want to prove is the following: given a potential function for the random walk at step  $t$

$$\Phi^t = \sum_{i,j} (\mathbf{p}_{j \rightarrow i}^t - 1/n)^2 = \sum_i \|\mathbf{p}_i^t - \mathbf{1}/n\|_2^2.$$

**Claim 13.4.2.** *In the algorithm `SPARSITYCERTIFYORCUT`( $\cdot$ ), we have  $\mathbb{E}[\Phi^t - \Phi^{(t+1)}] = \Omega(\Phi^t / \log n) - O(1/n^{-5})$ . Further, we have that  $\Phi^t - \Phi^{(t+1)}$  is always non-negative. The expectation is over the random vector  $\mathbf{r}$  chosen in the current round.*

**Corollary 13.4.3.** *For appropriate  $T = \Theta(\log^2 n)$ , the algorithm `SPARSITYCERTIFYORCUT`( $\cdot$ ) has  $\Phi^{(T+1)} \leq 4/n^2$  w.h.p. (i.e. with probability  $\geq 1 - 1/n$ ).*

To obtain the Corollary, one can simply set up a sequence of random 0-1 variables  $X^1, X^2, \dots, X^{(T)}$  where each  $X^{(t+1)}$  is 1 if and only if the decrease in potential is at least an  $\Omega(1/\log n)$ -fraction of the previous potential. Since the expectation is only over the current  $\mathbf{r}$  in each round and we choose these independently at random, one can then use a Chernoff bound to argue that after  $T$  rounds (for an appropriate hidden constant), one has at least  $\Omega(T)$  rounds during which the potential is decreased substantially (unless it is already tiny and the  $O(n^{-5})$  factor dominates).

We further observe that this implies that  $\{M_1, M_2, \dots, M_T\}$  is mixing (otherwise, the contribution of the violating coordinate  $\mathbf{p}_{j \rightarrow i}^T$  would contribute  $> 1/(2n)^2$  to the potential, but all terms are non-negative), and thereby we conclude the proof of our main theorem.

Let us now give the prove of Claim 13.4.2:

**Interpreting the Potential Drop.** Let us start by writing out the amount by which the potential decreases

$$\Phi^t - \Phi^{(t+1)} = \sum_i \|\mathbf{p}_i^t - \mathbf{1}/n\|_2^2 - \sum_i \|\mathbf{p}_i^{(t+1)} - \mathbf{1}/n\|_2^2$$

Considering now matching  $M_{t+1}$ , and an edge  $(i, j) \in M_{t+1}$ . We can re-write the former sum as  $\sum_i \|\mathbf{p}_i^t - \mathbf{1}/n\|_2^2 = \sum_{(i,j) \in M_{t+1}} \|\mathbf{p}_i^t - \mathbf{1}/n\|_2^2 + \|\mathbf{p}_j^t - \mathbf{1}/n\|_2^2$  as each vertex occurs as exactly one endpoint of a matching edge. We can do the same for the  $t+1$ -step walk probabilities. Further, recall that for  $(i, j) \in M_{t+1}$ , we have  $\mathbf{p}_i^{(t+1)} = \mathbf{p}_j^{(t+1)} = \frac{\mathbf{p}_i^t + \mathbf{p}_j^t}{2}$ . Thus,

$$\begin{aligned} \Phi^t - \Phi^{(t+1)} &= \sum_{(i,j) \in M_{t+1}} \|\mathbf{p}_i^t - \mathbf{1}/n\|_2^2 + \|\mathbf{p}_j^t - \mathbf{1}/n\|_2^2 - \|\mathbf{p}_i^{(t+1)} - \mathbf{1}/n\|_2^2 - \|\mathbf{p}_j^{(t+1)} - \mathbf{1}/n\|_2^2 \\ &= \sum_{(i,j) \in M_{t+1}} \|\mathbf{p}_i^t - \mathbf{1}/n\|_2^2 + \|\mathbf{p}_j^t - \mathbf{1}/n\|_2^2 - 2 \left\| \frac{\mathbf{p}_i^t + \mathbf{p}_j^t}{2} - \mathbf{1}/n \right\|_2^2. \end{aligned}$$

Finally, we can use the formula  $\|\mathbf{x}\|_2^2 + \|\mathbf{y}\|_2^2 - 2\|(\mathbf{x} + \mathbf{y})/2\|_2^2 = \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|_2^2$  term-wise to derive

$$\Phi^t - \Phi^{(t+1)} = \frac{1}{2} \sum_{(i,j) \in M_{t+1}} \|(\mathbf{p}_i^t - \mathbf{1}/n) - (\mathbf{p}_j^t - \mathbf{1}/n)\|_2^2 = \frac{1}{2} \sum_{(i,j) \in M_{t+1}} \|\mathbf{p}_i^t - \mathbf{p}_j^t\|_2^2.$$

The potential thus drops by a lot if vertices  $i$  and  $j$  are matched where  $\mathbf{p}_i^t$  and  $\mathbf{p}_j^t$  differ starkly. Note that this equality implies directly the remark in our claim that  $\Phi^t - \Phi^{(t+1)}$  is non-negative.

**Understanding the Random Projection.** Next, we want to further lower bound the potential drop using the random vector  $\mathbf{u}$ . This intuitively helps a lot in our analysis since we are matching vertices  $i, j$  with high value  $\mathbf{u}(i)$  and low value  $\mathbf{u}(j)$  (or vice versa). We will show that (w.p.  $\geq 1 - n^{-3}$ )

$$\Phi^t - \Phi^{(t+1)} = \frac{1}{2} \sum_{(i,j) \in M_{t+1}} \|\mathbf{p}_i^t - \mathbf{p}_j^t\|_2^2 \geq \frac{n-1}{64 \cdot \log n} \sum_{(i,j) \in M_{t+1}} |\mathbf{u}(i) - \mathbf{u}(j)|^2. \quad (13.1)$$

We prove this claim again term-wise, showing that for each pair of vertices  $i, j \in V$ , we have  $\|\mathbf{p}_i^t - \mathbf{p}_j^t\|_2^2 \geq \frac{n-1}{16 \cdot \log n} |\mathbf{u}(i) - \mathbf{u}(j)|^2$  w.h.p. It will then suffice to talk a union bound over all pairs  $i, j$ .

To this end, let us make the following observations: since  $\mathbf{u}(i) = \mathbf{p}_i^t \cdot \mathbf{r}$ , we have that  $\mathbf{u}(i) - \mathbf{u}(j) = (\mathbf{p}_i^t - \mathbf{p}_j^t) \cdot \mathbf{r}$  by linearity. Also note that since  $\sum_j \mathbf{p}_{j \rightarrow i}^t = 1$  for all  $i$  (since  $\Pi^t$  is doubly-stochastic), we further have that the projection  $(\mathbf{p}_i^t - \mathbf{p}_j^t)$  is orthogonal to  $\mathbf{1}$ .

We can now use the following statement about random vector  $\mathbf{r}$  to argue about the effect of projecting  $(\mathbf{p}_i^t - \mathbf{p}_j^t)$  onto  $\mathbf{r}$ . Below, we note that we have  $d = n - 1$  since  $\mathbf{r}$  is chosen from the  $(n - 1)$ -dimensional space orthogonal to  $\mathbf{1}$ .

**Theorem 13.4.4** (see for example [?] or any lecture notes on Johnson-Lindenstrauss). *If  $\mathbf{y}$  is a vector of length  $\ell$  in  $\mathbb{R}^d$ , and  $\mathbf{r}$  a unit random vector in  $\mathbb{R}^d$ , then*

- $\mathbb{E}[(\mathbf{y}^\top \mathbf{r})^2] = \frac{\ell^2}{d}$ , and
- for  $x \leq d/16$ , then  $\mathbb{P}[(\mathbf{y}^\top \mathbf{r})^2 \geq x\ell^2/d] \leq e^{-x/4}$

This allows us to pick  $x = 32 \cdot \log n$ , and we then obtain that

$$\mathbb{P}\left[((\mathbf{p}_i^t - \mathbf{p}_j^t) \cdot \mathbf{r})^2 \geq \frac{32 \log n}{n-1} \|\mathbf{p}_i^t - \mathbf{p}_j^t\|_2^2\right] \leq e^{-4 \log n} = n^{-8}. \quad (13.2)$$

Multiplying both sides of the event by  $(n-1)/(64 \log n)$ , we derive the claimed Inequality (13.1). We can further union bound over the  $n/2$  matching pairs to all satisfy this bound with probability  $\geq 1 - n^{-7}$ , as desired.

**Relating to the Lengths of the Projections.** Let  $\mu = \max_{i \in S} \mathbf{u}(i)$ , then we have by definition that  $\mathbf{u}(i) \leq \mu \leq \mathbf{u}(j)$  for all  $i \in S, j \in \bar{S}$ .

Now we can write

$$\begin{aligned} \Phi^t - \Phi^{(t+1)} &\geq \frac{n-1}{64 \cdot \log n} \sum_{(i,j) \in M_{t+1}} |\mathbf{u}(i) - \mathbf{u}(j)|^2 \\ &\geq \frac{n-1}{64 \cdot \log n} \sum_{(i,j) \in M_{t+1}} (\mathbf{u}(i) - \mu)^2 + (\mathbf{u}(j) - \mu)^2 \\ &= \frac{n-1}{64 \cdot \log n} \sum_{i \in V} (\mathbf{u}(i) - \mu)^2 \\ &= \frac{n-1}{64 \cdot \log n} \left( \sum_{i \in V} \mathbf{u}(i)^2 - 2\mu \cdot \sum_{i \in V} \mathbf{u}(i) + n\mu^2 \right) \end{aligned}$$

by standard calculations. We then observe that  $\sum_i \mathbf{u}(i) = \sum_i \mathbf{p}_i^t \cdot \mathbf{r} = \mathbf{1} \cdot \mathbf{r} = 0$  by the fact that  $\Pi^t$  is doubly-stochastic and since  $\mathbf{r}$  is orthogonal to the all-ones vector. We can therefore conclude

$$\frac{n-1}{64 \cdot \log n} \left( \sum_{i \in V} \mathbf{u}(i)^2 - 2\mu \cdot \sum_i \mathbf{u}(i) + n\mu^2 \right) \geq \frac{n-1}{64 \cdot \log n} \sum_{i \in V} \mathbf{u}(i)^2. \quad (13.3)$$

**Taking the Expectation.** Then, from the second fact of Theorem 13.4.4, we obtain that

$$\mathbb{E} \left[ \sum_{i \in V} \mathbf{u}(i)^2 \right] = \sum_{i \in V} \mathbb{E}[\mathbf{u}(i)^2] = \sum_{i \in V} \mathbb{E}[(\mathbf{p}_i^t \cdot \mathbf{r})^2] = \sum_{i \in V} \mathbb{E}[(\mathbf{p}_i^t - \mathbf{1}/n) \cdot \mathbf{r})^2] \quad (13.4)$$

$$= \sum_{i \in V} \frac{\|\mathbf{p}_i^t - \mathbf{1}/n\|_2^2}{n-1} = \frac{\Phi^t}{n-1} \quad (13.5)$$

where we used again that  $\mathbf{r}$  is orthogonal to  $\mathbf{1}$ .

Unfortunately, we cannot directly use this expectation since we already conditioned on the high probability events in Equation (13.2). But a simple trick allows us to recover: Let  $\mathcal{E}$  denote the union of all of these events. We have by the law of total expectation

$$\mathbb{E} \left[ \sum_{i \in V} \mathbf{u}(i)^2 \right] = \mathbb{P}[\neg \mathcal{E}] \cdot \mathbb{E} \left[ \sum_{i \in V} \mathbf{u}(i)^2 | \neg \mathcal{E} \right] + \mathbb{P}[\mathcal{E}] \cdot \mathbb{E} \left[ \sum_{i \in V} \mathbf{u}(i)^2 | \mathcal{E} \right]$$

But note that  $\mathbb{E} [\sum_{i \in V} \mathbf{u}(i)^2 | \mathcal{E}]$  has to be smaller than  $n$  because  $\sum_{i \in V} \mathbf{u}(i)^2 = \sum_{i \in V} (\mathbf{p}_i^t \cdot \mathbf{r})^2 \leq n$  with probability 1 (because each  $\mathbf{p}_i^t$  is a unit vector). Recall that we calculated  $\mathbb{P}[\mathcal{E}] \leq n^{-7}$ . Thus, we can conclude that  $\mathbb{E} [\sum_{i \in V} \mathbf{u}(i)^2 | \neg \mathcal{E}] \geq \frac{\Phi^t}{n-1} - n^{-6}$ .

It remains to combine our insights to conclude

$$\mathbb{E}[\Phi^t - \Phi^{(t+1)} | \neg \mathcal{E}] \geq \frac{n-1}{64 \cdot \log n} \left( \frac{\Phi^t}{n-1} - n^{-6} \right) = \Omega(\Phi^t / \log n) - O(1/n^{-5}).$$

Since again the event  $\mathcal{E}$  occurs with very low probability, and since  $\Phi^t - \Phi^{t+1}$  is non-negative always, we can then conclude that in unconditionally, in expectation, the potential decreases by  $\Omega(\Phi^t / \log n) - O(1/n^{-5})$ .

# **Part IV**

## **Further Topics in Convex Optimization**

# Chapter 14

## Separating Hyperplanes, Lagrange Multipliers, KKT Conditions, and Convex Duality

### 14.1 Overview

First part of this chapter introduces the concept of a separating hyperplane of two sets followed by a proof that for two closed, convex and disjoint sets a separating hyperplane always exists. This is a variant of the more general *separating hyperplane theorem*<sup>1</sup> due to Minkowski. Then *Lagrange multipliers*  $\mathbf{x}, \mathbf{s}$  of a convex optimization problem

$$\begin{aligned} \min_{\mathbf{y}} \quad & \mathcal{E}(\mathbf{y}) \\ \text{s.t.} \quad & \mathbf{A}\mathbf{y} = \mathbf{b} \\ & c(\mathbf{y}) \leq 0 \end{aligned}$$

are introduced and with that, the *Lagrangian*

$$L(\mathbf{y}, \mathbf{x}, \mathbf{s}) = \mathcal{E}(\mathbf{y}) + \mathbf{x}^\top (\mathbf{b} - \mathbf{A}\mathbf{y}) + \mathbf{s}^\top c(\mathbf{y})$$

is defined. Finally, we deal with the dual problem

$$\max_{\mathbf{x}, \mathbf{s}, \mathbf{s} \geq 0} L(\mathbf{x}, \mathbf{s}),$$

where  $L(\mathbf{x}, \mathbf{s}) = \min_{\mathbf{y}} L(\mathbf{y}, \mathbf{x}, \mathbf{s})$ . We show *weak duality*, i.e.  $L(\mathbf{y}, \mathbf{x}, \mathbf{s}) \leq \mathcal{E}(\mathbf{y})$  and that assuming *Slater's condition* the values of both the primal and dual is equal, which is referred to as *strong duality*.

---

<sup>1</sup>Wikipedia is good on this: [https://en.wikipedia.org/wiki/Hyperplane\\_separation\\_theorem](https://en.wikipedia.org/wiki/Hyperplane_separation_theorem)

## 14.2 Separating Hyperplane Theorem

Suppose we have two convex subsets  $A, B \subseteq \mathbb{R}^n$  that are disjoint ( $A \cup B = \emptyset$ ). We wish to show that there will always be a (hyper-)plane  $H$  that separates these two sets, i.e.  $A$  lies on one side, and  $B$  on the other side of  $H$ .

So what exactly do we mean by Hyperplane? Let's define it.

**Definition 14.2.1** (Hyperplane). A *hyperplane*  $H$  of dimension  $n$  is the subset  $H := \{\mathbf{x} \in \mathbb{R}^n : \langle \mathbf{n}, \mathbf{x} \rangle = \mu\}$ . We say  $H$  has *normal*  $\mathbf{n} \in \mathbb{R}^n$  and *threshold*  $\mu$ . It is required that  $\mathbf{n} \neq \mathbf{0}$ .

Every hyperplane divides  $\mathbb{R}^n$  into two halfspaces  $\{\mathbf{x} : \langle \mathbf{v}, \mathbf{x} \rangle \geq \mu\}$  and  $\{\mathbf{x} : \langle \mathbf{v}, \mathbf{x} \rangle \leq \mu\}$ . It separates two sets, if they lie in different halfspaces. We formally define separating hyperplane as follows.

**Definition 14.2.2** (Separating Hyperplane). We say a hyperplane  $H$  *separates* two sets  $A, B$  iff

$$\begin{aligned}\forall \mathbf{a} \in A : \langle \mathbf{n}, \mathbf{a} \rangle &\geq \mu \\ \forall \mathbf{b} \in B : \langle \mathbf{n}, \mathbf{b} \rangle &\leq \mu\end{aligned}$$

If we replace  $\geq$  with  $>$  and  $\leq$  with  $<$  we say  $H$  *strictly* separates  $A$  and  $B$ .

It is easy to see that there exists disjoint non-convex sets that can not be separated by a hyperplane (e.g. a point cannot be separated from a ring around it). But can two disjoint convex sets always be strictly separated by a hyperplane? The answer is no: consider the two-dimensional case depicted in Figure 14.1 with  $A = \{(x, y) : x \leq 0\}$  and  $B = \{(x, y) : x > 0 \text{ and } y \geq \frac{1}{x}\}$ . Clearly they are disjoint; however the only separating hyperplane is  $H = \{(x, y) : x = 0\}$  but it intersects  $A$ .

One can prove that there exists a non-strictly separating hyperplane for any two disjoint convex sets. We will prove that if we further require  $A, B$  to be closed and bounded, then a strictly separating hyperplane always exists. (Note in the example above how our choice of  $B$  is not bounded.)

**Theorem 14.2.3** (Separating Hyperplane Theorem; closed, bounded sets). *For two closed, bounded, and disjoint convex sets  $A, B \in \mathbb{R}^n$  there exists a strictly separating hyperplane  $H$ . One such hyperplane is given by normal  $\mathbf{n} = \mathbf{d} - \mathbf{c}$  and threshold  $\mu = \frac{1}{2} (\|\mathbf{d}\|_2^2 - \|\mathbf{c}\|_2^2)$ , where  $\mathbf{c} \in A$ ,  $\mathbf{d} \in B$  are the minimizers of the distance between  $A$  and  $B$*

$$dist(A, B) = \min_{\mathbf{a} \in A, \mathbf{b} \in B} \|\mathbf{a} - \mathbf{b}\|_2 > 0.$$

*Proof.* We omit the proof that  $dist(A, B) = \min_{\mathbf{a} \in A, \mathbf{b} \in B} \|\mathbf{a} - \mathbf{b}\|_2 > 0$ , which follows from  $A, B$  being disjoint, closed, and bounded. Now, we want to show that  $\langle \mathbf{n}, \mathbf{b} \rangle > \mu$  for all  $\mathbf{b} \in B$ ; then  $\langle \mathbf{n}, \mathbf{a} \rangle < \mu$  for all  $\mathbf{a} \in A$  follows by symmetry. Observe that

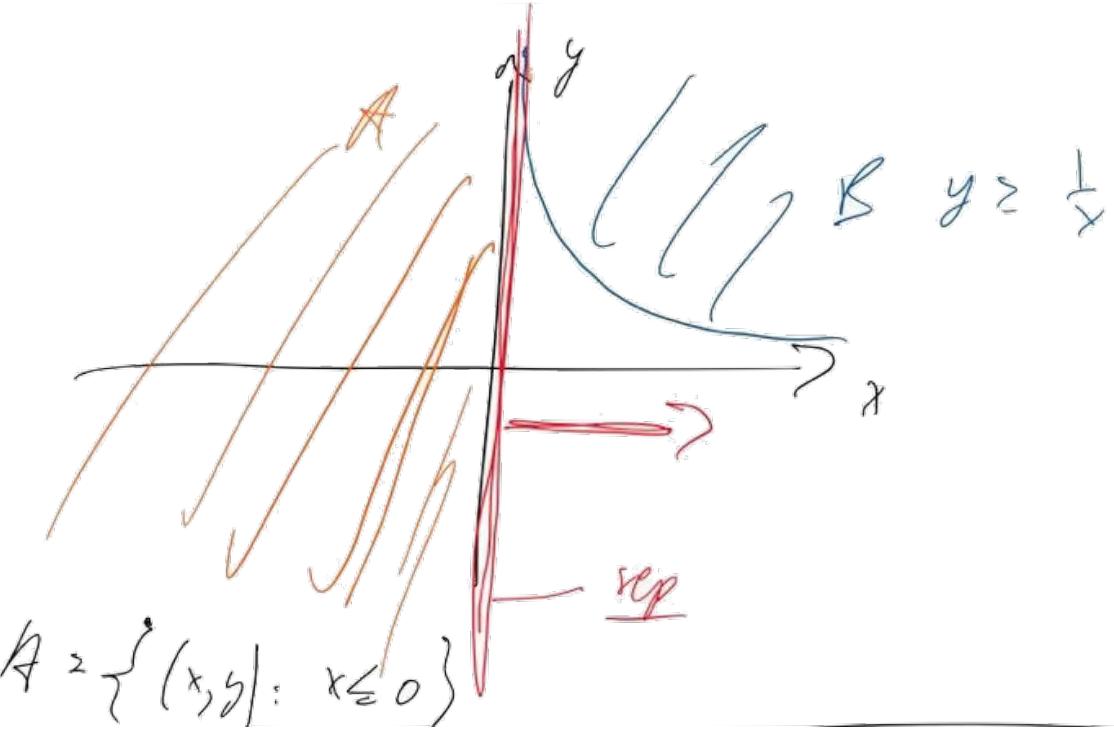


Figure 14.1: The sets  $A = \{(x, y) : x \leq 0\}$  and  $B = \{(x, y) : x > 0 \text{ and } y \geq \frac{1}{x}\}$  only permit a non-strictly separating hyperplane.

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

So suppose there exists  $\mathbf{u} \in B$  such that  $\langle \mathbf{n}, \mathbf{u} \rangle - \mu \leq 0$ . We now look at the line defined by the distance minimizer  $\mathbf{d}$  and the point on the “wrong side”  $\mathbf{u}$ . Define  $\mathbf{b}(\lambda) = \mathbf{d} + \lambda(\mathbf{u} - \mathbf{d})$ , and take the derivative of the distance between  $\mathbf{b}(\lambda)$  and  $\mathbf{c}$ . Evaluated at  $\lambda = 0$  (which is when  $\mathbf{b}(\lambda) = \mathbf{d}$ ), this yields

$$\frac{d}{d\lambda} \|\mathbf{b}(\lambda) - \mathbf{c}\|_2^2 \Big|_{\lambda=0} = 2 \langle \mathbf{d} - \lambda \mathbf{d} + \lambda \mathbf{u} - \mathbf{c}, \mathbf{u} - \mathbf{d} \rangle \Big|_{\lambda=0} = 2 \langle \mathbf{d} - \mathbf{c}, \mathbf{u} - \mathbf{d} \rangle.$$

However, this would imply that the gradient is strictly negative since

$$\begin{aligned} \langle \mathbf{n}, \mathbf{u} \rangle - \mu &= \langle \mathbf{d} - \mathbf{c}, \mathbf{u} \rangle - \langle \mathbf{d} - \mathbf{c}, \mathbf{d} \rangle + \langle \mathbf{d} - \mathbf{c}, \mathbf{d} \rangle - \mu \\ &= \langle \mathbf{d} - \mathbf{c}, \mathbf{u} - \mathbf{d} \rangle + \|\mathbf{d}\|_2^2 - \langle \mathbf{c}, \mathbf{d} \rangle - \frac{1}{2} \|\mathbf{d}\|_2^2 + \frac{1}{2} \|\mathbf{c}\|_2^2 \\ &= \langle \mathbf{d} - \mathbf{c}, \mathbf{u} - \mathbf{d} \rangle + \frac{1}{2} \|\mathbf{d} - \mathbf{c}\|_2^2 \leq 0. \end{aligned}$$

This contradicts the minimality of  $\mathbf{d}$  and thus concludes this proof.  $\square$

A more general separating hyperplane theorem holds even when the sets are not closed and bounded:

**Theorem 14.2.4** (Separating Hyperplane Theorem). *Given two disjoint convex sets  $A, B \in \mathbb{R}^n$  there exists a hyperplane  $H$  separating them.*

## 14.3 Lagrange Multipliers and Duality of Convex Problems

In this Section, we'll learn about *Lagrange Multipliers* and how they lead to convex duality. But first, let's see an example to help illustrate where these ideas come from.

Imagine you were to prove that for all  $\mathbf{x} \in \mathbb{R}^n$  we have  $\|\mathbf{x}\|_p \leq n^{\frac{1}{2} - \frac{1}{p}} \|\mathbf{x}\|_2$  for some  $1 \leq p \leq 2$ . We can look at this as optimizing  $\max_{\mathbf{x}} \|\mathbf{x}\|_p$  subject to  $\|\mathbf{x}\|_2$  being constant, e.g. simply  $\|\mathbf{x}\|_2 = 1$ . Similarly, we could also compute  $\max_{\|\mathbf{x}\|_2=1} \|\mathbf{x}\|_p^p$ . Then the statement above follows from a scaling argument.

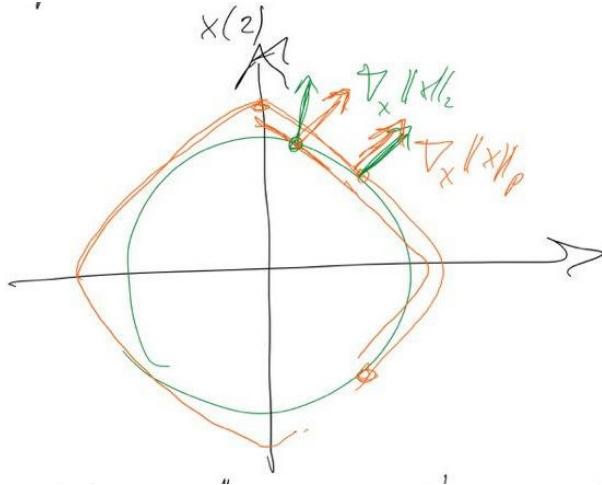


Figure 14.2: Looking at fixed  $\|\mathbf{x}\|_p = \alpha$  and  $\|\mathbf{x}\|_2 = 1$ . (Here,  $p = 1.5$ .)

If we move from  $\mathbf{x}$  to  $\mathbf{x} + \boldsymbol{\delta}$  with  $\boldsymbol{\delta} \perp \nabla_{\mathbf{x}} \|\mathbf{x}\|_2^2$  and  $\boldsymbol{\delta} \not\perp \nabla_{\mathbf{x}} \|\mathbf{x}\|_p^p$  means that for infinitesimally small  $\boldsymbol{\delta}$  the 2-norm stays constant but the  $p$ -norm changes. That means for either  $\mathbf{x} - \boldsymbol{\delta}$  or  $\mathbf{x} + \boldsymbol{\delta}$  the  $p$ -norm increases while the 2-norm stays constant. Hence at the maximum of  $\|\mathbf{x}\|_p$  the gradients of both norms have to be parallel, i.e.

$$\nabla_{\mathbf{x}} \left( \|\mathbf{x}\|_p^p - \lambda \|\mathbf{x}\|_2^2 \right) = 0.$$

This insight is the core idea of Lagrange multipliers (in this case  $\lambda$ ). Note that once we know this, we can conclude that there exists a maximizer where  $\mathbf{x} = \alpha \mathbf{1}_S$ , i.e.  $\mathbf{x}$  is the indicator of some set  $S$ , because otherwise the gradients cannot be parallel. Finally, we can argue that among such vectors, the maximum gap between  $\|\mathbf{x}\|_2$  and  $\|\mathbf{x}\|_p$  arises when  $\mathbf{x} = \mathbf{1}$ .

We caution that the whole argument is rather informal and it takes some work to make it formal.

Note that here the problem is not convex, because  $\{\mathbf{x} : \|\mathbf{x}\|_2^2 = 1\}$  is not convex and because we are asking to *maximize* a norm. In the following we will study Lagrange multipliers for general convex problems.

### 14.3.1 Karush-Kuhn Tucker Optimality Conditions for Convex Problems

A full formal treatment of convex duality would require us to be more careful about using inf and sup in place of min and max, as well as considering problems that have no feasible solutions. Today, we'll ignore these concerns.

Let us consider a general convex optimization problem with convex objective, linear equality constraints and convex inequality constraints

$$\begin{aligned} & \min_{y \in S} \mathcal{E}(y) \\ \text{s.t. } & \mathbf{A}y = \mathbf{b} \\ & \mathbf{c}(y) \leq \mathbf{0}, \end{aligned} \tag{14.1}$$

where  $\mathcal{E}(y) : S \rightarrow \mathbb{R}$  is defined on a convex subset  $S \subseteq \mathbb{R}^n$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{c}(y)$  is a vector of constraints  $\mathbf{c}(y) = (c_i(y))_{i \in [k]}$ . For every  $i \in [k]$  the function  $c_i : S \rightarrow \mathbb{R}$  should be convex, which implies the sublevel set  $\{y : c_i(y) \leq 0\}$  is convex. Given a solution  $y$ , we say an inequality constraint is *tight* at  $y$  if  $c_i(y) = 0$ . In the following we will denote by  $\alpha^* = \mathcal{E}(y^*)$  the optimal value of this program where  $y^*$  is a minimizer.

We will call this *the primal program* – and later we will see that we can associate another related convex program with any such convex program – and we will call this second program *the dual program*.

**Definition 14.3.1** (Primal feasibility). We say that  $y \in S$  is *primal feasible* if all constraints are satisfied, i.e.  $\mathbf{A}y = \mathbf{b}$  and  $\mathbf{c}(y) \leq \mathbf{0}$ .

Now, as we did in our example with the 2- and  $p$ -norms, we will try to understand the relationship between the gradient of the objective function and of the constraint functions at an optimal solution  $y^*$ .

**An (not quite true!) intuition.** Suppose  $y^*$  is an optimal solution to the convex program above. Let us additionally suppose that  $y^*$  is *not* on the boundary of  $S$ . Then, generally speaking, because we are at a constrained minimum of  $\mathcal{E}(y^*)$ , we must have that for any infinitesimal  $\delta$  s.t.  $y^* + \delta$  is also feasible,  $\delta^\top \nabla \mathcal{E}(y^*) \geq 0$ , i.e. the infinitesimal does not decrease the objective. We can also view this as saying that if  $\delta^\top \nabla \mathcal{E}(y^*) < 0$ , the update must be infeasible. But what kind of updates will make  $y^* + \delta$  infeasible? This will be true

if  $\mathbf{a}_j^\top \boldsymbol{\delta} \neq 0$  for some linear constraint  $j$  or, roughly speaking,  $\nabla c_i(\mathbf{y}^*)^\top \boldsymbol{\delta} \neq 0$  for some *tight* inequality constraint. But, if this is true for *all* directions that have a negative inner product with  $\nabla \mathcal{E}(\mathbf{y}^*)$ , then we must have that  $-\nabla \mathcal{E}(\mathbf{y}^*)$  can be written as a linear combination of  $\mathbf{a}_j$ , i.e. gradients of the linear constraint, and of gradients  $\nabla c_i(\mathbf{y}^*)$  of tight inequality constraints, and furthermore the coefficients of the gradients of these tight inequality constraints must be *positive*, so that moving along this direction will increase the function value and violate the constraint.

To recap, given coefficients  $\mathbf{x} \in \mathbb{R}^m$  and  $\mathbf{s} \in \mathbb{R}^k$  with  $\mathbf{s}(i) \geq 0$  if  $c_i(\mathbf{y}^*) = 0$ , and  $\mathbf{s}(i) = 0$  otherwise, we should be able to write

$$-\nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y}^*) = \sum_j \mathbf{x}(j) \mathbf{a}_j + \sum_i \mathbf{s}(i) \nabla c_i(\mathbf{y}^*) \quad (14.2)$$

Note that since  $\mathbf{y}^*$  is feasible, and hence  $\mathbf{c}(\mathbf{y}^*) \leq 0$ , we can write the condition that  $\mathbf{s}(i) \geq 0$  if  $c_i(\mathbf{y}^*) = 0$ , and  $\mathbf{s}(i) = 0$  otherwise, in a very slick way: namely as  $\mathbf{s} \geq \mathbf{0}$  and  $\mathbf{s}^\top \mathbf{c}(\mathbf{y}^*) = 0$ . Traditionally, the variables in  $\mathbf{s}$  are called *slack variables*, because of this, i.e. they are non-zero only if there is no *slack* in the constraint. This condition has a fancy name: when  $\mathbf{s}^\top \mathbf{c}(\mathbf{y}) = 0$  for some feasible  $\mathbf{y}$  and  $\mathbf{s} \geq \mathbf{0}$ , we say that  $\mathbf{y}$  and  $\mathbf{s}$  satisfy *complementary slackness*. We will think of the vectors  $\mathbf{s}$  and  $\mathbf{x}$  as variables that help us prove optimality of a current solution  $\mathbf{y}$ , and we call them *dual variables*.

**Definition 14.3.2** (Dual feasibility). We say  $(\mathbf{x}, \mathbf{s})$  is dual feasible if  $\mathbf{s} \geq \mathbf{0}$ . If additionally  $\mathbf{y}$  is primal feasible, we say  $(\mathbf{y}, \mathbf{x}, \mathbf{s})$  is primal-dual feasible.

Now, we have essentially argued that at any optimal solution  $\mathbf{y}^*$ , we must have that complementary slackness holds, and that Equation (14.2) holds for some  $\mathbf{x}$  and some  $\mathbf{s} \geq \mathbf{0}$ . However, while this intuitive explanation is largely correct, it turns out that it can fail for technical reasons in some weird situations<sup>2</sup>. Nonetheless, under some mild conditions, it is indeed true that the conditions we argued for above must hold at any optimal solution. These conditions have a name: The Karush-Kuhn-Tucker Conditions. For convenience, we will state the conditions using  $\nabla \mathbf{c}(\mathbf{y})$  to denote the matrix whose  $i$ th column is given by  $\nabla c_i(\mathbf{y})$ . This is sometimes called the Jacobian of  $\mathbf{c}$ .

**Definition 14.3.3** (The Karush-Kuhn-Tucker (KKT) Conditions). Given a convex optimization problem of form (14.1) where the domain  $S$  is *open*. Suppose  $\mathbf{y}, \mathbf{x}, \mathbf{s}$  satisfy the following conditions:

---

<sup>2</sup>Consider the following single-variable optimization problem

$$\begin{aligned} & \min_{x \in \mathbb{R}} x \\ & \text{s.t. } x^2 = 0. \end{aligned}$$

This has only a single feasible point  $x = 0$ , which must then be optimal. But at this point, the gradient of the constraint function is zero, while the gradient of the objective is non-zero. Thus our informal reasoning breaks down, because there exists an infeasible direction  $\delta$  we can move along where the constraint function grows, but at a rate of  $O(\delta^2)$ .

- $\mathbf{A}\mathbf{y} = \mathbf{b}$  and  $\mathbf{c}(\mathbf{y}) \leq 0$  (primal feasibility)
- $\mathbf{s} \geq 0$  (dual feasibility)
- $\nabla_{\mathbf{y}}\mathcal{E}(\mathbf{y}) + \mathbf{A}^\top \mathbf{x} + \nabla \mathbf{c}(\mathbf{y})\mathbf{s} = \mathbf{0}$  (KKT gradient condition, i.e. Eq. (14.2) restated)
- $\mathbf{s}(i) \cdot \mathbf{c}_i(\mathbf{y}) = 0$  for all  $i$  (complementary slackness)

Then we say that  $\mathbf{y}, \mathbf{x}, \mathbf{s}$  satisfy the *Karush-Kuhn-Tucker* (KKT) conditions. Note that because of primal-dual feasibility, we can also write complementary slackness as  $\mathbf{s}^T \mathbf{c}(\mathbf{y}) = 0$ .

Note that we tried to informally argued that the KKT conditions must hold at an optimality solution (i.e. KKT is necessary for optimality), but this is not quite true without additional assumptions, as shown by our simple counter-example  $\min_{x \in R, x^2 \leq 0} x$ . On the other hand, it turns out that KKT is always sufficient for optimality, i.e. ‘KKT at  $(\mathbf{y}, \mathbf{x}, \mathbf{s})$ ’  $\Rightarrow$  ‘ $\mathbf{y}$  is optimal’ – as we will prove shortly (Theorem 14.3.12).

### 14.3.2 Slater’s Condition

There exists many different mild technical conditions under which the KKT conditions do indeed hold at any optimal solution  $\mathbf{y}$ . The simplest and most useful is probably *Slater’s condition*.

**Definition 14.3.4** (Slater’s condition with full domain). A (primal) problem as defined in (14.1) with  $S = \mathbb{R}^n$  fulfills Slater’s condition if there exists a *strictly feasible* point, i.e. there exists  $\tilde{\mathbf{y}}$  s.t.  $\mathbf{A}\tilde{\mathbf{y}} = \mathbf{b}$  and  $\mathbf{c}(\tilde{\mathbf{y}}) < \mathbf{0}$ . This means that the strictly feasible point  $\tilde{\mathbf{y}}$  lies strictly inside the set  $\{\mathbf{y} : \mathbf{c}(\mathbf{y}) \leq \mathbf{0}\}$  defined by the inequality constraints.

One way to think about Slater’s condition is that your inequality constraints should not restrict the solution space to be lower-dimensional. This is a degenerate case as the sublevel sets of the inequality constraints are generically full-dimensional and you want to avoid this degenerate case.

We can also extend Slater’s condition to the case when the domain  $S$  is an open set. To extend Slater’s condition to this case, we need the notion of a “relative interior”.

**Definition 14.3.5** (Relative interior). Given a convex set  $S \subset \mathbb{R}^n$ , the *relative interior* of  $S$  is

$$\text{relint}(S) = \{\mathbf{x} \in S : \text{for all } \mathbf{y} \in S \text{ there exists } \epsilon > 0 \text{ such that } \mathbf{x} + \epsilon(\mathbf{x} - \mathbf{y}) \in S\}.$$

In other words,  $\mathbf{x} \in \text{relint}(S)$  if starting at  $\mathbf{x} \in S$  we can move “away” from any  $\mathbf{y} \in S$  by a little and still be in  $S$ . As an example, suppose  $S = \{(s, t) \in \mathbb{R}^2 \text{ such that } s \geq 0 \text{ and } t = 0\}$ . Then  $(0, 0) \in S$  but  $(0, 0) \notin \text{relint}(S)$ , while  $(1, 0) \in \text{relint}(S)$ .

Now, we can state a more general version of Slater’s condition.

**Definition 14.3.6** (Slater's condition). A (primal) problem as defined in (14.1) fulfills Slater's condition if there exists a *strictly feasible* point  $\tilde{\mathbf{y}} \in \text{relint}(S)$ . We require  $\mathbf{A}\tilde{\mathbf{y}} = \mathbf{b}$  and  $\mathbf{c}(\tilde{\mathbf{y}}) < \mathbf{0}$ . This means that the strictly feasible point  $\tilde{\mathbf{y}}$  lies strictly inside the set  $\{\mathbf{y} : \mathbf{c}(\mathbf{y}) \leq \mathbf{0}\}$  defined by the inequality constraints.

Finally, we end with a proposition that tells us that given Slater's condition, the KKT are indeed necessary for optimality of our convex programs.

**Proposition 14.3.7** (Given Slater's condition, KKT is necessary for optimality, i.e. optimal solution  $\Rightarrow$  KKT.). Consider a convex program in the form (14.1) that satisfies Slater's condition and has an open set  $S$  as its domain. Suppose  $\mathbf{y}$  is a primal optimal (feasible) solution, then  $\mathbf{y}, \mathbf{x}, \mathbf{s}$  satisfy the KKT conditions.

We will prove this proposition later, after developing some tools we will use in the proof. In fact, we will also see later that assuming Slater's condition, the KKT conditions are sufficient for optimality.

### 14.3.3 The Lagrangian and The Dual Program

Notice that we can also rewrite the KKT gradient condition as

$$\nabla_{\mathbf{y}} \left[ \underbrace{\mathcal{E}(\mathbf{y}) + \mathbf{x}^\top (\mathbf{b} - \mathbf{A}\mathbf{y}) + \mathbf{s}^\top \mathbf{c}(\mathbf{y})}_{(*)} \right] = \mathbf{0}$$

That is, we can write this condition as the gradient of the quantity  $(*)$  is zero. But what is this quantity  $(*)$ ? We call it *the Lagrangian* of the program.

**Definition 14.3.8.** Given a convex program (14.1), we define the *Lagrangian* of the program as

$$L(\mathbf{y}, \mathbf{x}, \mathbf{s}) = \mathcal{E}(\mathbf{y}) + \mathbf{x}^\top (\mathbf{b} - \mathbf{A}\mathbf{y}) + \mathbf{s}^\top \mathbf{c}(\mathbf{y}).$$

We can think of  $\mathbf{x}$  as assigning a price to violating the linear constraints, and of  $\mathbf{s}$  as assigning a price to violating the inequality constraints. The KKT gradient condition tells us that at the given prices, the there is no benefit gained from locally violating the constraints – i.e. changing the primal solution  $\mathbf{y}$  would not improve the cost.

Notice that if  $\mathbf{y}, \mathbf{x}, \mathbf{s}$  are primal-dual feasible, then

$$\begin{aligned} \mathcal{E}(\mathbf{y}) &= \mathcal{E}(\mathbf{y}) + \mathbf{x}^\top (\mathbf{b} - \mathbf{A}\mathbf{y}) && \text{as } \mathbf{b} - \mathbf{A}\mathbf{y} = \mathbf{0}. \\ &\geq \mathcal{E}(\mathbf{y}) + \mathbf{x}^\top (\mathbf{b} - \mathbf{A}\mathbf{y}) + \mathbf{s}^\top \mathbf{c}(\mathbf{y}) && \text{as } \mathbf{c}(\mathbf{y}) \leq \mathbf{0} \text{ and } \mathbf{s} \geq \mathbf{0}. \\ &= L(\mathbf{y}, \mathbf{x}, \mathbf{s}) && \end{aligned} \tag{14.3}$$

Thus, for primal-dual feasible variables, the Lagrangian is always a lower bound on the objective value.

The primal problem can be written in terms of the Lagrangian.

$$\alpha^* = \min_{\mathbf{y}} \max_{\mathbf{x}; \mathbf{s} \geq 0} L(\mathbf{y}, \mathbf{x}, \mathbf{s}) \quad (14.4)$$

This is because for a minimizing  $\mathbf{y}$  all constraints have to be satisfied and the Lagrangian simplifies to  $L(\mathbf{y}, \mathbf{x}, \mathbf{s}) = \mathcal{E}(\mathbf{y})$ . If  $\mathbf{Ax} - \mathbf{b} = \mathbf{0}$  was violated, making  $\mathbf{x}$  large sends  $L(\mathbf{y}, \mathbf{x}, \mathbf{s}) \rightarrow \infty$ . And if  $\mathbf{c}(\mathbf{y}) \leq \mathbf{0}$  is violated, we can make  $L(\mathbf{y}, \mathbf{x}, \mathbf{s}) \rightarrow \infty$  by choosing large  $\mathbf{s}$ .

Note that we require  $\mathbf{s} \geq \mathbf{0}$ , as we only want to penalize the violation of the inequality constraints in one direction, i.e. when  $\mathbf{c}(\mathbf{y}) > 0$ .

We also define a Lagrangian only in terms of the dual variables by minimizing over  $\mathbf{y}$  as

$$L(\mathbf{x}, \mathbf{s}) = \min_{\mathbf{y}} L(\mathbf{y}, \mathbf{x}, \mathbf{s}).$$

When  $\mathbf{s} \geq 0$ , we have that  $L(\mathbf{y}, \mathbf{x}, \mathbf{s})$  is a convex function of  $\mathbf{y}$ . For each  $\mathbf{y}$ , the Lagrangian  $L(\mathbf{y}, \mathbf{x}, \mathbf{s})$  is linear in  $(\mathbf{x}, \mathbf{s})$  and hence also concave in them. Hence  $L(\mathbf{x}, \mathbf{s})$  is a concave function, because it is the pointwise minimum (over  $\mathbf{y}$ ), of a collection of concave functions in  $(\mathbf{x}, \mathbf{s})$ .

$L(\mathbf{x}, \mathbf{s})$  is defined by minimizing  $L(\mathbf{y}, \mathbf{x}, \mathbf{s})$  over  $\mathbf{y}$ , i.e. what is the worst case value of the lower bound  $L(\mathbf{y}, \mathbf{x}, \mathbf{s})$  across all  $\mathbf{y}$ . We can think of this as computing how good the given ‘‘prices’’ are at approximately enforcing the constraints. This naturally leads to a new optimization problem: How can we choose our prices  $\mathbf{x}, \mathbf{s}$  to get the best (highest) possible lower bound?

**Definition 14.3.9** (Dual problem). We define the *dual problem* as

$$\max_{\substack{\mathbf{x}, \mathbf{s} \\ \mathbf{s} \geq 0}} \min_{\mathbf{y}} L(\mathbf{y}, \mathbf{x}, \mathbf{s}) = \max_{\substack{\mathbf{x}, \mathbf{s} \\ \mathbf{s} \geq 0}} L(\mathbf{x}, \mathbf{s}) \quad (14.5)$$

and denote the optimal dual value by  $\beta^*$ .

The dual problem is really a convex optimization problem in disguise, because we can flip the sign of  $-L(\mathbf{x}, \mathbf{s})$  to get a convex function and minimizing this is equivalent to maximizing  $L(\mathbf{x}, \mathbf{s})$ .

$$\max_{\substack{\mathbf{x}, \mathbf{s} \\ \mathbf{s} \geq 0}} L(\mathbf{x}, \mathbf{s}) = - \min_{\substack{\mathbf{x}, \mathbf{s} \\ \mathbf{s} \geq 0}} -L(\mathbf{x}, \mathbf{s})$$

When  $\mathbf{x}$  and  $\mathbf{s}$  are optimal for the dual program, we say they are dual optimal. And for convenience, when we also have a primal optimal  $\mathbf{y}$ , altogether, we will say that  $(\mathbf{y}, \mathbf{x}, \mathbf{s})$  are primal-dual optimal.

For any primal-dual feasible  $\mathbf{y}, \mathbf{x}, \mathbf{s}$  we have  $L(\mathbf{y}, \mathbf{x}, \mathbf{s}) \leq \mathcal{E}(\mathbf{y})$  (see Equation (14.3)) and hence also  $L(\mathbf{x}, \mathbf{s}) = \min_{\mathbf{y}} L(\mathbf{y}, \mathbf{x}, \mathbf{s}) \leq \mathcal{E}(\mathbf{y})$ .

In other words  $\max_{\mathbf{x}; \mathbf{s} \geq 0} L(\mathbf{x}, \mathbf{s}) = \beta^* \leq \alpha^*$ . This is referred to as *weak duality*.

Using the forms in Equations (14.5) and (14.4), we can also state this as

**Theorem 14.3.10** (The Weak Duality Theorem). *For any convex program (14.5) and its dual, we have*

$$\alpha^* = \min_{\mathbf{y}} \max_{\mathbf{x}; \mathbf{s} \geq \mathbf{0}} L(\mathbf{y}, \mathbf{x}, \mathbf{s}) \geq \max_{\mathbf{x}; \mathbf{s} \geq \mathbf{0}} \min_{\mathbf{y}} L(\mathbf{y}, \mathbf{x}, \mathbf{s}) = \beta^*.$$

### 14.3.4 Strong Duality and KKT

So now that we have proved weak duality  $\beta^* \leq \alpha^*$ , what is strong duality?  $\beta^* = \alpha^*$ ? The answer is yes, but strong duality only holds under some conditions. Again, a simple sufficient condition is Slater's condition (Definition 14.3.6).

**Theorem 14.3.11.** *For a program (14.1) satisfying Slater's condition, strong duality holds, i.e.  $\alpha^* = \beta^*$ . In other words, the optimal value of the primal problem  $\alpha^*$  is equal to the optimal value of the dual.*

Note that at primal optimal  $\mathbf{y}^*$  and dual optimal  $\mathbf{x}^*, \mathbf{s}^*$ , we have

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

Thus, we can conclude that  $L(\mathbf{y}^*, \mathbf{x}^*, \mathbf{s}^*) = \alpha^* = \beta^*$ .

**KKT is sufficient for optimality.** We introduced the KKT conditions by informally thinking about conditions that should be true at a local extremum. In fact, for convex problems (with sufficient differentiability), the KKT conditions are sufficient for optimality, as the next theorem shows. As a bonus, we also get strong duality when they hold.

**Theorem 14.3.12.** *Consider a convex program (14.1) with an open domain set  $S$ . Then if the KKT conditions hold at  $(\tilde{\mathbf{y}}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}})$ , they must be primal-dual optimal, and strong duality must hold.*

*Proof.*  $\tilde{\mathbf{y}}$  is global minimizer of  $\mathbf{y} \mapsto L(\mathbf{y}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}})$ , since this function is convex with vanishing gradient at  $\tilde{\mathbf{y}}$ . Hence,

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

On the other hand, due to primal feasibility and complementary slackness,

$$L(\tilde{\mathbf{y}}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}}) = \mathcal{E}(\tilde{\mathbf{y}}) + \tilde{\mathbf{x}}^\top (\mathbf{b} - \mathbf{A}^\top \tilde{\mathbf{y}}) + \tilde{\mathbf{s}}^\top \mathbf{c}(\tilde{\mathbf{y}}) = \mathcal{E}(\tilde{\mathbf{y}}) \geq \alpha^*.$$

Thus,  $\beta^* \geq \alpha^*$ . But also  $\beta^* \leq \alpha^*$  by weak duality. Therefore,  $\beta^* = \alpha^*$  and  $\tilde{\mathbf{y}}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}}$  are primal/dual optimal.  $\square$

**When strong duality holds, KKT is necessary for optimality.** When strong duality holds, the KKT conditions must necessarily hold at any optimal solution.

Consider a convex program (14.1) with an open set  $S$  as its domain. Suppose  $\mathbf{y}^*$  is an optimizer of the primal problem and  $\mathbf{x}^*, \mathbf{s}^*$  for the dual, and suppose that strong duality holds. We thus have

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

Because  $L(\mathbf{y}, \mathbf{x}^*, \mathbf{s}^*)$  is a convex function in  $\mathbf{y}$ , it also follows that as  $\mathcal{E} : S \rightarrow \mathbb{R}$  and  $\mathbf{c}$  are differentiable then we must have that the gradient w.r.t.  $\mathbf{y}$  is zero, i.e.

$$\nabla_{\mathbf{y}} L(\mathbf{y}, \mathbf{x}^*, \mathbf{s}^*)|_{\mathbf{y}=\mathbf{y}^*} = 0 \quad (14.6)$$

This says exactly that the KKT gradient condition holds at  $(\mathbf{y}^*, \mathbf{x}^*, \mathbf{s}^*)$ .

Next, we want to confirm that complementary slackness holds for our optimal pair  $(\mathbf{y}^*, \mathbf{x}^*, \mathbf{s}^*)$ . We can see that

$$\mathcal{E}(\mathbf{y}^*) = \alpha^* = \mathcal{E}(\mathbf{y}^*) + \mathbf{x}^\top (\mathbf{b} - \mathbf{A}\mathbf{y}^*) + \mathbf{s}^\top \mathbf{c}(\mathbf{y}^*) = \mathcal{E}(\mathbf{y}^*) + \mathbf{s}^\top \mathbf{c}(\mathbf{y}^*)$$

and hence when the  $i$ -th convex constraint is not active, i.e.  $c_i(\mathbf{y}^*) < 0$  the slack must be zero, i.e.  $\mathbf{s}(i) = 0$ . Conversely if the slack is non-zero, that is  $\mathbf{s}(i) \neq 0$  implies that the constraint is active, i.e.  $c_i(\mathbf{y}^*) = 0$ . This says precisely that the complementary slackness condition holds at primal-dual optimal  $(\mathbf{y}^*, \mathbf{x}^*, \mathbf{s}^*)$ . Combined with our previous observation Equation (14.6), we get the following result.

**Theorem 14.3.13.** *Consider a convex program (14.1) with an open domain set  $S$  and whose dual satisfies strong duality. Then KKT conditions necessarily hold at primal-dual optimal  $(\mathbf{y}^*, \mathbf{x}^*, \mathbf{s}^*)$ .*

Theorem 14.3.13 combined with Theorem 14.3.11 immediately imply Proposition 14.3.7.

A good reference for basic convex duality theory is Boyd's free online book "Convex optimization" (linked to on the course website). It provides a number of different interpretations of duality. One particularly interesting one comes from economics: economists see the slack variables  $\mathbf{s}$  as prices for violating the constraints.

### 14.3.5 Proof that Slater's Condition Implies Strong Duality

In this section, we'll prove Theorem 14.3.11. But, before we prove the theorem, let's make a few observations to get us warmed up. If you get bored, skip ahead to the proof.

It is sufficient to prove that  $\alpha^* \leq \beta^*$ , as the statement then follows in conjunction with weak duality. We define the set

$$G = \{(\mathcal{E}(\mathbf{y}), \mathbf{A}\mathbf{y} - \mathbf{b}, \mathbf{c}(\mathbf{y})) : \mathbf{y} \in S\},$$

where  $S \subseteq \mathbb{R}^n$  is the domain of  $\mathcal{E}$ .

Immediately, we observe that we can write the optimal primal value as

$$\alpha^* = \min\{t : (t, \mathbf{v}, \mathbf{u}) \in G, \mathbf{v} = \mathbf{0}, \mathbf{u} \leq \mathbf{0}\}.$$

Similarly, we can write the Lagrangian (after minimizing over  $\mathbf{y}$ )

$$L(\mathbf{x}, \mathbf{s}) = \min_{(t, \mathbf{v}, \mathbf{u}) \in G} (1, \mathbf{x}, \mathbf{s})^\top (t, \mathbf{v}, \mathbf{u}).$$

This is equivalent to the inequality, for  $(t, \mathbf{v}, \mathbf{u}) \in G$ ,

$$(1, \mathbf{x}, \mathbf{s})^\top (t, \mathbf{v}, \mathbf{u}) \geq L(\mathbf{x}, \mathbf{s}).$$

which defines a hyperplane with  $\mathbf{n} = (1, \mathbf{x}, \mathbf{s})$  and  $\mu = L(\mathbf{x}, \mathbf{s})$  such that  $G$  is on one side.

To establish strong duality, we would like to show the existence of a hyperplane such that for  $(t, \mathbf{v}, \mathbf{u}) \in G$

$$\mathbf{n}^\top (t, \mathbf{v}, \mathbf{u}) \geq \alpha^* \text{ and } \mathbf{n} = (1, \widehat{\mathbf{x}}, \widehat{\mathbf{s}}) \text{ with } \widehat{\mathbf{s}} \geq \mathbf{0}.$$

Then we would immediately get

$$\beta^* \geq L(\widehat{\mathbf{x}}, \widehat{\mathbf{s}}) = \min_{(t, \mathbf{v}, \mathbf{u}) \in G} (1, \mathbf{x}, \mathbf{s})^\top (t, \mathbf{v}, \mathbf{u}) \geq \alpha^*.$$

Perhaps not surprisingly, we will use the Separating Hyperplane Theorem. What are the challenges we need to deal with?

- We need to replace  $G$  with a convex set (which we will call  $A$ ) and separate  $A$  from some other convex set (which we will call  $B$ ).
- We need to make sure the hyperplane normal  $\mathbf{n}$  has 1 in the first coordinate and  $\mathbf{s} \geq \mathbf{0}$ , and the hyperplane threshold is  $\alpha^*$ .

*Proof of Theorem 14.3.11.* For simplicity, our proof will assume that  $S = \mathbb{R}^n$ , but only a little extra work is required to handle the general case.

Let's move to on finding two convex disjoists sets  $A, B$  to enable the use of the separating hyperplane Theorem 14.2.4.

First set we define  $A$ , roughly speaking, as a multi-dimensional epigraph of  $G$ . More precisely

$$A = \{(t, \mathbf{v}, \mathbf{u}) : \exists \mathbf{y} \in S, t \geq \mathcal{E}(\mathbf{y}), \mathbf{v} = \mathbf{A}\mathbf{y} - \mathbf{b}, \mathbf{u} \geq \mathbf{c}(\mathbf{y})\}.$$

Note that  $A$  is a convex set. The proof is similar to the proof that the epigraph of a convex function is a convex set. The optimal value of the primal program can be now written as

$$\alpha^* = \min_{(t, \mathbf{0}, \mathbf{0}) \in A} t.$$

And we define another set  $B$  of the same dimensionality as  $A$  by

$$B := \{(r \in \mathbb{R}, \mathbf{0} \in \mathbb{R}^m, \mathbf{0} \in \mathbb{R}^k) : r < \alpha^*\}.$$

This set  $B$  is convex, as it is a ray. An example of two such sets  $A, B$  is illustrated in Figure 14.3.

We show that  $A \cap B = \emptyset$  by contradiction. Suppose  $A, B$  are not disjoint; then there exists  $\mathbf{y}$  such that

$$(\mathcal{E}(\mathbf{y}), \mathbf{A}\mathbf{y} - \mathbf{b}, \mathbf{c}(\mathbf{y})) = (r, \mathbf{0}, \mathbf{u})$$

with  $\mathbf{u} \leq \mathbf{0}$ . But this means that  $\mathbf{y}$  is feasible and  $\mathcal{E}(\mathbf{y}) = r < \alpha^*$ ; contradicting the optimality of  $\alpha^*$ .

To make things simpler, we assume that our linear constraint matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , has full row rank and  $m < n$  (but very little extra work is required to deal with the remaining cases, which we omit).

As we just proved,  $A$  and  $B$  are convex and disjoint sets and hence the separating hyperplane theorem (Theorem 14.2.4) we introduced earlier in this chapter implies the existence a separating hyperplane. This means there exists a normal  $\mathbf{n} = (\tilde{\rho}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}})$  and threshold  $\mu$  and with  $A$  on one side, i.e.

$$(t, \mathbf{v}, \mathbf{u}) \in A \implies (t, \mathbf{v}, \mathbf{u})^\top (\tilde{\rho}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}}) \geq \mu \quad (14.7)$$

and the set  $B$  on the other side:

$$(t, \mathbf{v}, \mathbf{u}) \in B \implies (t, \mathbf{v}, \mathbf{u})^\top (\tilde{\rho}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}}) \leq \mu. \quad (14.8)$$

Now, we claim that  $\tilde{\mathbf{s}} \geq 0$ . Suppose  $\tilde{\mathbf{s}}(i) < 0$ , then for  $\mathbf{u}(i) \rightarrow \infty$  the threshold would grow unbounded, i.e.  $\mu \rightarrow -\infty$  contradicting that the threshold  $\mu$  is finite by the separating hyperplane theorem. Similarly we claim  $\tilde{\rho} \geq 0$ , as if this were not the case, having  $t \rightarrow \infty$  implies that  $\mu \rightarrow -\infty$  again contradicting the finiteness of  $\mu$ .

From Equation (14.8) it follows that  $t\tilde{\rho} \leq \mu$  for all  $t < \alpha^*$  which implies that  $t\tilde{\rho} \leq \mu$  for  $t = \alpha^*$  by taking the limit. Hence we have  $\alpha^*\tilde{\rho} \leq \mu$ . From  $(t, \mathbf{v}, \mathbf{u}) \in A$  we get from Equation (14.7)

$$(\tilde{\rho}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}})^\top (t, \mathbf{v}, \mathbf{u}) \geq \mu \geq \alpha^*\tilde{\rho}$$

and thus

$$(\tilde{\rho}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}})^\top (\mathcal{E}(\mathbf{y}), \mathbf{A}\mathbf{y} - \mathbf{b}, \mathbf{c}(\mathbf{y})) \geq \alpha^*\tilde{\rho}. \quad (14.9)$$

Now we consider two cases; starting with the “good” case where  $\tilde{\rho} > 0$ . Dividing Equation (14.9) by  $\tilde{\rho}$  gives

$$\mathcal{E}(\mathbf{y}) + \frac{\tilde{\mathbf{x}}^\top}{\tilde{\rho}} (\mathbf{A}\mathbf{y} - \mathbf{b}) + \frac{\tilde{\mathbf{s}}^\top}{\tilde{\rho}} \mathbf{c}(\mathbf{y}) \geq \alpha^*.$$

Noting that the left hand side above is  $L(\mathbf{y}, \frac{\tilde{\mathbf{x}}}{\tilde{\rho}}, \frac{\tilde{\mathbf{s}}}{\tilde{\rho}})$  and that the equation holds for arbitrary  $\mathbf{y}$ ; therefore also for the minimum we get

$$\min_{\mathbf{y}} L \left( \mathbf{y}, \frac{\tilde{\mathbf{x}}}{\tilde{\rho}}, \frac{\tilde{\mathbf{s}}}{\tilde{\rho}} \right) \geq \alpha^*$$

and hence via definition of  $\beta^*$  finally

$$\beta^* \geq L\left(\frac{\tilde{\mathbf{x}}}{\tilde{\rho}}, \frac{\tilde{\mathbf{s}}}{\tilde{\rho}}\right) \geq \alpha^*.$$

Next consider the “bad” case  $\tilde{\rho} = 0$ . As  $\alpha^* \tilde{\rho} \leq \mu$ , we have  $0 \leq \mu$ . From Equation (14.7) we get

$$\mathbf{c}(\mathbf{y})^\top \mathbf{s} + \mathbf{x}^\top (\mathbf{b} - \mathbf{A}\mathbf{y}) \geq \mu \geq 0.$$

As Slater’s condition holds, there is an interior point  $\tilde{\mathbf{y}}$ , i.e. it satisfies  $\mathbf{b} - \mathbf{A}\tilde{\mathbf{y}} = \mathbf{0}$  and  $\mathbf{c}(\tilde{\mathbf{y}}) < \mathbf{0}$ . Together with the equation above this yields

$$\mathbf{c}(\tilde{\mathbf{y}})^\top \tilde{\mathbf{s}} + \tilde{\mathbf{x}}^\top \mathbf{0} \geq 0$$

which implies  $\mathbf{c}(\tilde{\mathbf{y}})^\top \tilde{\mathbf{s}} \geq 0$  and as  $\mathbf{c}(\tilde{\mathbf{y}}) < \mathbf{0}$  this means  $\tilde{\mathbf{s}} = \mathbf{0}$ .

As the normal  $(\tilde{\rho}, \tilde{\mathbf{s}}, \tilde{\mathbf{x}})$  of the hyperplane can not be all zeroes, this means the last “component”  $\tilde{\mathbf{x}}$  must contain a non-zero entry, i.e.  $\tilde{\mathbf{x}} \neq \mathbf{0}$ . Furthermore  $\tilde{\mathbf{x}}^\top (\mathbf{b} - \mathbf{A}\tilde{\mathbf{y}}) = \mathbf{0}$ ,  $\mathbf{c}(\tilde{\mathbf{y}}) < \mathbf{0}$  and  $\mathbf{A}$  has full row rank, hence there exists  $\delta$  such that

$$\tilde{\mathbf{x}}^\top (\mathbf{b} - \mathbf{A}(\tilde{\mathbf{y}} + \delta)) < \mathbf{0} \text{ and } \mathbf{c}(\tilde{\mathbf{y}} + \delta) < 0.$$

This, however, means that there is a point in  $A$  on the wrong side of the hyperplane, as

$$(\tilde{\rho}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}})^\top (\mathcal{E}(\tilde{\mathbf{y}} + \delta), \mathbf{b} - \mathbf{A}(\tilde{\mathbf{y}} + \delta), \mathbf{c}(\tilde{\mathbf{y}} + \delta)) < 0$$

but the threshold is  $\mu \geq 0$ . □

**Remark.** Note that our reasoning about why  $\mathbf{s} \geq \mathbf{0}$  in the proof above is very similar to our reasoning for why the primal program can be written as Problem (14.4).

**Example.** As an example of  $A$  and  $B$  as they appear in the above proof, consider

$$\min_{\substack{y \in (0, \infty) \\ 1/y - 1 \leq 0}} y^2$$

This leads to  $\alpha^* = 1$ ,  $y^* = 1$ , and  $A = \{(t, u) : y \in (0, \infty) \text{ and } t > y^2 \text{ and } u \geq 1/y - 1\}$ , and  $B = \{(t, 0) : t < 1\}$  and the separating hyperplane normal is  $\mathbf{n} = (1, 2)$ . These two sets  $A, B$  are illustrated in Figure 14.3.

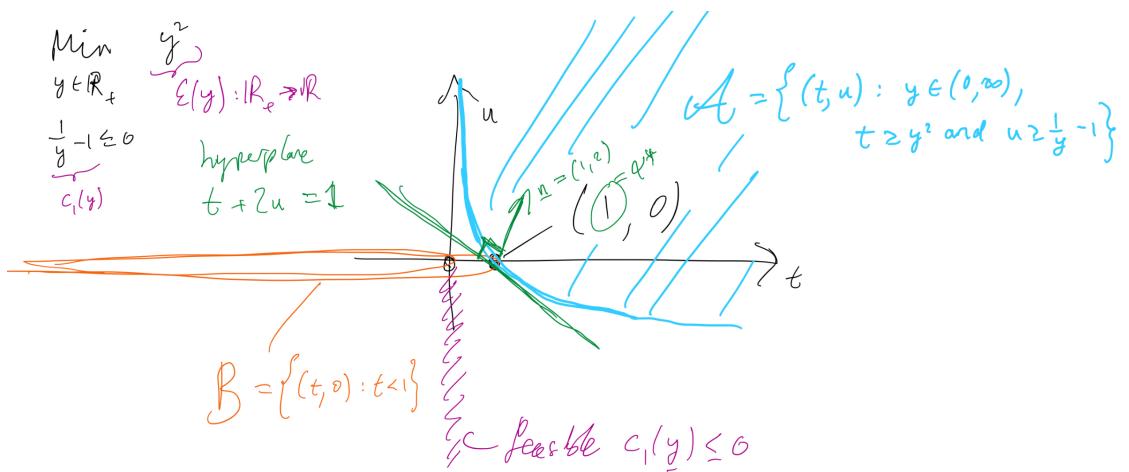


Figure 14.3: Example of the convex sets  $A$  and  $B$  we wish to separate by hyperplane.

# Chapter 15

## Fenchel Conjugates and Newton's Method

### 15.1 Lagrange Multipliers and Convex Duality Recap

Recall the convex optimization program we studied last chapter,

$$\begin{aligned} \min \quad & \mathcal{E}(\mathbf{y}) \\ \text{s.t.} \quad & \mathbf{A}\mathbf{y} = \mathbf{b} \\ & \mathbf{c}(\mathbf{y}) \leq \mathbf{0}, \end{aligned} \tag{15.1}$$

where  $\mathcal{E}(\mathbf{y}) : S \rightarrow \mathbb{R}$  is defined on a subset  $S \subseteq \mathbb{R}^n$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{c}(\mathbf{y})$  is a vector of constraints  $\mathbf{c}(\mathbf{y}) = (c_i(\mathbf{y}))_{i \in [k]}$ . For every  $i \in [k]$  the function  $c_i : S \rightarrow \mathbb{R}$  is convex. We call (15.1) the primal (program) and denote its optimal value by  $\alpha^*$ .

The associated Lagrangian is defined by

$$L(\mathbf{y}, \mathbf{x}, \mathbf{s}) = \mathcal{E}(\mathbf{y}) + \mathbf{x}^T(\mathbf{b} - \mathbf{A}\mathbf{y}) + \mathbf{s}^T\mathbf{c}(\mathbf{y}).$$

where  $\mathbf{x} \in \mathbb{R}^m$ ,  $\mathbf{s} \in \mathbb{R}^k$  are dual variables. The dual (program) is given by

$$\max_{\substack{\mathbf{x}, \mathbf{s} \\ \mathbf{s} \geq \mathbf{0}}} L(\mathbf{x}, \mathbf{s}) \tag{15.2}$$

whose optimal value is denoted by  $\beta^*$ . The dual is always a convex optimization program even though the primal is non-convex. The optimal value of the primal (15.1) can also be written as

$$\alpha^* = \inf_{\mathbf{y}} \sup_{\mathbf{x}; \mathbf{s} \geq \mathbf{0}} L(\mathbf{y}, \mathbf{x}, \mathbf{s}), \tag{15.3}$$

where no constraint is imposed on the primal variable  $\mathbf{y}$ . The optimal value of the dual (15.2) is

$$\beta^* = \sup_{\mathbf{x}; \mathbf{s} \geq \mathbf{0}} \inf_{\mathbf{y}} L(\mathbf{y}, \mathbf{x}, \mathbf{s}). \tag{15.4}$$

Note the only difference between (15.3) and (15.4) is that the positions of “inf” and “sup” are swapped. The weak duality theorem states that the dual optimal value is a lower bound of the primal optimal value, i.e.  $\beta^* \leq \alpha^*$ .

The Slater’s condition for (15.1) open domain  $S$  requires the existence of a *strictly feasible* point, i.e. there exists  $\tilde{\mathbf{y}} \in S$  s.t.  $\mathbf{A}\tilde{\mathbf{y}} = \mathbf{b}$  and  $\mathbf{c}(\tilde{\mathbf{y}}) < \mathbf{0}$ . This means that the strictly feasible point  $\tilde{\mathbf{y}}$  lies inside the interior of the set  $\{\mathbf{y} : \mathbf{c}(\mathbf{y}) \leq \mathbf{0}\}$  defined by the inequality constraints. If the domain  $S$  is not open, Slater’s condition also requires that a strictly feasible point is in the relative interior of  $S$  (NB: when  $S$  is open,  $S$  is equal to its relative interior). The strong duality theorem says that Slater’s condition implies strong duality,  $\beta^* = \alpha^*$ .

We were also introduced to the KKT conditions, and we saw that for our convex programs with continuously differentiable objective and constraints functions, when the domain is open, the conditions are sufficient to imply strong duality and primal-dual optimality of the points that satisfy them,  $\text{KKT} \implies (\mathbf{y}, \mathbf{x}, \mathbf{s})$  primal-dual optimal and we have strong duality’. Finally, we saw that when strong duality holds, KKT must hold at the primal-dual optimal solutions.

**In summary:** *Slater’s condition  $\implies$  strong duality and strong duality  $\iff$  KKT.*

**Example.** In Chapter 12, we gave a combinatorial proof of the min-cut max-flow theorem, and showed that the min-cut program can be expressed as a linear program. Now, we will use the strong duality theorem to give an alternative proof, and directly find the min-cut linear program is the dual program to our maximum flow linear program.

We will assume that Slater’s condition holds for our primal program. Since scaling the flow down enough will always ensure that capacity constraints are strictly satisfied i.e.  $\mathbf{f} < \mathbf{c}$ , the only concern is to make sure that non-negativity constraints are satisfied. This means that there is an  $s$ - $t$  flow that sends a non-zero flow on every edge. In fact, this may not always be possible, but it is easy to detect such edges and remove them without changing the value of the program: an edge  $(u, v)$  should be removed if there is no path  $s$  to  $u$  or no path  $v$  to  $t$ . We can identify all such edges using a BFS from  $s$  along the directed edges and a BFS along reversed directed edges from  $t$ .

$$\begin{aligned} \min_{\substack{F \in \mathbb{R} \\ Bf = Fb_{s,t} \\ 0 \leq f \leq c}} -F &= \min_{F; f \geq \mathbf{0}} \max_{\mathbf{x}; \mathbf{s} \geq \mathbf{0}} -F + \mathbf{x}^\top (F\mathbf{b}_{s,t} - \mathbf{B}\mathbf{f}) + (\mathbf{f} - \mathbf{c})^\top \mathbf{s} \\ &\quad (\text{Slater's condition } \implies \text{strong duality}) \\ - \max_{\substack{F \in \mathbb{R} \\ Bf = Fb_{s,t} \\ 0 \leq f \leq c}} F &= \max_{\mathbf{x}; \mathbf{s} \geq \mathbf{0}} \min_{F; f \geq \mathbf{0}} F(\mathbf{b}_{s,t}^\top \mathbf{x} - 1) + \mathbf{f}^\top (\mathbf{s} - \mathbf{B}^\top \mathbf{x}) - \mathbf{c}^\top \mathbf{s} \\ &= \max_{\substack{\mathbf{x}; \mathbf{s} \geq \mathbf{0} \\ \mathbf{b}_{s,t}^\top \mathbf{x} = 1 \\ \mathbf{s} \geq \mathbf{B}^\top \mathbf{x}}} -\mathbf{c}^\top \mathbf{s} \end{aligned}$$

Thus switching signs gives us

$$\begin{array}{ll} \max_{\substack{F \in \mathbb{R} \\ Bf = Fb_{s,t} \\ 0 \leq f \leq c}} F = & \min_{\substack{\mathbf{x}; \mathbf{s} \geq \mathbf{0} \\ \mathbf{b}_{s,t}^\top \mathbf{x} = 1 \\ \mathbf{s} \geq \mathbf{B}^\top \mathbf{x}}} \mathbf{c}^\top \mathbf{s} \end{array} \quad (15.5)$$

The LHS of (15.5) is exactly the LP formulation of max-flow, while the RHS is exactly the LP formulation of min-cut. Note that we treated the “constraint”  $\mathbf{0} \leq \mathbf{f}$  as a restriction on the domain of  $\mathbf{f}$  rather than a constraint with a dual variable associated with it. We always have this kind of flexibility when deciding how to compute a dual, and some choices may lead to a simpler dual program than others.

## 15.2 Fenchel Conjugates

In this section we will learn about Fenchel conjugates. This is a notion of dual function that is closely related to duality of convex programs, and we will learn more about how dual programs behave by studying these functions.

**Definition 15.2.1** (Fenchel conjugate). Given a (convex) function  $\mathcal{E} : S \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ , its *Fenchel conjugate* is a function  $\mathcal{E}^* : \mathbb{R}^n \rightarrow \mathbb{R}$  defined as

$$\mathcal{E}^*(\mathbf{z}) = \sup_{\mathbf{y} \in S} \langle \mathbf{z}, \mathbf{y} \rangle - \mathcal{E}(\mathbf{y}).$$

**Remark 15.2.2.**  $\mathcal{E}^*$  is a convex function whether  $\mathcal{E}$  is convex or not, since  $\mathcal{E}^*(\mathbf{z})$  is pointwise supremum of a family of convex (here, affine) functions of  $\mathbf{z}$ .

In this course, we have only considered convex functions that are real-valued and continuous and defined on a convex domain. For any such  $\mathcal{E}$ , we have  $\mathcal{E}^{**} = \mathcal{E}$ , i.e. the Fenchel conjugate of the Fenchel conjugate is the original function. This is a consequence of the Fenchel-Moreau theorem, which establishes this under slightly more general conditions. We will not prove this generally, but as part of Theorem 15.2.3 below, we sketch a proof under more restrictive assumptions.

**Example.** Let  $\mathcal{E}(\mathbf{y}) = \frac{1}{p} \|\mathbf{y}\|_p^p$  ( $p > 1$ ). We want to evaluate its Fenchel conjugate  $\mathcal{E}^*$  at any given point  $\mathbf{z} \in \mathbb{R}^n$ . Since  $\mathcal{E}$  is convex and differentiable, the supremum must be achieved at some  $\mathbf{y}$  with vanishing gradient

$$\nabla_{\mathbf{y}} \langle \mathbf{z}, \mathbf{y}^* \rangle - \nabla \mathcal{E}(\mathbf{y}^*) = \mathbf{z} - \nabla \mathcal{E}(\mathbf{y}^*) = \mathbf{0} \iff \mathbf{z} = \nabla \mathcal{E}(\mathbf{y}^*).$$

It's not difficult to see, for all  $i$ ,

$$z(i) = \text{sgn}(\mathbf{y}(i)) |\mathbf{y}(i)|^{p-1}.$$

Then,

$$\begin{aligned}
\mathcal{E}^*(\mathbf{z}) &= \langle \mathbf{z}, \mathbf{y} \rangle - \mathcal{E}(\mathbf{y}) \\
&= \sum_i |\mathbf{z}(i)|^{\frac{1}{p-1}+1} - \frac{1}{p} |\mathbf{z}(i)|^{\frac{p}{p-1}} \\
&\quad (\text{define } q \text{ s.t. } \frac{1}{q} + \frac{1}{p} = 1) \\
&= \frac{1}{q} \|\mathbf{z}\|_q^q.
\end{aligned}$$

More generally, given a convex and differentiable function  $\mathcal{E} : S \rightarrow \mathbb{R}$ , if there exists  $\mathbf{y} \in S$  s.t.  $\mathbf{z} = \nabla \mathcal{E}(\mathbf{y})$ , then  $\mathcal{E}^*(\mathbf{z}) = (\mathbf{y}^*)^\top \nabla \mathcal{E}(\mathbf{y}^*) - \mathcal{E}(\mathbf{y}^*)$ . The Fenchel conjugate and Lagrange duality are closely related, which is demonstrated in the following example.

**Example.** Consider a convex optimization program with only linear constraints,

$$\begin{aligned}
\min_{\mathbf{y} \in \mathbb{R}^n} \quad & \mathcal{E}(\mathbf{y}) \\
\text{s.t.} \quad & \mathbf{A}\mathbf{y} = \mathbf{b}
\end{aligned}$$

where  $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$  is a convex function and  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Then the corresponding dual program is

$$\begin{aligned}
\sup_{\mathbf{x} \in \mathbb{R}^n} \inf_{\mathbf{y} \in \mathbb{R}^n} \mathcal{E}(\mathbf{y}) + \mathbf{x}^\top (\mathbf{b} - \mathbf{A}\mathbf{y}) &= \sup_{\mathbf{x} \in \mathbb{R}^n} \mathbf{b}^\top \mathbf{x} - \sup_{\mathbf{y} \in \mathbb{R}^n} (\mathbf{x}^\top \mathbf{A}\mathbf{y} - \mathcal{E}(\mathbf{y})) \\
&= \sup_{\mathbf{x} \in \mathbb{R}^n} \mathbf{b}^\top \mathbf{x} - \mathcal{E}^*(\mathbf{A}^\top \mathbf{x})
\end{aligned}$$

**Theorem 15.2.3** (Properties of the Fenchel conjugate). *Consider a strictly convex function  $\mathcal{E} : S \rightarrow \mathbb{R}$  where  $S \subseteq \mathbb{R}^n$  is an open convex set. When  $\mathcal{E}$  is differentiable with a Hessian that is positive definite everywhere and its gradient  $\nabla \mathcal{E}$  is surjective onto  $\mathbb{R}^n$ , we have the following three properties:*

1.  $\nabla \mathcal{E}(\nabla \mathcal{E}^*(\mathbf{z})) = \mathbf{z}$  and  $\nabla \mathcal{E}^*(\nabla \mathcal{E}(\mathbf{y})) = \mathbf{y}$
2.  $(\mathcal{E}^*)^* = \mathcal{E}$ , i.e. the Fenchel conjugate of the Fenchel conjugate is the original function.
3.  $\mathbf{H}_{\mathcal{E}^*}(\nabla \mathcal{E}(\mathbf{y})) = \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y})$

primal point $\mathbf{y}$	$\xrightarrow{\nabla_y \mathcal{E}}$ $\xleftarrow{\nabla_z \mathcal{E}^*}$	dual point $\mathbf{z}$
gradient $\nabla_y \mathcal{E}(\mathbf{y}) = \mathbf{z}$		gradient $\nabla_z \mathcal{E}^*(\mathbf{z}) = \mathbf{y}$
Hessian $\mathbf{H}_{\mathcal{E}}(\mathbf{y}) = \mathbf{H}_{\mathcal{E}^*}^{-1}(\mathbf{z})$		Hessian $\mathbf{H}_{\mathcal{E}^*}(\mathbf{y}) = \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y})$

Figure 15.1: Properties of Fenchel conjugate

*Proof sketch. Part 1.* Because the gradient  $\nabla_y \mathcal{E}$  is surjective onto  $\mathbb{R}^n$ , given any  $\mathbf{z} \in \mathbb{R}^n$ , there exists a  $\mathbf{y}$  such that  $\nabla \mathcal{E}(\mathbf{y}) = \mathbf{z}$ . Let  $\mathbf{y}(\mathbf{z})$  be a  $\mathbf{y}$  s.t.  $\nabla \mathcal{E}(\mathbf{y}) = \mathbf{z}$ . It can be shown that because  $\mathcal{E}$  is strictly convex,  $\mathbf{y}(\mathbf{z})$  is unique.

The function  $\mathbf{y} \mapsto \langle \mathbf{z}, \mathbf{y} \rangle - \mathcal{E}(\mathbf{y})$  is concave in  $\mathbf{y}$  and has gradient  $\mathbf{z} - \nabla \mathcal{E}(\mathbf{y})$  and is hence maximized at  $\mathbf{y} = \mathbf{y}(\mathbf{z})$ . This follows because we know a differentiable convex function is minimized when its gradient is zero and so a differentiable concave function is maximized when its gradient is zero.

Then, using the product rule and composition rule of derivatives,

$$\begin{aligned}\nabla \mathcal{E}^*(\mathbf{z}) &= \nabla_z (\langle \mathbf{z}, \mathbf{y}(\mathbf{z}) \rangle - \mathcal{E}(\mathbf{y}(\mathbf{z}))) \\ &= \mathbf{y}(\mathbf{z}) + (\nabla_z (\mathbf{y}(\mathbf{z})^\top)) \mathbf{z} - (\nabla_z (\mathbf{y}(\mathbf{z})^\top)) \underbrace{\nabla_y \mathcal{E}(\mathbf{y}(\mathbf{z}))}_{=-\mathbf{z}} \\ &= \mathbf{y}(\mathbf{z})\end{aligned}$$

Thus we have  $\nabla_y \mathcal{E}(\mathbf{y}(\mathbf{z})) = \mathbf{z}$  and  $\nabla_z \mathcal{E}^*(\mathbf{z}) = \mathbf{y}(\mathbf{z})$ . Combining the two, we have  $\nabla \mathcal{E}(\nabla \mathcal{E}^*(\mathbf{z})) = \mathbf{z}$ .

We can also see that for any  $\mathbf{y}$ , there exists a  $\mathbf{z}$  such that  $\nabla_z \mathcal{E}^*(\mathbf{z}) = \mathbf{y}$ , namely, this is attained by  $\mathbf{z} = \nabla_y \mathcal{E}(\mathbf{y})$ . Thus,  $\nabla \mathcal{E}^*(\nabla \mathcal{E}(\mathbf{y})) = \mathbf{y}$ .

*Part 2.* Observe that

$$\mathcal{E}^{**}(\mathbf{u}) = \sup_{\mathbf{z} \in \mathbb{R}^n} \langle \mathbf{u}, \mathbf{z} \rangle - \mathcal{E}^*(\mathbf{z})$$

and let  $\mathbf{z}(\mathbf{u})$  denote the  $\mathbf{z}$  obtaining the supremum, in the above program. We then have  $\mathbf{u} = \nabla \mathcal{E}^*(\mathbf{z}(\mathbf{u}))$ . Letting  $\mathbf{y}(\mathbf{z})$  be defined as in Part 1, we get  $\mathbf{y}(\mathbf{z}(\mathbf{u})) = \nabla_z \mathcal{E}^*(\mathbf{z}(\mathbf{u})) = \mathbf{u}$

$$\mathcal{E}^{**}(\mathbf{u}) = \langle \mathbf{u}, \mathbf{z}(\mathbf{u}) \rangle - (\langle \mathbf{z}(\mathbf{u}), \mathbf{y}(\mathbf{z}(\mathbf{u})) \rangle - \mathcal{E}(\mathbf{y}(\mathbf{z}(\mathbf{u})))) = \mathcal{E}(\mathbf{u}).$$

*Part 3.* Now we add two infinitesimals  $\boldsymbol{\tau}$  and  $\boldsymbol{\delta}$  to  $\mathbf{z}$  and  $\mathbf{y}$  respectively s.t.

$$\nabla_z \mathcal{E}^*(\mathbf{z} + \boldsymbol{\tau}) = \mathbf{y} + \boldsymbol{\delta}, \quad \nabla_y \mathcal{E}(\mathbf{y} + \boldsymbol{\delta}) = \mathbf{z} + \boldsymbol{\tau}.$$

Then,

$$\nabla_y \mathcal{E}(\mathbf{y} + \boldsymbol{\delta}) - \nabla_y \mathcal{E}(\mathbf{y}) = \boldsymbol{\tau}, \quad \nabla_z \mathcal{E}^*(\mathbf{z} + \boldsymbol{\tau}) - \nabla_z \mathcal{E}^*(\mathbf{z}) = \boldsymbol{\delta}.$$

Since  $\mathbf{H}_{\mathcal{E}}(\mathbf{y})$  measures the change of  $\nabla_{\mathbf{y}}\mathcal{E}(\mathbf{y})$  when  $\mathbf{y}$  changes by an infinitesimal  $\boldsymbol{\delta}$ , then

$$\begin{aligned} \nabla_{\mathbf{y}}\mathcal{E}(\mathbf{y} + \boldsymbol{\delta}) - \nabla_{\mathbf{y}}\mathcal{E}(\mathbf{y}) &\approx \mathbf{H}_{\mathcal{E}}(\mathbf{y})\boldsymbol{\delta} \\ \iff \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y})(\nabla_{\mathbf{y}}\mathcal{E}(\mathbf{y} + \boldsymbol{\delta}) - \nabla_{\mathbf{y}}\mathcal{E}(\mathbf{y})) &\approx \boldsymbol{\delta} \\ \iff \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y})\boldsymbol{\tau} &\approx \boldsymbol{\delta} = \nabla\mathcal{E}^*(\mathbf{z} + \boldsymbol{\tau}) - \nabla\mathcal{E}^*(\mathbf{z}) \\ \iff \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y})\boldsymbol{\tau} &\approx \nabla\mathcal{E}^*(\mathbf{z} + \boldsymbol{\tau}) - \nabla\mathcal{E}^*(\mathbf{z}) \end{aligned} \quad (15.6)$$

Similarly,

$$\mathbf{H}_{\mathcal{E}^*}(\mathbf{z})\boldsymbol{\tau} \approx \nabla_{\mathbf{z}}\mathcal{E}^*(\mathbf{z} + \boldsymbol{\tau}) - \nabla_{\mathbf{z}}\mathcal{E}^*(\mathbf{z}) \quad (15.7)$$

Comparing (15.6) and (15.7), it is easy to see

$$\mathbf{H}_{\mathcal{E}^*}(\mathbf{z}) = \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}) \iff \mathbf{H}_{\mathcal{E}^*}(\nabla\mathcal{E}(\mathbf{y})) = \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}).$$

□

**Remark 15.2.4.** Theorem 15.2.3 can be generalized to show that the Fenchel conjugate has similar nice properties under much more general conditions, e.g. see [BV04].

## 15.3 Newton's Method

### 15.3.1 Warm-up: Quadratic Optimization

First, let us play with a toy example, minimizing a quadratic function

$$\mathcal{E}(\mathbf{y}) = \frac{1}{2}\mathbf{y}^\top \mathbf{A}\mathbf{y} + \mathbf{b}^\top \mathbf{y} + c$$

where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is positive definite. By setting the gradient w.r.t.  $\mathbf{y}$  to zero,

$$\nabla\mathcal{E}(\mathbf{y}) = \mathbf{A}\mathbf{y} + \mathbf{b} = \mathbf{0},$$

we obtain the global minimizer

$$\mathbf{y}^* = -\mathbf{A}^{-1}\mathbf{b}.$$

To make it more like gradient descent, let us start at some “guess” point  $\mathbf{y}$  and take a step  $\boldsymbol{\delta}$  to move to the new point  $\mathbf{y} + \boldsymbol{\delta}$ . Then we try to minimize  $\mathcal{E}(\mathbf{y} + \boldsymbol{\delta})$  by setting the gradient w.r.t.  $\boldsymbol{\delta}$  to zero,

$$\begin{aligned} \nabla_{\boldsymbol{\delta}}\mathcal{E}(\mathbf{y} + \boldsymbol{\delta}) &= \mathbf{A}(\mathbf{y} + \boldsymbol{\delta}) + \mathbf{b} = \mathbf{0} \\ \boldsymbol{\delta} &= -\mathbf{y} - \mathbf{A}^{-1}\mathbf{b} \\ \mathbf{y} + \boldsymbol{\delta} &= -\mathbf{A}^{-1}\mathbf{b} \end{aligned}$$

This gives us exactly global minimizer in just one step. However, the situation changes when the function is not quadratic anymore and thus we do not have a constant Hessian. But taking a step which tries to set the gradient to zero might still be a good idea.

### 15.3.2 $K$ -stable Hessian

Next, consider a convex function  $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$  whose Hessian is “nearly constant”. Recall the Hessian  $\mathbf{H}_{\mathcal{E}}(\mathbf{y})$  aka  $\nabla^2 \mathcal{E}(\mathbf{y})$  at a point  $\mathbf{y}$  is just a matrix of pairwise 2nd order partial derivatives  $\frac{\partial^2 \mathcal{E}(\mathbf{y})}{\partial \mathbf{y}_i \partial \mathbf{y}_j}$ . We say  $\mathcal{E}$  has a  $k$ -stable Hessian if there exists a constant matrix  $\mathbf{A}$  s.t. for all  $\mathbf{y}$

$$\mathbf{H}_{\mathcal{E}}(\mathbf{y}) \approx_K \mathbf{A} \iff \frac{1}{1+K} \mathbf{A} \preceq \mathbf{H}_{\mathcal{E}}(\mathbf{y}) \preceq (1+K) \mathbf{A}.$$

Note that we just require the existence of  $\mathbf{A}$  and do not assume we know  $\mathbf{A}$ . Then a natural question is to ask what convergence rate can be achieved if we take a gradient step “guided” by the Hessian, which is called a “Newton step”. Such method is also known as the 2nd order method. Note that this is very similar to preconditioning.

Now, let us make our setting precise. We want to minimize a convex function  $\mathcal{E}$  with  $k$ -stable Hessian  $\mathbf{A} \succ \mathbf{0}$ . And  $\mathbf{y}^*$  is a global minimizer of  $\mathcal{E}$ . Start from some initial point  $\mathbf{y}_0$ . The update rule is

$$\mathbf{y}_{i+1} = \mathbf{y}_i - \alpha \cdot \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}_i) \nabla \mathcal{E}(\mathbf{y}_i),$$

where  $\alpha$  is the step size and it will be decided later.

**Theorem 15.3.1.**  $\mathcal{E}(\mathbf{y}_k) - \mathcal{E}(\mathbf{y}^*) \leq \epsilon (\mathcal{E}(\mathbf{y}_0) - \mathcal{E}(\mathbf{y}^*))$  when  $k > (K+1)^6 \log(1/\epsilon)$ .

*Proof.* By Taylor’s theorem, there exists  $\tilde{\mathbf{y}} \in [\mathbf{y}, \mathbf{y} + \boldsymbol{\delta}]$  s.t.

$$\begin{aligned} \mathcal{E}(\mathbf{y} + \boldsymbol{\delta}) &= \mathcal{E}(\mathbf{y}) + \nabla \mathcal{E}(\mathbf{y})^\top \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^\top \mathbf{H}_{\mathcal{E}}(\tilde{\mathbf{y}}) \boldsymbol{\delta} \\ &\leq \underbrace{\mathcal{E}(\mathbf{y}) + \nabla \mathcal{E}(\mathbf{y})^\top \boldsymbol{\delta} + \frac{(K+1)^2}{2} \boldsymbol{\delta}^\top \mathbf{H}_{\mathcal{E}}(\mathbf{y}) \boldsymbol{\delta}}_{=: f(\boldsymbol{\delta})} \end{aligned} \tag{15.8}$$

where the inequality comes from the  $K$ -stability of the Hessian,

$$\mathbf{H}_{\mathcal{E}}(\tilde{\mathbf{y}}) \preceq (1+K) \mathbf{A} \preceq (1+K)^2 \mathbf{H}_{\mathcal{E}}(\mathbf{y}).$$

Observe that  $f(\boldsymbol{\delta})$  is a convex quadratic function in  $\boldsymbol{\delta}$ . By minimizing it, or equivalently setting  $\nabla_{\boldsymbol{\delta}} f(\boldsymbol{\delta}^*) = \mathbf{0}$ , we get

$$\boldsymbol{\delta}^* = -\frac{1}{(K+1)^2} \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}) \nabla_y \mathcal{E}(\mathbf{y}) \tag{15.9}$$

Here, the step size  $\alpha$  is equal to  $(K+1)^{-2}$ . Then, plugging (15.9) into (15.8),

$$\begin{aligned} \mathcal{E}(\mathbf{y} + \boldsymbol{\delta}^*) &\leq \mathcal{E}(\mathbf{y}) - \frac{1}{2(K+1)^2} \nabla_y \mathcal{E}(\mathbf{y})^\top \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}) \nabla_y \mathcal{E}(\mathbf{y}) \\ &\leq \mathcal{E}(\mathbf{y}) - \frac{1}{2(K+1)^3} \nabla_y \mathcal{E}(\mathbf{y})^\top \mathbf{A}^{-1} \nabla_y \mathcal{E}(\mathbf{y}) \\ &\quad (\text{subtract } \mathcal{E}(\mathbf{y}^*) \text{ on both sides}) \end{aligned}$$

$$\mathcal{E}(\mathbf{y} + \boldsymbol{\delta}^*) - \mathcal{E}(\mathbf{y}^*) \leq \mathcal{E}(\mathbf{y}) - \mathcal{E}(\mathbf{y}^*) - \frac{1}{2(K+1)^3} \underbrace{\nabla_y \mathcal{E}(\mathbf{y})^\top \mathbf{A}^{-1} \nabla_y \mathcal{E}(\mathbf{y})}_{=: \sigma} \tag{15.10}$$

where the second inequality is due to  $K$ -stability of the inverse Hessian,

$$\frac{1}{1+K} \mathbf{A}^{-1} \preceq \mathbf{H}_\varepsilon(\mathbf{y})^{-1} \preceq (1+K) \mathbf{A}^{-1}.$$

Meanwhile, by convexity, we also have

$$\mathcal{E}(\mathbf{y}) - \mathcal{E}(\mathbf{y}^*) \leq \underbrace{(\mathbf{y} - \mathbf{y}^*)^\top \nabla \mathcal{E}(\mathbf{y})}_{=:\gamma}$$

Next, our task is reduced to comparing  $\sigma$  and  $\gamma$ .

We define  $\mathbf{z}_s = \nabla \mathcal{E}(\mathbf{y}^*) + s(\nabla \mathcal{E}(\mathbf{y}) - \nabla \mathcal{E}(\mathbf{y}^*))$  and then  $d\mathbf{z}_s = \nabla \mathcal{E}(\mathbf{y}) ds$ . Using *Theorem 15.2.3*, we have

$$\mathbf{y} - \mathbf{y}^* = \int_0^1 \mathbf{H}_{\varepsilon^*}(\mathbf{z}_s) \nabla \mathcal{E}(\mathbf{y}) ds.$$

Then,

$$\begin{aligned} \nabla \mathcal{E}(\mathbf{y})^\top (\mathbf{y} - \mathbf{y}^*) &= \int_0^1 \nabla \mathcal{E}(\mathbf{y})^\top \mathbf{H}_{\varepsilon^*}(\mathbf{z}_s) \nabla \mathcal{E}(\mathbf{y}) ds \\ &\leq (K+1) \int_0^1 \nabla \mathcal{E}(\mathbf{y})^\top \mathbf{A}^{-1} \nabla \mathcal{E}(\mathbf{y}) ds \\ &\leq (K+1)\sigma \end{aligned} \tag{15.11}$$

Combining (15.10) and (15.11) yields

$$\mathcal{E}(\mathbf{y} + \boldsymbol{\delta}^*) - \mathcal{E}(\mathbf{y}^*) \leq (\mathcal{E}(\mathbf{y}) - \mathcal{E}(\mathbf{y}^*)) \left(1 - \frac{1}{2(K+1)^4}\right).$$

□

**Remark 15.3.2.** The basic idea of relating  $\sigma$  and  $\gamma$  in the above proof is writing the same quantity,  $\nabla \mathcal{E}(\mathbf{y})^\top (\mathbf{y} - \mathbf{y}^*)$ , as two integrations along different lines.  $(K+1)^4$  can be reduced to  $(K+1)^2$  and even to  $(K+1)$  with more care. In some settings, Newton's method converges in  $\log \log(1/\epsilon)$  steps.

### 15.3.3 Linearly Constrained Newton's Method

Let us apply Newton's method to convex optimization programs with only linear constraints,

$$\begin{aligned} \min_{\mathbf{f} \in \mathbb{R}^m} \quad & \mathcal{E}(\mathbf{f}) \\ \text{s.t.} \quad & \mathbf{B}\mathbf{f} = \mathbf{d} \end{aligned}$$

where  $\mathcal{E} : \mathbb{R}^m \rightarrow \mathbb{R}$  is a convex function and  $\mathbf{B} \in \mathbb{R}^{n \times m}$ . Wlog, let  $\mathbf{d} = \mathbf{0}$ , since otherwise we can equivalently deal the following program with  $\mathbf{B}\mathbf{f}_0 = \mathbf{d}$ ,

$$\begin{aligned} \min_{\boldsymbol{\rho} \in \mathbb{R}^m} \quad & \mathcal{E}(\mathbf{f}_0 + \boldsymbol{\rho}) \\ \text{s.t.} \quad & \mathbf{B}\boldsymbol{\rho} = \mathbf{0} \end{aligned}$$

It is useful to think of the variable  $\mathbf{f} \in \mathbb{R}^m$  as a flow in a graph. Define  $C := \{\mathbf{f} : \mathbf{B}\mathbf{f} = \mathbf{0}\}$  which is the kernel space of  $\mathbf{B}$ .  $C$  is also called the “cycle space” as it is the set of cycle flows when treating  $\mathbf{f}$  as flows.

**Analyzing the convergence of Newton’s method with linear constraints.** It is not immediately obvious, but our analysis of Newton step’s and their convergence on objectives with a  $K$ -stable Hessian carries over directly to linearly constrained convex optimization problems. We will only sketch a proof of this. Firstly, we should notice that instead of thinking of our objective function  $\mathcal{E}$  as defined on  $\mathbb{R}^m$  and then constrained to inputs  $\mathbf{f} \in C$ , we can think of a new function  $\hat{\mathcal{E}} : C \rightarrow \mathbb{R}$  defined such that for  $\mathbf{f} \in C$  we have  $\hat{\mathcal{E}}(\mathbf{f}) = \mathcal{E}(\mathbf{f})$ . But,  $C$  is a linear subspace and is isomorphic<sup>1</sup> to  $\mathbb{R}^{\dim(C)}$ . This means that our previous analysis can be directly applied, if we can compute the gradient and Hessian of the function viewed as an *unconstrained* function on  $C$  (or equivalently  $\mathbb{R}^{\dim(C)}$ ). We now have two important questions to answer:

1. What does the gradient and Hessian, and hence Newton steps, of  $\hat{\mathcal{E}}(\mathbf{f})$  look like?
2. Does the  $K$ -stability of the Hessian of  $\mathcal{E}$  carry over to the function  $\hat{\mathcal{E}}$ ?

The gradient and Hessian of  $\hat{\mathcal{E}}$  should live in  $\mathbb{R}^{\dim(C)}$  and  $\mathbb{R}^{\dim(C) \times \dim(C)}$  respectively. Let  $\Pi_C$  be the orthogonal projection matrix onto  $C$ , meaning  $\Pi_C$  is symmetric and  $\Pi_C\boldsymbol{\delta} = \boldsymbol{\delta}$  for any  $\boldsymbol{\delta} \in C$ . Given any  $\mathbf{f} \in C$ , add to it an infinitesimal  $\boldsymbol{\delta} \in C$ , then

$$\begin{aligned}\hat{\mathcal{E}}(\mathbf{f} + \boldsymbol{\delta}) &= \mathcal{E}(\mathbf{f} + \boldsymbol{\delta}) \approx \mathcal{E}(\mathbf{f}) + \langle \nabla \mathcal{E}(\mathbf{f}), \boldsymbol{\delta} \rangle \\ &= \mathcal{E}(\mathbf{f}) + \langle \nabla \mathcal{E}(\mathbf{f}), \Pi_C \boldsymbol{\delta} \rangle \\ &= \mathcal{E}(\mathbf{f}) + \langle \Pi_C \nabla \mathcal{E}(\mathbf{f}), \boldsymbol{\delta} \rangle\end{aligned}$$

From this, we can deduce that the gradient of  $\hat{\mathcal{E}}$  at a point  $\mathbf{f} \in C$  is essentially equal (up to a fixed linear transformation independent of  $\mathbf{f}$ ) to the projection of gradient of  $\nabla \mathcal{E}$  at  $\mathbf{f}$  onto the subspace  $C$ . Similarly,

$$\hat{\mathcal{E}}(\mathbf{f} + \boldsymbol{\delta}) = \mathcal{E}(\mathbf{f} + \boldsymbol{\delta}) \approx \mathcal{E}(\mathbf{f}) + \langle \Pi_C \nabla \mathcal{E}(\mathbf{f}), \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \Pi_C \mathbf{H}_{\mathcal{E}}(\mathbf{f}) \Pi_C \boldsymbol{\delta} \rangle$$

Again from this, we can deduce that the Hessian of  $\hat{\mathcal{E}}$  at a point  $\mathbf{f} \in C$  is essentially equal (again up to a fixed linear transformation independent of  $\mathbf{f}$ ) to the matrix  $\Pi_C \mathbf{H}_{\mathcal{E}}(\mathbf{f}) \Pi_C$ . Note that  $\mathbf{X} \preceq \mathbf{Y}$  implies  $\Pi_C \mathbf{X} \Pi_C \preceq \Pi_C \mathbf{Y} \Pi_C$ , and from this we can see that the Hessian of  $\hat{\mathcal{E}}$  is  $K$ -stable if the Hessian of  $\mathcal{E}$  is. Also note that we were not terribly formal in the discussion above. We can be more precise by replacing  $\Pi_C$  with a linear map from  $\mathbb{R}^{\dim(C)}$  to  $\mathbb{R}^m$  which maps any vector in  $\mathbb{R}^{\dim(C)}$  to a vector in  $C$  and then going through a similar chain of reasoning.

---

<sup>1</sup>You don’t need to know the formal definition of isomorphism on vector spaces. In this context, it means equivalent up to a transformation by an invertible matrix. In fact in our case, the isomorphism is given by an orthonormal matrix.

What is a Newton step  $\boldsymbol{\delta}^*$  w.r.t.  $\hat{\mathcal{E}}$ ? It turns out that for actually computing the Newton step, it is easier to think again of  $\mathcal{E}$  with a constraint that the Newton step must lie in the subspace  $C$ . One can show that this is equivalent to the Newton step of  $\hat{\mathcal{E}}$ , but we omit this.

In the constrained view,  $\boldsymbol{\delta}^*$  should be a minimizer of

$$\begin{aligned} & \min_{\substack{\boldsymbol{\delta} \in \mathbb{R}^m \\ \mathbf{B}\boldsymbol{\delta} = \mathbf{0}}} \underbrace{\langle \nabla \mathcal{E}(\mathbf{f}), \boldsymbol{\delta} \rangle}_{=:g} + \frac{1}{2} \langle \boldsymbol{\delta}, \underbrace{\mathbf{H}_{\mathcal{E}}(\mathbf{f}) \boldsymbol{\delta}}_{=:H} \rangle \\ & \quad (\text{Lagrange duality}) \\ \iff & \max_{\mathbf{x} \in \mathbb{R}^n} \min_{\boldsymbol{\delta} \in \mathbb{R}^m} \underbrace{\langle \mathbf{g}, \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \mathbf{H} \boldsymbol{\delta} \rangle - \mathbf{x}^\top \mathbf{B} \boldsymbol{\delta}}_{\text{Lagrangian } L(\boldsymbol{\delta}, \mathbf{x})} \end{aligned} \quad (15.12)$$

Applying the KKT optimality conditions, one has

$$\begin{aligned} \mathbf{B} \boldsymbol{\delta} &= \mathbf{0}, \\ \nabla_{\boldsymbol{\delta}} L(\boldsymbol{\delta}, \mathbf{x}) &= \mathbf{g} + \mathbf{H} \boldsymbol{\delta} - \mathbf{B}^\top \mathbf{x} = \mathbf{0}, \end{aligned}$$

from which we get

$$\begin{aligned} \boldsymbol{\delta} + \mathbf{H}^{-1} \mathbf{g} &= \mathbf{H}^{-1} \mathbf{B}^\top \mathbf{x} \\ \underbrace{\mathbf{B} \boldsymbol{\delta}}_{=0} + \mathbf{B} \mathbf{H}^{-1} \mathbf{g} &= \mathbf{B} \mathbf{H}^{-1} \mathbf{B}^\top \mathbf{x} \\ \mathbf{B} \mathbf{H}^{-1} \mathbf{g} &= \underbrace{\mathbf{B} \mathbf{H}^{-1} \mathbf{B}^\top}_{=:L} \mathbf{x} \end{aligned}$$

Finally, the solutions to (15.12) are

$$\begin{cases} \mathbf{x}^* &= \mathbf{L}^{-1} \mathbf{B} \mathbf{H}^{-1} \mathbf{g} \\ \boldsymbol{\delta}^* &= -\mathbf{H}^{-1} \mathbf{g} + \mathbf{H}^{-1} \mathbf{B}^\top \mathbf{x}^* \end{cases}$$

It is easy to verify that  $\mathbf{B} \boldsymbol{\delta}^* = \mathbf{0}$ . Thus, our update rule is  $\mathbf{f}_{i+1} = \mathbf{f}_i + \boldsymbol{\delta}^*$ . And we have the following convergence result.

**Theorem 15.3.3.**  $\hat{\mathcal{E}}(\mathbf{f}_k) - \hat{\mathcal{E}}(\mathbf{f}^*) \leq \epsilon \cdot (\hat{\mathcal{E}}(\mathbf{f}_0) - \hat{\mathcal{E}}(\mathbf{f}^*))$  when  $k > 2(K+1)^4 \log(1/\epsilon)$ .

**Remark 15.3.4.** Note if  $\mathcal{E}(\mathbf{f}) = \sum_{i=1}^m \mathcal{E}_i(\mathbf{f}(i))$ , then  $\mathbf{H}_{\mathcal{E}}(\mathbf{f})$  is diagonal. Thus,  $\mathbf{L} = \mathbf{B} \mathbf{H}^{-1} \mathbf{B}^\top$  is indeed a Laplacian provided that  $\mathbf{B}$  is an incidence matrix. Therefore, the linear equations we need to solve to apply Newton's method in a network flow setting are Laplacians, which means we can solve them very quickly.

# Chapter 16

## Interior Point Methods for Maximum Flow

### Background and Notation

In this chapter, we'll learn about interior point methods for solving maximum flow, which is a rich and active area of research [DS08, Mad13, LS20b, LS20a].

We're going to frequently need to refer to vectors arising from elementwise operations combining other vectors.

To that end, given two vector  $\mathbf{a} \in \mathbb{R}^m$ , and  $\mathbf{b} \in \mathbb{R}^m$ , we will use  $\overrightarrow{(\mathbf{a}(i)\mathbf{b}(i))}$  to denote the vector  $\mathbf{z}$  with  $\mathbf{z}(i) = \mathbf{a}(i)\mathbf{b}(i)$  and so on.

Throughout this chapter, when we are working in the context of some given graph  $G$  with vertices  $V$  and edges  $E$ , we will let  $m = |E|$  and  $n = |V|$ .

The plots in this chapter were made using Mathematica, which is available to ETH students for download through the ETH IT Shop.

### 16.1 An Interior Point Method

**The Maximum Flow problem in undirected graphs.**

$$\begin{aligned} & \max_{\mathbf{f} \in \mathbb{R}^E} F \\ \text{s.t. } & \mathbf{Bf} = F\mathbf{b}_{st} && \text{“The Undirected Maximum Flow Problem”} \\ & -\mathbf{c} \leq \mathbf{f} \leq \mathbf{c} \end{aligned} \tag{16.1}$$

We use  $\text{val}(\mathbf{f})$  to denote  $F$  when  $B\mathbf{f} = F\mathbf{b}_{st}$ .

As we develop algorithms for this problem, we will assume that we know the maximum flow value  $F^*$ . Let  $\mathbf{f}^*$  denote some maximum flow, i.e. a flow with  $-\mathbf{c} \leq \mathbf{f} \leq \mathbf{c}$  can  $\text{val}(\mathbf{f}^*) = F^*$ .

In general, an a lower bound  $F \leq F^*$  will allow us to find a flow with value  $F$ , and because of this, we can use a binary search to approximate  $F^*$ .

### 16.1.1 A Barrier Function and an Algorithm

$$V(\mathbf{f}) = \sum_e -\log(\mathbf{c}(e) - \mathbf{f}(e)) - \log(\mathbf{c}(e) + \mathbf{f}(e))$$

We assume the optimal value of Program (16.1) is  $F^*$ . Then for a given  $0 \leq \alpha < 1$  we define a program

$$\begin{aligned} \min_{\mathbf{f} \in \mathbb{R}^E} \quad & V(\mathbf{f}) \\ \text{s.t. } \mathbf{Bf} = \alpha F^* \mathbf{b}_{st} \quad & \text{“The Barrier Problem”} \end{aligned} \tag{16.2}$$

This problem makes sense for any  $0 \leq \alpha < 1$ . When  $\alpha = 0$ , we are not routing any flow yet. This will be our starting point. For any  $0 \leq \alpha < 1$ , the scaled-down maximum flow  $\alpha \mathbf{f}^*$  strictly satisfies the capacities  $-\mathbf{c} < \alpha \mathbf{f}^* < \mathbf{c}$ , and  $\mathbf{B}\alpha \mathbf{f}^* = \alpha F^* \mathbf{b}_{st}$ . Hence  $\alpha \mathbf{f}^*$  is a feasible flow for this value of  $\alpha$  and hence  $V(\alpha \mathbf{f}^*) < \infty$  and so the optimal flow for the Barrier Problem at this  $\alpha$  must also have objective value strictly below  $\infty$ , and hence in turn strictly satisfy the capacity constraints. Thus, if we can find the optimal flow for Program (16.2) for  $\alpha = 1 - \epsilon$ , we will have a feasible flow with Program (16.1), the Undirected Maximum Flow Problem, routing  $(1 - \epsilon)F^*$ . This is how we will develop an algorithm for computing the maximum flow.

Program (16.2) has the Lagrangian

$$\mathcal{L}(\mathbf{f}, \mathbf{x}) = V(\mathbf{f}) + \mathbf{x}^\top (\alpha F^* \mathbf{b}_{st} - \mathbf{Bf})$$

And we have optimality when

$$\mathbf{Bf} = \alpha F^* \mathbf{b}_{st} \text{ and } -\mathbf{c} \leq \mathbf{f} \leq \mathbf{c} \tag{16.3}$$

“Barrier feasibility”

and  $\nabla_{\mathbf{f}} \mathcal{L}(\mathbf{f}, \mathbf{x}) = \mathbf{0}$ , i.e.

$$\nabla V(\mathbf{f}) = \mathbf{B}^\top \mathbf{x} \tag{16.4}$$

” Barrier Lagrangian gradient optimality”

Let  $\mathbf{f}_\alpha^*$  denote the optimal solution to Problem 16.2 for a given  $0 \leq \alpha < 1$ , and let  $\mathbf{x}_\alpha^*$  be optimal dual voltages such that  $\nabla V(\mathbf{f}_\alpha^*) = \mathbf{B}^\top \mathbf{x}_\alpha^*$ .

It turns out that, if we have a solution  $\mathbf{f}_\alpha^*$  to this problem for some  $\alpha < 1$ , then we can find a solution  $\mathbf{f}_{\alpha+\alpha'}$  for some  $\alpha' < 1 - \alpha$ . And, we can compute  $\mathbf{f}_{\alpha+\alpha'}$  using a small number of Newton steps, each of which will only require a Laplacian linear equation solve, and hence is computable in  $\tilde{O}(m)$  time. Concretely, for any  $0 \leq \alpha < 1$ , given the optimal flow at this  $\alpha$ , we will be able to compute the optimal flow at  $\alpha_{\text{new}} = \alpha + (1 - \alpha) \frac{1}{150\sqrt{m}}$ . This means that after  $T = 150\sqrt{m} \log(1/\epsilon)$  updates, we have a solution for  $\alpha \geq 1 - \epsilon$ .

We can state the update problem as

$$\begin{aligned} \min_{\boldsymbol{\delta} \in \mathbb{R}^E} \quad & V(\boldsymbol{\delta} + \mathbf{f}) \\ \text{s.t. } \mathbf{B}\boldsymbol{\delta} = \alpha' F^* \mathbf{b}_{st} \end{aligned} \tag{16.5}$$

“The Update Problem”

### 16.1.2 Updates using Divergence

It turns out that for the purposes of analysis, it will be useful to ensure that our “Update Problem” uses an objective function that is minimized at  $\boldsymbol{\delta} = \mathbf{0}$ .

This leads to a variant of the Update Problem, which we call the “Divergence Update Problem”. We obtain our new problem by switching from  $V(\boldsymbol{\delta} + \mathbf{f})$  as our objective to  $V(\boldsymbol{\delta} + \mathbf{f}) - (V(\mathbf{f}) + \langle \nabla V(\mathbf{f}), \boldsymbol{\delta} \rangle)$  as our objective, and this is called the *divergece* of  $V$  w.r.t.  $\boldsymbol{\delta}$  based at  $\mathbf{f}$ .

$$\begin{aligned} \min_{\boldsymbol{\delta} \in \mathbb{R}^E} \quad & V(\boldsymbol{\delta} + \mathbf{f}) - (V(\mathbf{f}) + \langle \nabla V(\mathbf{f}), \boldsymbol{\delta} \rangle) \\ \text{s.t. } \mathbf{B}\boldsymbol{\delta} = \alpha' F^* \mathbf{b}_{st} \end{aligned} \tag{16.6}$$

“The Divergence Update Problem”

Now, for any flow  $\boldsymbol{\delta}$  such that  $\mathbf{B}\boldsymbol{\delta} = \alpha' F^* \mathbf{b}_{st}$ , using the Lagrangian gradient condition (16.4), we have  $\langle \nabla V(\mathbf{f}_\alpha^*), \boldsymbol{\delta} \rangle = \langle \mathbf{x}_\alpha^*, \alpha' F^* \mathbf{b}_{st} \rangle$ . Hence, for such  $\boldsymbol{\delta}$ , we have

$$V(\boldsymbol{\delta} + \mathbf{f}_\alpha^*) - (V(\mathbf{f}_\alpha^*) + \langle \nabla V(\mathbf{f}_\alpha^*), \boldsymbol{\delta} \rangle) = V(\boldsymbol{\delta} + \mathbf{f}_\alpha^*) - (V(\mathbf{f}_\alpha^*) + \langle \mathbf{x}_\alpha^*, \alpha' F^* \mathbf{b}_{st} \rangle)$$

We conclude that the objectives of the Update Problem (16.5) and the Divergence Update Problem (16.6) have the same minimizer, which we denote  $\boldsymbol{\delta}_{\alpha'}^*$ , although, to be precise, it is also a function of  $\alpha$ .

Thus  $\mathbf{f}_\alpha^* + \boldsymbol{\delta}_{\alpha'}^*$  is optimal for the optimization problem

$$\begin{aligned} \min_{\mathbf{f} \in \mathbb{R}^E} \quad & V(\mathbf{f}) \\ \text{s.t. } \mathbf{B}\mathbf{f} = (\alpha + \alpha') F^* \mathbf{b}_{st} \end{aligned} \tag{16.7}$$

**Lemma 16.1.1.** *Suppose  $\mathbf{f}_\alpha^*$  is the minimizer of Problem (16.2) (the Barrier Problem with parameter  $\alpha$ ) and  $\boldsymbol{\delta}_{\alpha'}^*$  is the minimizer of Problem (16.6) (the Update Problem with parameters  $\mathbf{f}_\alpha^*$  and  $\alpha'$ ), then  $\mathbf{f}_\alpha^* + \boldsymbol{\delta}_{\alpha'}^*$  is optimal for Problem (16.2) with parameter  $\alpha + \alpha'$  (i.e. a new instance of the Barrier problem).*

---

**Algorithm 11: INTERIOR POINT METHOD**


---

```

1  $\mathbf{f} \leftarrow \mathbf{0}$ ;
2  $\alpha \leftarrow 0$ ;
3 while  $\alpha < 1 - \epsilon$  do
4    $a' \leftarrow \frac{1-\alpha}{20\sqrt{m}}$ ;
5   Compute  $\boldsymbol{\delta}$ , the minimizer of Problem (16.6);
6   Let  $\mathbf{f} \leftarrow \mathbf{f} + \boldsymbol{\delta}$  and  $\alpha \leftarrow \alpha + a'$ ;
7 end
8 return  $\mathbf{f}$ 

```

---

**Pseudotheorem 16.1.1.** *Let  $\mathbf{f}$  be the minimizer of Problem (16.2). Then, when  $a' \leq \frac{1-\alpha}{20\sqrt{m}}$ , the minimizer  $\boldsymbol{\delta}$  of Problem (16.7) can be computed in  $\tilde{O}(m)$  time.*

The key insight in this type of interior point method is that when the update  $\alpha'$  is small enough,

**Theorem 16.1.2.** *Algorithm 11 returns a flow  $\mathbf{f}$  that is feasible for Problem (16.1) in time  $\tilde{O}(m^{1.5} \log(1/\epsilon))$ .*

*Proof Sketch.* First note that for  $\alpha = 0$ , the minimizer of Problem (16.2) is  $\mathbf{f} = \mathbf{0}$ . The proof now essentially follows by Lemma (16.1.1), and Pseudotheorem 16.1.1. Note that  $1 - \alpha$  shrinks by a factor  $(1 - \frac{1}{20\sqrt{m}})$  in each iteration of the while-loop, and so after  $20\sqrt{m} \log(1/\epsilon)$  iterations, we have  $1 - \alpha \leq \epsilon$ , at which point the loop terminates. To turn this into a formal proof, we need to take care of the fact the proper theorem corresponding to Pseudotheorem 16.1.1 only gives a highly accurate but not exact solution  $\delta$  to the “Update Problem”. But it’s possible to show that this is good enough (even though both  $\mathbf{f}$  and  $\boldsymbol{\delta}$  end up not being exactly optimal in each iteration).  $\square$

**Remark 16.1.3.** For the maximum flow problem, when capacities are integral and polynomially bounded, if we choose  $\epsilon = m^{-c}$  for some large enough constant  $c$ , given a feasible flow with  $\text{val}(\mathbf{f}) = 1 - \epsilon$ , is it possible to compute an exact maximum flow in nearly linear time. Thus Theorem 16.1.2 can also be used to compute an exact maximum flow in  $\tilde{O}(m)$  time, but we omit the proof. The idea is to first round to an almost optimal, feasible integral flow (which requires a non-trivial combinatorial algorithm), and then to recover the exact flow using Ford-Fulkerson. See [Mad13] for details.

**Remark 16.1.4.** It is possible to reduce an instance of directed maximum flow to an instance of undirected maximum flow in nearly-linear time, in such a way that if we can *exactly* solve the undirected instance, then in nearly-linear time we can recover an exact solution to the directed maximum flow problem. Thus Theorem (16.1.2) can also be used to solve directed maximum flow. We will ask you to develop this reduction in Graded Homework 2.

**Remark 16.1.5.** For sparse graphs with  $m = \tilde{O}(n)$  and large capacities, this running time is the best known, and improving it is major open problem.

### 16.1.3 Understanding the Divergence Objective

Note that if  $V(x) = -\log(1 - x)$ , then  $D(x) = V(x) - (V(0) + V'(0)x)$ .

```
Plot[{Log[1 / (1 - x)], x}, {x, -1, 1}]
```

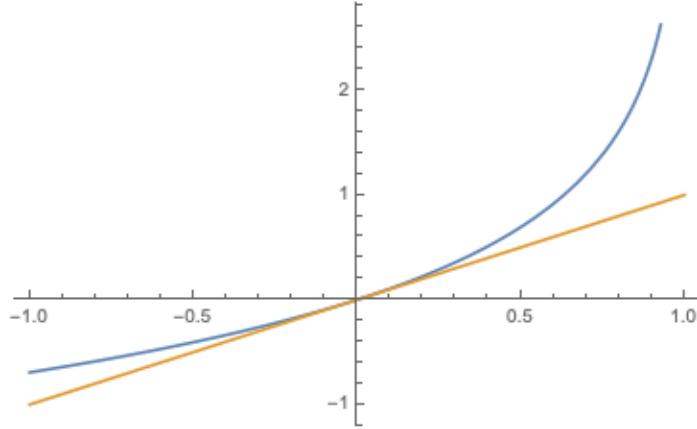


Figure 16.1: Plot showing  $V(x) = -\log(1 - x)$  and then linear approximation  $V(0) + V'(0)x$ .

```
Plot[Log[1 / (1 - x)] - x, {x, -1, 1}]
```

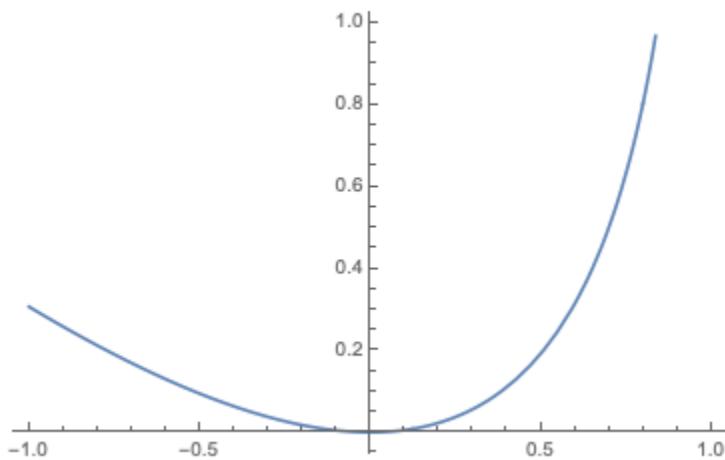


Figure 16.2: Plot showing  $D(x) = V(x) - (V(0) + V'(0)x)$ .

We let

$$\mathbf{c}_+(e) = \mathbf{c}(e) - \mathbf{f}(e) \text{ and } \mathbf{c}_-(e) = \mathbf{c}(e) + \mathbf{f}(e)$$

So then

$$\begin{aligned}
D_V(\boldsymbol{\delta}) &= V(\boldsymbol{\delta} + \mathbf{f}) - (V(\mathbf{f}) + \langle \nabla V(\mathbf{f}), \boldsymbol{\delta} \rangle) \\
&= \sum_e -\log \left( \frac{\mathbf{c}(e) - (\boldsymbol{\delta}(e) + \mathbf{f}(e))}{\mathbf{c}(e) - \mathbf{f}(e)} \right) - \frac{\boldsymbol{\delta}(e)}{\mathbf{c}(e) - \mathbf{f}(e)} \\
&\quad - \log \left( \frac{\mathbf{c}(e) + (\boldsymbol{\delta}(e) + \mathbf{f}(e))}{\mathbf{c}(e) + \mathbf{f}(e)} \right) + \frac{\boldsymbol{\delta}(e)}{\mathbf{c}(e) + \mathbf{f}(e)} \\
&= \sum_e D \left( \frac{\boldsymbol{\delta}(e)}{\mathbf{c}(e) - \mathbf{f}(e)} \right) + D \left( -\frac{\boldsymbol{\delta}(e)}{\mathbf{c}(e) + \mathbf{f}(e)} \right) \\
&= \sum_e D \left( \frac{\boldsymbol{\delta}(e)}{\mathbf{c}_+(e)} \right) + D \left( -\frac{\boldsymbol{\delta}(e)}{\mathbf{c}_-(e)} \right)
\end{aligned}$$

Note that we can express Problem (16.6) as

$$\begin{aligned}
\min_{\boldsymbol{\delta} \in \mathbb{R}^E} \quad & D_V(\boldsymbol{\delta}) \\
\text{s.t. } \mathbf{B}\boldsymbol{\delta} = \alpha' F^* \mathbf{b}_{st} \quad & \text{The Update Problem, restated}
\end{aligned} \tag{16.8}$$

Note that  $D_V(\boldsymbol{\delta})$  is strictly convex over the feasible set, so the argmin is unique.

#### 16.1.4 Quadratically Smoothing Divergence and Local Agreement

$$\tilde{D}_\epsilon(x) = \begin{cases} -\log(1-x) - x & \text{if } |x| \leq \epsilon \\ D(\epsilon) + D'(\epsilon)(x-\epsilon) + \frac{D''(\epsilon)}{2}(x-\epsilon)^2 & \text{if } x > \epsilon \\ D(-\epsilon) + D'(-\epsilon)(x+\epsilon) + \frac{D''(-\epsilon)}{2}(x+\epsilon)^2 & \text{if } x < -\epsilon \end{cases}$$

For brevity, we define

$$\tilde{D}(x) = \tilde{D}_{0.1}(x)$$

**Lemma 16.1.6.**

1.  $1/2 \leq \tilde{D}''(\mathbf{x}) \leq 2$ .
2. For  $x \geq 0$ , we have  $x/2 \leq \tilde{D}'(\mathbf{x}) \leq 2x$  and  $-2x \leq \tilde{D}'(-\mathbf{x}) \leq -x/2$ .
3.  $x^2/4 \leq \tilde{D}(\mathbf{x}) \leq x^2$ .

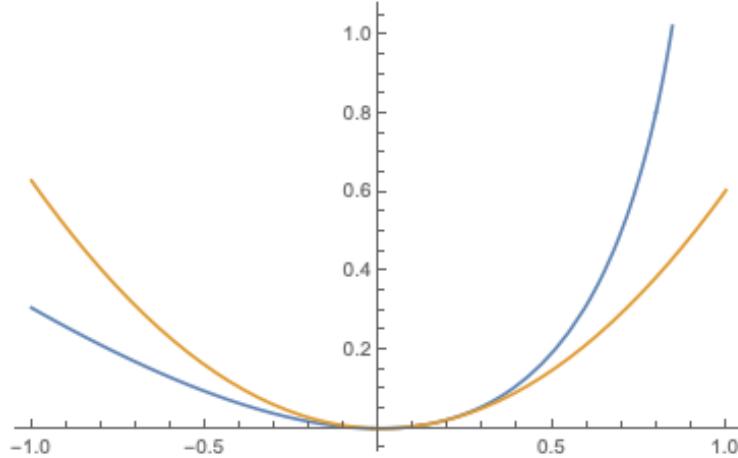
What's happening here? We glue together  $D(x)$  for small  $x$  with its quadratic approximation for  $|x| > \epsilon$ . For  $x > \epsilon$ , we "glue in" a Taylor series expansion based at  $x = \epsilon$ .

```
NumberForm[Series[Log[1/(1-x)] - x, {x, 0.1, 2}], 3]
```

NumberForm=

$$0.00536 + 0.111(x - 0.1) + 0.617(x - 0.1)^2 + O[(x - 0.1)^3]$$

```
Plot[{Log[1/(1-x)] - x, 0.00536 + 0.111(x - 0.1) + 0.617(x - 0.1)^2}, {x, -1, 1}]
```



```
Plot[{Log[1/(1-x)] - x, 0.00536 + 0.111(x - 0.1) + 0.617(x - 0.1)^2}, {x, 0, 0.2}]
```

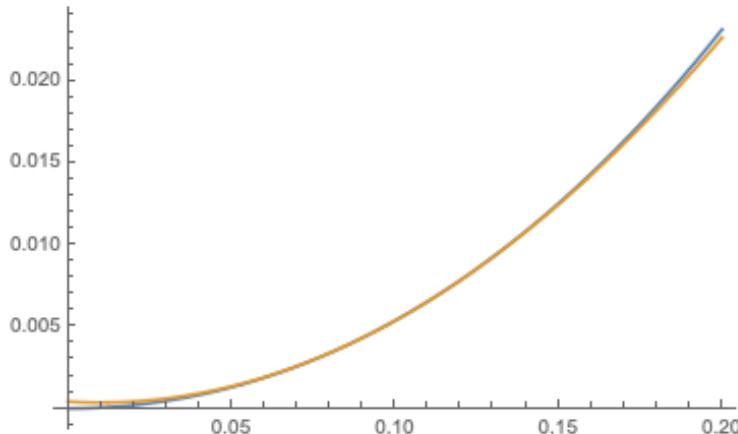


Figure 16.3: Plot showing  $D(x) = -\log(1-x)$  and the quadratic approximation based at  $x = 0.1$ .

We also define

$$\tilde{D}_V(\boldsymbol{\delta}) = \sum_e \tilde{D}\left(\frac{\boldsymbol{\delta}(e)}{\mathbf{c}_+(e)}\right) + \tilde{D}\left(-\frac{\boldsymbol{\delta}(e)}{\mathbf{c}_-(e)}\right)$$

We can now introduce the smoothed optimization problem

$$\begin{aligned} & \min_{\boldsymbol{\delta} \in \mathbb{R}^E} \quad \tilde{D}_V(\boldsymbol{\delta}) \\ \text{s.t. } & \mathbf{B}\boldsymbol{\delta} = \boldsymbol{\alpha}' F^* \mathbf{b}_{st} \end{aligned} \tag{16.9}$$

“The Smoothed Update Problem”

Note that  $\tilde{D}_V(\boldsymbol{\delta})$  is strictly convex over the feasible set, so the argmin is unique.

**Pseudoclaim 16.1.7.** *We can compute the argmin  $\boldsymbol{\delta}^*$  of Problem (16.9), the Smoothed Update Problem, using the Newton-Steps for K-stable Hessian convex functions that we saw in the previous chapter, in  $\tilde{O}(m)$  time.*

*Sketch of proof.* Problem (16.9) fits the class of problems for which we showed in the previous chapter that (appropriately scaled) Newton steps converge. This is true because the Hessian is always a 2-spectral approximation of the Hessian at  $\tilde{D}_V(\boldsymbol{\delta}^*)$ , as can be shown from Lemma 16.1.6. Because the Hessian of  $\tilde{D}_V(\boldsymbol{\delta})$  is diagonal, and the constraints are flow constraints, each Newton step boils down to solving a Laplacian linear system, which can be done to high accuracy  $\tilde{O}(m)$  time.  $\square$

**Remark 16.1.8.** There are three things we need to modify to turn the pseudoclaim into a true claim, addressing the errors arising from both Laplacian solvers and Newton steps:

1. We need to rephrase the claim to so that we only claim  $\boldsymbol{\delta}^*$  has been computed to high accuracy, rather than exactly.
2. We need to show that we can construct an initial guess to start off Newton's method  $\boldsymbol{\delta}_0$  for which the value  $\tilde{D}_V(\boldsymbol{\delta}_0)$  is not too large. (This is easy).
3. We need show that Newton steps converge despite using a Laplacian solver that doesn't give exact solutions, only high accuracy solutions. (Takes a bit of work, but is ultimately not too difficult).

Importantly, to ensure our overall interior point method still works, we also need to show that it converges, even if we're using approximate solutions everywhere. This also takes some work to show, again is not too difficult.

### Local Agreement Implies Same Optimum.

**Lemma 16.1.9.** *Suppose  $S \subseteq \mathbb{R}^n$  is a convex set, and let  $f, g : S \rightarrow \mathbb{R}$  be convex functions. Let  $\mathbf{x}^* = \arg \min_{\mathbf{x} \in S} f(\mathbf{x})$ . Suppose  $f, g$  agree on a neighborhood of  $\mathbf{x}^*$  in  $S$  (i.e. an open set containing  $\mathbf{x}^*$ ). Then  $\mathbf{x}^* = \arg \min_{\mathbf{x} \in S} g(\mathbf{x})$ .*

*Proof Sketch.* We sketch the proof in the case when both  $f, g$  are differentiable: Observe that  $\mathbf{0} = \nabla f(\mathbf{x}^*) = \nabla g(\mathbf{x}^*)$ , and hence  $g(\mathbf{x})$  is also minimized at  $\mathbf{x}^*$ .  $\square$

We define

$$\hat{\mathbf{c}}(e) = \min(\mathbf{c}_+(e), \mathbf{c}_-(e)) \quad (16.10)$$

**Lemma 16.1.10.** *Suppose  $\boldsymbol{\delta}^*$  is the argmin of Problem (16.9), the Smoothed Update Problem, and  $\left\| \overrightarrow{(\boldsymbol{\delta}^*(e)/\hat{\mathbf{c}}(e))} \right\|_\infty < 0.1$ . Then  $\boldsymbol{\delta}^*$  is the argmin of Problem (16.8).*

*Proof.* We observe that if  $\left\| \overrightarrow{(\delta^*(e)/\widehat{\mathbf{c}}(e))} \right\|_\infty < 0.1$ , then  $\tilde{D}_V(\delta^*) = D_V(\delta^*)$ , and, for all  $\tau \in \mathbb{R}^m$  with norm

$$\left\| \overrightarrow{(\tau(e)/\widehat{\mathbf{c}}(e))} \right\|_\infty < 0.1 - \left\| \overrightarrow{(\delta^*(e)/\widehat{\mathbf{c}}(e))} \right\|_\infty$$

we have that  $\tilde{D}_V(\delta^* + \tau) = D_V(\delta^* + \tau)$ . Thus  $\tilde{D}_V$  and  $D_V$  agree on a neighborhood around  $\delta^*$  and hence by Lemma 16.1.9, we have that  $\delta^*$  is the argmin of Problem (16.8).  $\square$

### 16.1.5 Step size for divergence update

**Definition 16.1.11** (*s-t well-conditioned* graph). An undirected, capacitated multi-graph  $G = (V, E, \mathbf{c})$  with source  $s$  and sink  $t$  is *s-t well-conditioned* if, letting  $U$  denote the maximum edge capacity  $U = \|\mathbf{c}\|_\infty$ , we have at least  $\frac{2}{3}m$  multi-edges of capacity  $U$  going directly from  $s$  to  $t$ .

**Remark 16.1.12.** It is straightforward to make a graph *s-t* well-conditioned. We just add  $2m$  new edges of capacity  $U$  directly between  $s$  and  $t$ . Given an exact maximum flow in the new graph, it is trivial to get one in the original graph: Just remove the flow on the new edges.

**Definition 16.1.13.** Given a *directed* graph  $G = (V, E, \mathbf{c})$ , the *symmetrization* of  $G$  is the undirected  $\widehat{G} = (V, \widehat{E}, \widehat{\mathbf{c}})$  is the undirected graph given by

$$\{a, b\} \in \widehat{E} \text{ if } (a, b) \in E \text{ AND } (b, a) \in E$$

and

$$\widehat{\mathbf{c}}(\{a, b\}) = \min(\mathbf{c}(a, b), \mathbf{c}(b, a)).$$

Note that when  $\widehat{G}_f$  is the symmetrization of the residual graph  $G_f$  (which we defined in Chapter 11), then  $\widehat{\mathbf{c}}$  matches exactly the definition of  $\widehat{\mathbf{c}}$  in Equation (16.10).

**Lemma 16.1.14.** Let  $G$  be an undirected, capacitated multi-graph  $G = (V, E, \mathbf{c})$  which is *s-t well-conditioned*. Let  $\mathbf{f}$  be the minimizer of Program (16.2). Let  $\widehat{G}_f$  be the symmetrization of the residual graph  $G_f$  (in the sense of Lecture 10). Then there exists a flow  $\widehat{\delta}$  which satisfies  $\mathbf{B}\widehat{\delta} = \frac{1-\alpha}{5}F^*\mathbf{b}_{st}$  and is feasible in  $\widehat{G}_f$ . Note that we can also state the feasibility in  $\widehat{G}_f$  as

$$\left\| \overrightarrow{(\widehat{\delta}(e)/\widehat{\mathbf{c}}(e))} \right\|_\infty \leq 1$$

*Proof.* We recall since  $\mathbf{f}$  is the minimizer of Program (16.2), there exists dual-optimal voltages  $\mathbf{x}$  such that

$$\mathbf{B}^\top \mathbf{x} = \nabla V(\mathbf{f}) = \overrightarrow{\left( \frac{1}{\mathbf{c}(e) - \mathbf{f}(e)} - \frac{1}{\mathbf{c}(e) + \mathbf{f}(e)} \right)}$$

From Lecture 10, we know that there is flow  $\bar{\delta}$  that is feasible with respect to the residual graph capacities of the graph  $G_f$  such that  $\mathbf{B}\bar{\delta} = (1 - \alpha)F^*\mathbf{b}_{st}$ . Note when treating  $\bar{\delta}$

as an undirected flow, feasibility in the residual graph means that  $\bar{\delta}(e) < \mathbf{c}(e) - \mathbf{f}(e)$  and  $-\bar{\delta}(e) < \mathbf{c}(e) + \mathbf{f}(e)$ . Thus,

$$(1 - \alpha)F^* \mathbf{b}_{st}^\top \mathbf{x} = \bar{\delta} \mathbf{B}^\top \mathbf{x} = \sum_e \frac{\bar{\delta}}{\mathbf{c}(e) - \mathbf{f}(e)} - \frac{\bar{\delta}}{\mathbf{c}(e) + \mathbf{f}(e)} \leq m$$

Now, because the graph is  $s$ - $t$  well-conditioned, there are at  $\frac{2}{3}m$  edges directly from  $s$  to  $t$  with capacity  $U$  and each of these  $e$  satisfy by the Lagrangian gradient optimality condition (16.4)

$$\mathbf{b}_{st}^\top \mathbf{x} = \frac{1}{U - \mathbf{f}(e)} - \frac{1}{U + \mathbf{f}(e)}$$

Note that  $\frac{2}{3}mU \leq F^* \leq mU$  because the graph is  $s$ - $t$  well-conditioned. To complete the analysis, we consider three cases.

*Case 1:*  $|\mathbf{f}(e)| \leq \frac{2}{3}U$ . Then the capacity on each of these edges in the symmetrized residual graph  $\widehat{G}_f$  is at least  $U/3$ . As there are  $\frac{2}{3}m$  of them, we get that there is a feasible flow in  $\widehat{G}_f$  of value at least  $\frac{2}{9}mU \geq \frac{1}{10}F^*$ .

*Case 2:*  $\mathbf{f}(e) < -\frac{2}{3}U$ . By the gradient condition, we have the same flow on all of the  $\frac{2}{3}m$   $s$ - $t$  edges, adding up to at least  $\frac{2}{3}mU$  going from  $t$  to  $s$ . This means that we must have at least  $\frac{2}{3}mU$  flow going from  $s$  to  $t$  via the remaining edges. But, their combined capacity is at most  $\frac{1}{3}mU$ , so that cannot happen. Thus we can rule out this case entirely.

*Case 3:*  $\mathbf{f}(e) > \frac{2}{3}U$ . Then

$$\frac{m}{(1 - \alpha)F^*} \geq \mathbf{b}_{st}^\top \mathbf{x} \geq \frac{1}{U - \mathbf{f}(e)} - \frac{1}{U + \mathbf{f}(e)} \geq \frac{4/5}{U - \mathbf{f}(e)}$$

So

$$U - \mathbf{f}(e) \geq \frac{4}{5} \frac{(1 - \alpha)F^*}{m} \geq \frac{1}{2}(1 - \alpha)U$$

In this case, the capacity on each of the  $\frac{2}{3}m$   $s$ - $t$  edges with capacity  $U$  in  $G$  will have capacity  $(1 - \alpha)U/2$  in  $\widehat{G}_f$ . This guarantees that there is feasible flow in  $\widehat{G}_f$  of value at least  $\frac{1}{3}(1 - \alpha)mU \geq \frac{1}{3}(1 - \alpha)F^*$ .  $\square$

**Lemma 16.1.15.** *Let  $0 < \alpha' \leq \frac{1-\alpha}{150\sqrt{m}}$ . Then the minimizer  $\delta^*$  of Problem (16.9) satisfies  $\left\| \overrightarrow{(\delta^*(e)/\widehat{c}(e))} \right\|_\infty < 0.1$ .*

*Proof.* By Lemma 16.1.14, there exists a flow  $\widehat{\delta}$  which satisfies  $\mathbf{B}\widehat{\delta} = \frac{1-\alpha}{5}F^*\mathbf{b}_{st}$  and  $\left\| \overrightarrow{(\widehat{\delta}(e)/\widehat{c}(e))} \right\|_\infty \leq 1$ . Hence for any  $0 < \alpha' \leq \frac{1-\alpha}{150\sqrt{m}}$ , the flow  $\tilde{\delta} = \alpha'\frac{5}{1-\alpha}\widehat{\delta}$  satisfies

$\mathbf{B}\tilde{\boldsymbol{\delta}} = \alpha' F^* \mathbf{b}_{st}$  and  $\left\| \overrightarrow{(\tilde{\boldsymbol{\delta}}(e)/\widehat{\mathbf{c}}(e))} \right\|_\infty \leq \frac{1}{30\sqrt{m}}$ . This means that

$$\begin{aligned}\tilde{D}_V(\tilde{\boldsymbol{\delta}}) &= \sum_e \tilde{D} \left( \frac{\tilde{\boldsymbol{\delta}}(e)}{\mathbf{c}_+(e)} \right) + \tilde{D} \left( -\frac{\tilde{\boldsymbol{\delta}}(e)}{\mathbf{c}_-(e)} \right) \\ &\leq \sum_e 4 \left( \frac{\tilde{\boldsymbol{\delta}}(e)}{\mathbf{c}_+(e)} \right)^2 + 4 \left( -\frac{\tilde{\boldsymbol{\delta}}(e)}{\mathbf{c}_-(e)} \right)^2 \\ &\leq \sum_e 8 \left( \frac{\tilde{\boldsymbol{\delta}}(e)}{\widehat{\mathbf{c}}(e)} \right)^2 \\ &\leq 8/900 < 1/100.\end{aligned}$$

This then means that the minimizer  $\boldsymbol{\delta}^*$  of Problem (16.9) also satisfies  $\tilde{D}_V(\tilde{\boldsymbol{\delta}}) < 1/100$ .

$$\begin{aligned}\left\| \overrightarrow{(\boldsymbol{\delta}^*/\widehat{\mathbf{c}}(e))} \right\|_\infty^2 &\leq \sum_e \left( \frac{\boldsymbol{\delta}^*(e)}{\mathbf{c}_+(e)} \right)^2 + \left( -\frac{\boldsymbol{\delta}^*(e)}{\mathbf{c}_-(e)} \right)^2 \\ &\leq \sum_e \tilde{D} \left( \frac{\boldsymbol{\delta}^*(e)}{\mathbf{c}_+(e)} \right) + \tilde{D} \left( -\frac{\boldsymbol{\delta}^*(e)}{\mathbf{c}_-(e)} \right) \quad \text{By Lemma 16.1.6.} \\ &= \tilde{D}_V(\tilde{\boldsymbol{\delta}}) < 1/100.\end{aligned}$$

Hence  $\left\| \overrightarrow{(\boldsymbol{\delta}^*/\widehat{\mathbf{c}}(e))} \right\|_\infty < 0.1$ . □

## Part V

# Further Topics in Combinatorial Graph Algorithms

# Chapter 17

## Low-Diameter Decompositions

In this chapter, we learn about an algorithm that clusters vertices in graphs by their distance. This clustering is often referred to as *low-diameter decomposition*. We define this concept below and give one of the most versatile and robust algorithms.

**Definition 17.0.1.** Given an undirected graph  $G = (V, E, w)$  with non-negative weights  $w \in \mathbb{R}_{\geq 0}$ . We say that a partition  $\mathcal{X}$  of  $V$  and a set of edges  $E_{del}$  form a  $(D, \alpha)$ -low-diameter decomposition of  $G$  if

1. for every  $i \in [1, k]$  and  $u, v \in X \in \mathcal{X}$ , we have  $\text{dist}_{G[X]}(u, v) \leq \alpha D$ , and
2.  $\forall e \in E$ ,  $\mathbb{P}[e \in E_{del}] \leq w(e)/D$ , and
3.  $E(G/\mathcal{X}) \subseteq E_{del}$ .

Above we use  $G/\mathcal{X}$  to denote the graph obtained from  $G$  by contracting all vertices in the same set in  $\mathcal{X}$  into a super-vertex and removing all self-loops. Thus the final property expresses that  $E_{del}$  contains all edges that cross clusters in  $\{X_1, X_2, \dots, X_k\}$ .

**Remark 17.0.2.** The definition above is with respect to an implicit distribution  $\mathcal{D}$  over partition sets  $\mathcal{X}$ . A more precise statement would be to say that for a distribution  $\mathcal{D}$  over partition sets of  $V$ , we say that  $\mathcal{D}$  is an  $(D, \alpha)$ -LDD-distribution if sampling a partition from  $\mathcal{D}$  yields the properties stated above. We will refrain from that here, however.

**Remark 17.0.3.** Instead of the second property, LDDs are often defined to have the weight of edges in  $E_{del}$  to be a small fraction of the edges. For unweighted graphs, the best one can hope for is essentially at most  $m/D$  intercluster edges  $E_{del}$  in total. The above definition allows us to obtain this property in expectation. Defined this way, we also have already seen a reasonable algorithm to obtain such a clustering:  $\phi$ -expander decompositions have few intercluster edges and  $\phi$ -expanders have diameter  $O(\log(n)/\phi)$  as you have shown in an exercise.

However, for applications, the probabilistic guarantee is very useful and algorithms to find LDDs are generally much simpler than computing expander decompositions. LDDs have been used as a subroutine in many algorithms recently, we will see a truly exciting application of LDDs in the next chapter. In the first half of this lecture, we analyze a simple algorithm to obtain an LDD. In the second half of the lecture, we extend LDDs and our algorithm to directed graphs.

## 17.1 Low-Diameter Decomposition in Undirected Graphs

The main result of this half of the lecture is the following.

**Theorem 17.1.1.** *Given an undirected graph  $G = (V, E, \mathbf{w})$  with  $\mathbf{w} \in \mathbb{R}_{\geq 0}$  and  $D > 0$ , there is an algorithm  $\text{COMPUTELDD}(G, D)$  that computes an  $(D, 8 \log(n))$ -low-diameter-decomposition in time  $O(m \log n)$  with probability at least  $1 - n^{-3}$ .*

**Exponential Distribution.** We denote by  $\text{EXPDISTR}(D)$  a function that samples a real according to the exponential distribution with parameter  $D$ . Formally, we sample from a distribution with cumulative distribution function

$$F(x, D) = \begin{cases} 1 - e^{-x/D} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Note that we have the property that  $F(D \cdot 4 \log n, D) = 1 - e^{-4 \log n} = 1 - 1/n^4$ . A really cool property of the exponential distribution is further that it is memoryless which means that for any randomly sampled integer  $X \sim \text{EXPDISTR}(D)$  and fixed positive reals  $x, s$ , we have  $\mathbb{P}[X > x + s | X > s] = \mathbb{P}[X > x]$ .

**A Simple LDD Algorithm.** Below, we present a simple algorithm to compute a low-diameter decomposition  $\mathcal{X}$ . The algorithm simply selects an arbitrary vertex  $r$  in the graph, then chooses a random radius  $\tilde{D}$  and adds the ball  $B_G(r, \tilde{D})$  as a cluster to  $\mathcal{X}$ . The process is then repeated on the graph  $G[V \setminus \mathcal{X}]$  (where we abuse notation and really mean  $G[V \setminus (\cup_{X \in \mathcal{X}} X)]$ ), i.e. we repeat on the graph induced by all non-clustered vertices. Clearly, this outputs a partition of the vertex set  $V$  in the end.

**Analysis.** We now analyze the algorithm.

**Claim 17.1.2.** *Given a clustering  $\mathcal{X}$  returned by Algorithm 12. We have with probability  $\geq 1 - n^{-3}$ , that for any  $u, v \in X \in \mathcal{X}$ ,  $\text{dist}_{G[X]}(u, v) \leq 8D \log(n)$ .*

---

**Algorithm 12:** COMPUTELDD( $G, D$ )

---

```

1  $\mathcal{X} = \{\}$ .
2 while  $\mathcal{X}$  is not a full partition of  $V$  do
3   Let  $r$  be an arbitrary vertex in  $V \setminus \mathcal{X}$ .
4    $\tilde{D} \leftarrow \text{EXPDISTR}(D)$ .
5    $\mathcal{X} \leftarrow \mathcal{X} \cup \{B_G(r, \tilde{D})\}$ .
6 return  $(\mathcal{X}, E_{\text{del}} = \{(u, v) \in E \mid u \in X \in \mathcal{X}, v \in Y \in \mathcal{X}, X \neq Y\})$ 

```

---

*Proof.* From our discussion of the exponential distribution, we have by a simple union bound over at most  $n$  iterations of the algorithm, that none of the sampled radii  $\tilde{D}$  exceeds  $4D \log(n)$  with probability  $1 - n^{-3}$ . Let us condition on this event.

Next, consider any iteration, where initially we have  $\mathcal{X}$  and then we find a new vertex  $r \in G[V \setminus \mathcal{X}]$  and a ball  $X = B_G(r, \tilde{D})$  that is added to  $\mathcal{X}$  at the end of the iteration. We clearly have that  $G[X]$  includes a path from each vertex  $u \in X$  to  $r$  of length at most  $4D \log(n)$ . But then by the triangle inequality, for every pair  $u, v \in X$ ,  $\text{dist}_{G[X]}(u, v) \leq \text{dist}_{G[X]}(u, r) + \text{dist}_{G[X]}(r, v) \leq 8D \log(n)$ .  $\square$

**Claim 17.1.3.** *Given a clustering  $\mathcal{X}$  returned by Algorithm 12.  $\forall e = (u, v) \in E$ ,  $\mathbb{P}[e \in E_{\text{del}}] \leq w(e)/D$ .*

*Detailed Proof Sketch.* Consider the first iteration where either  $u$  or  $v$  is in a ball  $B_G(r, \tilde{D})$ . By assumption, we have that in this iteration  $\tilde{D} \geq \min_{x \in \{u, v\}} \text{dist}_G(r, x)$ , i.e. we condition on the event  $\mathcal{E}$  that  $\tilde{D}$  satisfies this property.

Now, we have that  $e$  becomes intercluster and thus is in  $E_{\text{del}}$  if and only if one of the endpoints is not in the ball, i.e. that  $\tilde{D} < \max_{x \in \{u, v\}} \text{dist}_G(r, x)$ . We obtain the following inequalities

$$\begin{aligned}
& \mathbb{P}[e \text{ is an intercluster edge} \mid \mathcal{E}] \\
&= \mathbb{P}[\tilde{D} < \max_{x \in \{u, v\}} \text{dist}_G(r, x) \mid \mathcal{E}] \\
&\leq \mathbb{P}[\tilde{D} < \min_{x \in \{u, v\}} \text{dist}_G(r, x) + w(e) \mid \tilde{D} \geq \min_{x \in \{u, v\}} \text{dist}_G(r, x)] \\
&= \mathbb{P}[\tilde{D} < w(e)] \\
&= F(w(e), D) - F(0, D) = (1 - e^{-w(e)/D}) - (1 - e^0) = 1 - e^{-w(e)/D} \\
&\leq 1 - (1 - w(e)/D) = w(e)/D.
\end{aligned}$$

where we use the triangle inequality, that yields  $\max_{x \in \{u, v\}} \text{dist}_G(r, x) \leq \min_{x \in \{u, v\}} \text{dist}_G(r, x) + w(e)$ . Thus, the event that  $\tilde{D} < \min_{x \in \{u, v\}} \text{dist}_G(r, x) + w(e)$  contains the event that  $e$  becomes intercluster, and it suffices to upper bound its probability. We then exploit the memoryless property, the definition of the exponential distribution, and then use  $1 + x \leq e^x$  to simplify.

Note that we only analyzed the probability conditioned on  $\mathcal{E}$ . A more rigorous proof has to argue that in the event  $\neg\mathcal{E}$ , the probability of  $e$  becoming intercluster is 0. This requires some more careful set-up.  $\square$

It remains to bound the runtime to obtain Theorem 17.1.1. To this end, observe that every iteration can be executed running Dijkstra's algorithm from the vertex  $r$  on the induced graph where vertices are only relaxed if their distance estimate drops below  $\tilde{D}$ . It is not hard to show that each iteration then takes time  $O(|E(B_G(r, \tilde{D}))| \log(n))$  and since these balls are vertex disjoint, we have that the total runtime is at most  $O(m \log n)$ , as desired.

## 17.2 Generalizing LDDs to Directed Graphs

For our most exciting application, we need to generalize LDDs to directed graphs. We will start out with a new definition, and then argue about why this is a reasonable generalization (other generalizations exist but this one allows to obtain LDDs efficiently and turns out very useful in applications).

**Definition 17.2.1.** Given a **directed graph**  $G = (V, E, w)$  with non-negative weights  $w \in \mathbb{R}_{\geq 0}$ , we say that a partition  $X_1, X_2, \dots, X_k$  of  $V$  and a set of edges  $E_{del} \subseteq E$  form a  $(D, \alpha)$ -low-diameter decomposition of  $G$  if

1. for every  $i \in [1, k]$  and  $u, v \in X_i$ , we have  $\text{dist}_G(u, v) \leq \alpha D$ , and
2.  $\forall e \in E$ ,  $\mathbb{P}[e \in E_{del}] \leq O(w(e) \log n / D)$ .
3. the graph  $(G / \{X_1, X_2, \dots, X_k\}) \setminus E_{del}$ , that is the graph where each cluster  $X_i$  is contracted into a super-vertex, is a directed acyclic graph (DAG).

Note in particular the subtle difference in Property 1 that clusters have small distances in  $G$  and not in the graph induced by the cluster. In Property 2, we gave the probability an  $O(\log n)$  slack. While this makes the definition slightly unclean, it makes it easier to state our result.

The third property might initially look a bit strange. Let us briefly argue that it is a reasonable property to have. Note first, that in directed graphs, we might have a highly asymmetric relationship in distances, i.e.  $\text{dist}_G(u, v) \gg \text{dist}_G(v, u)$  for some, or even all the vertices. One example that really stresses this is the cycle graph. While it is very easy to decompose the undirected cycle graph, the directed cycle graph only allows for decomposition into singletons if one insists on having clusters of small diameter  $D$ . What the third property expresses is that in this case, while our clusters are not very meaningful, we learn something about the structure of  $G$ : removing a few edges (in this case a single one) makes the directed graph acyclic.

In the next lecture, we will see an application of this subroutine in an algorithm that basically does a case split for these two extreme cases: it exploits the small distances in clusters to

make progress on the clusters, and it exploits the structural closeness to a DAG on the remainder of the graph.

The rest of this lecture is concerned with proving the following theorem.

**Theorem 17.2.2.** *Given a **directed** graph  $G = (V, E, \mathbf{w})$  with  $\mathbf{w} \in \mathbb{R}_{\geq 0}$ , there is an algorithm COMPUTEDIRECTEDLDD( $G, D$ ) that computes an  $(D, 8 \log(n))$ -low-diameter-decomposition in time  $\tilde{O}(m)$  with probability at least  $1 - n^{-3}$ .*

We first give a slightly simpler algorithm that is not quite efficient but allows us to generalize the ideas from the undirected setting cleanly. We then show how to adapt the algorithm to make it efficient.

**A Directed LDD Algorithm.** Let us show how to implement the algorithm. Again, we are simply carving out balls, one after another, until we know that all remaining distances have small distance to each other. An important detail in this algorithm is that while the ball-carving is executed on the graph  $G[V \setminus \mathcal{X}]$ , which is the graph with all carved out vertices removed, the eligibility criterion of the while-loop is with respect to the initial input graph  $G$ . We also point out that  $|V(G)|$  here might be different from  $n$ , as  $n$  denotes the size of the global input graph, while  $|V(G)|$  denotes the size of the input graph to the current procedure. This is important for recursive calls.

---

**Algorithm 13:** COMPUTEDIRECTEDLDD( $G, D$ )

---

```

1  $\mathcal{X} = \{\}; E_{del} \leftarrow \emptyset.$ 
2 while there is a vertex  $r \in V \setminus \mathcal{X}$  with  $|B_G^{out}(r, 4D \log n)| \leq \frac{1}{2}|V(G)|$  or
    $|B_G^{in}(r, 4D \log n)| \leq \frac{1}{2}|V(G)|$  do
3    $\tilde{D} \leftarrow \text{EXPDISTR}(D).$ 
4    $\# \leftarrow \underset{\#\in\{\text{out,in}\}}{\text{argmin}} |B_G^{\#}(r, \tilde{D})|.$ 
5    $X \leftarrow B_{G[V \setminus \mathcal{X}]}^{\#}(r, \tilde{D}).$ 
6   Add all  $\#$ -edges of  $X$  in  $G[V \setminus \mathcal{X}]$  to  $E_{del}.$ 
7    $(\mathcal{X}', E'_{del}) \leftarrow \text{COMPUTEDIRECTEDLDD}(G[X], D);$  // Recurse.
8    $\mathcal{X} \leftarrow \mathcal{X} \cup \{X\} \cup \mathcal{X}'; E_{del} \leftarrow E_{del} \cup E'_{del}.$ 
9    $\mathcal{X} \leftarrow \mathcal{X} \cup \{V \setminus \mathcal{X}\};$  // Add new cluster.
10 return  $(\mathcal{X}, E_{del})$ 

```

---

**Analysis.** Here, we thoroughly discuss the above algorithm but only sketch how to implement it efficiently. A rigorous exposition of these arguments can be found at [here<sup>1</sup>](#).

---

<sup>1</sup>Lecture Notes by Danupon Nanongkai: <https://hackmd.io/@U0nm1XUhREKPYLt1eUmE6g/Sycpovkiq>. The notes only present the arguments for obtaining a weak diameter guarantee for each cluster, however, this can be salvaged by slightly adapting the algorithm and analysis.

**Claim 17.2.3.** *Given a clustering  $\mathcal{X}$  returned by Algorithm 13. We have for any  $u, v \in X \in \mathcal{X}$ ,  $\text{dist}_G(u, v) \leq 8D \log(n)$ .*

*Proof.* Note that a new cluster is only added to  $\mathcal{X}$  in Line 11 (all other clusters stem from recursive calls that themselves find clusters by adding them in this line). Consider any two vertices  $u, v$  in this new cluster, we have that  $u$  reaches more than half the vertices in  $G$  in its out-ball of radius  $4D \log n$  and  $v$  reaches more than half the vertices of  $G$  in its in-ball of the same radius. Thus, these balls have an intersection and there is a path of weight at most  $8D \log n$ .  $\square$

**Claim 17.2.4.** *Given a clustering  $\mathcal{X}$  returned by Algorithm 12.  $\forall e = (u, v) \in E, \mathbb{P}[e \in E_{\text{del}}] = O(w(e) \log n / D)$ .*

*Proof Sketch.* It is not hard to see that Algorithm 12 only calls itself recursively on disjoint vertex subsets of at most half the size of the original graph (technically, this is only true when conditioned on  $\tilde{D}$  never exceeding  $4D \log n$ , but we will be imprecise here). It is therefore straightforward to see that each vertex  $u$  participates in at most  $O(\log n)$  calls.

The remainder of the analysis is analogous to the undirected case. Here, we exploit that an out-ball only adds the boundary edges leaving the out-ball to  $E_{\text{rem}}$  and an in-ball only the boundary edges entering the in-ball.  $\square$

**Claim 17.2.5.** *Given a clustering  $\mathcal{X}$  and edge set  $E_{\text{del}}$  returned by Algorithm 12. We have that  $(G/\mathcal{X}) \setminus E_{\text{del}}$  is a DAG.*

*Proof.* Let us first prove that the algorithm maintains the invariant that there is no cycle in  $G \setminus E_{\text{del}}$  that visits more than one cluster in  $\mathcal{X}$ .

We prove by induction on the size of  $V(G)$ . If  $|V(G)| \leq 1$ , the while loop is not entered and the algorithm returns  $\mathcal{X} = \{V\}$ .

For larger sizes of  $V(G)$ , whenever we remove a set  $X$  from  $V \setminus \mathcal{X}$  to recurse on  $G[X]$ , we add either all in- or out-edges of  $X$  to  $E_{\text{del}}$ . By the induction hypothesis, the recursive call returns a partition  $\mathcal{X}'$  of  $X$  such that no cycle in  $G[X] \setminus E'_{\text{del}}$  visits any two clusters. But this means that any such cycle that does, must use an in- **and** an out-edge of  $X$ , but we remove either all in- or all out-edges of  $X$ . Similarly, no future subcluster of  $V \setminus (\mathcal{X} \cup \{X\} \cup \mathcal{X}')$  can have a cycle with a cluster in  $\{X\} \cup \mathcal{X}'$  as again it requires an in- and an out-edge from  $X$ .

Finally, each cluster in  $\mathcal{X}$  has small diameter by Theorem 17.2.3 and thus, when we contract sets  $\mathcal{X}$ , we do not create new cycles. Thus  $(G/\mathcal{X}) \setminus E_{\text{del}}$  is DAG.  $\square$

This completes the proof of correctness.

**An Efficient Implementation.** Finally, we want to argue about runtime. The bottleneck in Algorithm 13 is the computation of the information in the while-loop condition. Unfortunately, it is not clear how to obtain this information in reasonable time.

---

**Algorithm 14:** COMPUTEDIRECTEDLDD( $G, D$ )

---

```

1  $\mathcal{X} = \{\}; E_{del} \leftarrow \emptyset.$ 
2 Sample each vertex  $v \in V$  into set  $S$  with probability  $p = \frac{324 \log(n)}{|V(G)|}$ .
3 Compute in- and out-SSSP for all  $s \in S$ .
4 while there is a vertex  $r \in V \setminus \mathcal{X}$  with  $|B_G^{out}(r, 4D \log n) \cap S| \leq \frac{2}{3}|V(G)|/p$  or
    $|B_G^{in}(r, 4D \log n) \cap S| \leq \frac{2}{3}|V(G)|/p$  do
5    $\tilde{D} \leftarrow \text{EXPDISTR}(D).$ 
6    $\# \leftarrow \underset{\# \in \{out, in\}}{\operatorname{argmin}} |B_G^{\#}(r, 4D \log n) \cap S|.$ 
7    $X \leftarrow B_G^{\#}(r, \tilde{D}).$ 
8   Add all  $\#$ -edges of  $X$  in  $G[V \setminus \mathcal{X}]$  to  $E_{del}$ .
9    $(\mathcal{X}', E'_{del}) \leftarrow \text{COMPUTEDIRECTEDLDD}(G[X], D);$  // Recurse.
10   $\mathcal{X} \leftarrow \mathcal{X} \cup \{X\} \cup \mathcal{X}'; E_{del} \leftarrow E_{del} \cup E'_{del}.$ 
11  $\mathcal{X} \leftarrow \mathcal{X} \cup \{V \setminus \mathcal{X}\};$  // Add new cluster.
12 return  $(\mathcal{X}, E_{del})$ 

```

---

Instead, we resort to weakening the while-loop condition slightly to obtain an efficient algorithm: we initially sample a set  $S$  consisting of  $\Theta(\log n)$  vertices uniformly at random. We then run Dijkstra from each vertex in  $S$  on the graph  $G$  and the reverse graph. This allows us to obtain for each vertex  $v \in V$ , the information how large the intersection of  $S$  is with its out- and in-ball of radius  $4D \log n$ . This now gives an estimator that w.h.p. can detect whether a vertex  $v$  has an out- and in-ball of size more than  $\frac{1}{2}|V(G)|$  or otherwise guarantees that both balls are of size at most  $\frac{5}{6}|V(G)|$ . We make this more precise with the claim below.

**Claim 17.2.6.** *For any call of Algorithm 14, with probability at least  $1 - n^{-2}$ , the sampling of  $S$  is such that for every vertex  $r \in V$  and  $\# \in \{in, out\}$ , we have that if  $|B_G^{\#}(r, 4D \log n) \cap S| \leq \frac{2}{3}|V(G)|/p$ , then  $|B_G^{\#}(r, 4D \log n)| \leq \frac{5}{6}|V(G)|$ . And if  $|B_G^{\#}(r, 4D \log n) \cap S| \geq \frac{2}{3}|V(G)|/p$ , then  $|B_G^{\#}(r, 4D \log n)| \geq \frac{1}{2}|V(G)|$ .*

*Proof.* By Chernoff, we have that for independent random variables  $Y_1, Y_2, \dots, Y_k \in \{0, 1\}$  with  $\mathbb{E}[Y_i] = p$  for all  $i$ , we have  $\mathbb{P}[|\sum_{i=1, \dots, k} Y_i - pk| \geq \delta pk] \leq 2e^{-\delta^2 pk/3}$  if  $0 < \delta < 1$ .

In our case, we choose  $\delta = \frac{1}{6}$  which yields a failure probability of our claim of at most  $2e^{-(1/6)^2 p|V(G)|/3} = 2e^{-p|V(G)|/108} = 2e^{-3 \log(n)} < n^{-2}$ .  $\square$

Thus, it now suffices to check whether a  $\frac{2}{3}$ -fraction of the vertices in  $S$  intersect which yields a good estimator on these two values.

Since balls in graph  $G[V \setminus \mathcal{X}]$  are subsets of balls in graph  $G$  (for the same center and radius), we have that every recursive call is indeed on a graph that has the number of vertices decreased by a constant fraction compared to the original graph. For the remaining vertices that are added as a cluster  $V \setminus \mathcal{X}$  in the current call, we have that they reach

half the vertices in  $G$  in in- and out-ball w.h.p. and thus, we can use an argument as in Theorem 17.2.3, we have that they are at distance at most  $8D \log n$ .

It is not hard to bound via the above guarantees the number of calls that every vertex participates in by  $O(\log n)$  which yields the two other properties of LDDs. And further, since sampling is efficient, and SSSP computations take time near-linear in the number of edges of the input graph, we can bound the runtime now by  $\tilde{O}(m)$ . We leave it to the reader to verify this bound.

# Chapter 18

## Negative Single-Source Shortest Paths

In this lecture, we are interested in the following result. It says that even in graphs with negative edge weights, we can find the single-source shortest paths in near-linear time! This result has been obtained extremely recently and the algorithm is a beautiful combination of some of the most useful modern graph algorithmic techniques.

**Theorem 18.0.1.** *Given a directed graph  $G = (V, E, \mathbf{w})$  with  $\mathbf{w} \in \mathbb{R}^V$  such that there exists no negative cycle in  $G$ , and a dedicated source vertex  $s \in V$ , the distances  $\text{dist}_G(s, t)$  can be computed in expected time  $\tilde{O}(m \log(W))$  where  $W = \|\mathbf{w}\|_\infty$ , i.e.  $W$  is the largest absolute weight (and where we assume that the smallest absolute non-zero weight is at least 1).*

Henceforth, we assume that any graph  $G$  under consideration does not contain a negative weight cycle.

### 18.1 A Not-So-Brief Recap: Classic Algorithms for Handling Negative Weights

**Bellman-Ford.** The classic algorithm that should come to mind when computing single-source shortest paths with negative weights is Bellman-Ford's algorithm. The algorithm works by initializing distance estimates  $\hat{d}(s, u)$  with infinity and  $\hat{d}(s, s)$  with 0. It then proceeds for  $n$  rounds by relaxing the edge set in every round.

The following theorem summarizes the guarantees of the algorithm.

**Theorem 18.1.1.** *Algorithm 15 takes as input a directed, weighted graph  $G = (V, E, \mathbf{w})$  where weights are allowed to be negative, but that contains no negative cycle; and a source vertex  $s \in V$ . It runs in time  $O(mn)$  and outputs the correct distances from  $s$  to every vertex  $v \in V$ .*

*Proof.* The algorithm runs for  $n$  iterations and relaxes  $m$  edges in each iteration, each in  $O(1)$  time. This establishes the runtime.

---

**Algorithm 15:** BELLMANFORD( $G, s$ )

---

```

1  $\forall u \in V \setminus \{s\}, \hat{d}(s, u) \leftarrow \infty; \hat{d}(s, s) \leftarrow 0.$ 
2 for  $i = 1, 2, \dots, n$  do
3   foreach  $e = (u, v) \in E$  do
4      $\hat{d}(s, v) \leftarrow \min\{\hat{d}(s, v), \hat{d}(s, u) + w(e)\}$ 
5 return  $\{\hat{d}(s, v)\}_v$ 

```

---

We prove correctness by induction. We prove that after the  $h$ -th iteration of the algorithm, for all vertices  $v$  that have a shortest path consisting of at most  $h$  edges have  $\hat{d}(s, v)$  is equal to the real  $s$  to  $v$  distance.

The base case  $h = 0$  is trivial because  $\hat{d}(s, s)$  is initialized to 0. For the inductive step  $h \mapsto h + 1$ , we use the subpath property of shortest paths, that is, for any  $v$  that has a shortest  $s$  to  $v$  shortest path consisting of exactly  $h + 1$  edges, we have that for the vertex  $u$  preceding  $v$  on this path, we have that the subpath from  $s$  to  $u$  is a shortest path itself and has  $< h + 1$  edges (it is easy to establish the subpath property via proof by contradiction). Thus in the  $h + 1$ -th iteration, when the edge  $e = (u, v)$  is relaxed, we have  $\hat{d}(s, v) \leftarrow \hat{d}(s, u) + w(e) = \text{dist}_G(s, u) + w(e) = \text{dist}_G(s, v)$ .  $\square$

**A Heuristic for Speeding-Up Bellman-Ford.** If you think a bit longer about Bellman-Ford and why you cannot use Dijkstra, you might stumble upon the following idea that might improve runtime in practice: consider that there is a shortest path from  $s$  to  $v$  that consists of at most  $h$  negative edges. Then, we could try to relax negative edges by Bellman-Ford, and relax all non-negative edges between two consecutive edges on the shortest path using Dijkstra's algorithm. We know that after  $h$  rounds of this algorithm,  $v$ 's distance estimate becomes the distance and we therefore never need to relax after round  $h + 1$  from  $v$ .

Let us make that intuition precise.

**Definition 18.1.2.** Given a directed, weighted graph  $G = (V, E, \mathbf{w})$  with negative weights but no negative cycle and a source  $s \in V$ . For any vertex  $v \in V$ , we let  $\eta(s, v)$  be the minimum number of negative edges over all shortest  $s$  to  $v$  paths.

The pseudo-code below now formalizes our intuitive approach. It reduces the runtime to  $O((\sum_v \eta(s, v) \cdot \deg^{\text{out}}(v) + m) \log n)$  which is always upper bounded by  $O(nm \log n)$ .

We can now prove the following theorem which proves that the algorithm above is faster at least for the special case when  $\eta(G)$  is much smaller than  $n$ .

**Theorem 18.1.3.** *Algorithm HEURISTICBELLMANFORD( $G, s$ ) runs in time  $O((\sum_v \eta(s, v) \cdot \deg^{\text{out}}(v) + m) \log n)$  time and outputs the correct  $s$  to  $v$  distances.*

*Proof Sketch.* The runtime is trivially analyzed. To avoid a tedious case analysis, let us assume that all out-going edges of  $s$  are negative (this can be enforced w.l.o.g.).

---

**Algorithm 16:** HEURISTICBELLMANFORD( $G, s$ )

---

```

1  $\forall u \in V \setminus \{s\}, \hat{d}(s, u) \leftarrow \infty; \hat{d}(s, s) \leftarrow 0.$ 
2 for  $i = 1, 2, \dots, n$  do
3   /* Relaxation via Bellman-Ford. */ *
4   foreach  $e = (u, v) \in E$  where  $\hat{d}(u)$  was decreased last iteration do
5      $\hat{d}(s, v) \leftarrow \min\{\hat{d}(s, v), \hat{d}(s, u) + w(e)\}.$ 
6   /* Relaxation via Dijkstra's Algorithm. */ *
7   Let  $Q$  be the collection of all vertices  $u \in V$  where  $\hat{d}(s, u)$  was decreased in the
8   current for-loop iteration.
9   while  $Q \neq \emptyset$  do
10    Extract the vertex  $u$  from  $Q$  that has smallest  $\hat{d}(s, u).$ 
11    foreach  $e = (u, v) \in E, w(e) \geq 0$  do
12      if  $\hat{d}(s, v) < \hat{d}(s, u) + w(e)$  then
13         $\hat{d}(s, v) \leftarrow \hat{d}(s, u) + w(e);$  add  $v$  to  $Q.$ 
14
15 return  $\{\hat{d}(s, v)\}_v$ 

```

---

Then correctness follows by induction on  $\eta(s, v)$ , i.e. after the  $h$ -th iteration of the algorithm, we have for all vertices  $v \in V$  with  $\eta(s, v) \leq h$ , that  $\hat{d}(s, v) = \text{dist}_G(s, v)$ . The base case is again trivial. For the inductive step  $h \mapsto h + 1$ , we have that for any  $v$  with  $\eta(s, v) = h + 1$ , by the subpath property, we have that the vertex  $u$  just before the last negative edge on the shortest  $s$  to  $v$  path of  $h + 1$  edges, that  $\hat{d}(u)$  is equal to the real distance after iteration  $h$ . But then  $u$  is relaxed in the round after learning the real distance by Bellman-Ford which then decreases the distance of the next vertex on the path. Dijkstra's algorithm then relaxes all non-negative edges in the correct order on the path which leaves  $v$ 's estimate decreased to the real distance.  $\square$

**All-Pairs Shortest Paths in Negative-Weight Graphs via Price Function.** While we are at reviewing classic algorithm, let us also briefly review Johnson's algorithms. This algorithm allows to compute APSP in directed graphs even when weights are negative (but again no negative cycles).

The key idea used in Johnson's algorithm is the concept of a price function.

**Definition 18.1.4** (Price Function). For any vector  $\phi \in \mathbb{R}$ , we say  $\phi$  is a *price function* if  $\mathbf{w}_\phi = \mathbf{w} + \mathbf{B}^\top \phi \geq \mathbf{0}$ , i.e. where for each edge  $e = (u, v) \in E$  directed from  $u$  to  $v$ , we have that  $\mathbf{w}(e) + \phi(u) - \phi(v) \geq 0$ . We write  $G_\phi$  to mean the graph  $(V, E, \mathbf{w}_\phi)$ .

A price function is a nice tool since once one obtains a price function, one can run Dijkstra's algorithm on the graph  $G_\phi$  and extract the real distances efficiently. Let us briefly prove that.

**Claim 18.1.5.** *Given a directed graph  $G = (V, E, \mathbf{w})$  and any function  $\phi$  over the vertices. Then we have for any vertices  $s, t \in V$  that  $\text{dist}_G(s, t) = \text{dist}_{G_\phi}(s, t) + \phi(s) - \phi(t)$ .*

*Proof.* For any  $s$  to  $t$  path  $P$  in  $G$ , we have that  $\mathbf{w}_\phi(P) = \sum_{e=(u,v) \in \pi} \mathbf{w}_\phi(e) = \sum_{e=(u,v) \in \pi} \mathbf{w}(e) + \phi(u) - \phi(v) = \mathbf{w}(\pi) + \sum_{e=(u,v) \in \pi} \phi(u) - \phi(v) = \mathbf{w}(\pi) = \phi(s) - \phi(t)$ . But this means that all  $s$  to  $t$  paths are additively increased by  $\phi(s) - \phi(t)$  thus the shortest path in  $G$  w.r.t.  $\mathbf{w}$  is still a shortest path in  $G_\phi$ .  $\square$

To see that price functions even exist consider the following algorithm which computes it via Bellman-Ford.

---

**Algorithm 17:** COMPUTEPRICEFUNCTION( $G$ )

---

- 1 Let  $G'$  be the graph  $G$  with an additional dummy source  $d$ , and edges  $(d, v)$  of weight 0 for every  $v \in V$ .
  - 2 Run HEURISTICBELLMANFORD( $G', d$ ).
  - 3 **return**  $\phi = \{\text{dist}_{G'}(d, v)\}_{v \in V}$
- 

**Claim 18.1.6.** *Algorithm 17 returns a price function  $\phi$  for  $G$ . It runs in time  $O((\sum_v \eta(d, v) \cdot \deg^{\text{out}}(v) + m) \log n)$ .*

*Proof.* In fact, we can prove that  $\phi$  extended by the value  $\phi(d) = 0$  for the dummy source  $d$  is a price function for  $G'$ .

To see this, observe that the distance from  $d$  to any vertex  $v$  in  $G$  weighted by  $\mathbf{w}_\phi$  is 0 by definition. Thus, the shortest path from  $d$  to any vertex  $v$  is of weight exactly 0 w.r.t.  $\phi$ . But assume that there is an edge  $e = (x, v)$  with  $\mathbf{w}_\phi(e) < 0$ . Then, the path obtained from the shortest path from  $d$  to  $x$  appended with edge  $(x, v)$  is of negative weight. But this contradicts that  $v$  is at distance 0 from  $d$ .  $\square$

Once a price function is found, one can obtain All-Pairs Shortest Paths by running Dijkstra's algorithm from each vertex on the graph  $G_\phi$ . The runtime of this algorithm is  $\tilde{O}(mn)$ .

## 18.2 Scaling for the Negative-Weight SSSP Problem

Next, we use a technique that is often referred to as *scaling*. It appeared first in the context of obtaining max-flow algorithms. The idea is that we can pre-condition/ normalize a problem and then only try to improve the situation slightly instead of completely solving the problem. To be precise, we want to prove the following statement.

**Theorem 18.2.1.** *Given an algorithm PRICEFUNCTIONONNORMALIZEDGRAPH( $G$ ) that may assume that the input graph  $G = (V, E, \mathbf{w})$  has  $\mathbf{w} \geq -2 \cdot \mathbf{1}$ , and that then computes*

a function  $\phi$  such that  $G_\phi$  has its most negative weight being at least  $-1$  and runs in time  $\mathcal{T}(m, n)$ .

Then, there exists an algorithm NEGATIVESSSP( $G, s$ ) for any graph  $G = (V, E, \mathbf{w})$  that solves the negative SSSP problem in time  $O(\mathcal{T}(m, n) \log(W_{\max}/W_{\min}))$  where  $W_{\max}$  is the largest absolute weight in  $G$  and  $W_{\min}$  is the smallest absolute non-zero weight.

Using scaling initially, we can and will assume henceforth that initially  $W_{\min} \geq 1$  (one simply needs to scale by  $1/W_{\min}$  to make this true).

**A Scaling Algorithm.** Let us now present such an algorithm. The idea is that one can initially scale the graph down to satisfy the condition that the smallest weight is at least  $-2$ , and then using the produced "price" function, scale up by a factor 2 in the next round when inputting the graph reweighted by this first price function. The algorithm below proposes to keep going with this process until one is at a precision where one can round the "price function" into a real price function  $\phi$  for  $G$ . It then remains to run Dijkstra's algorithm.

---

**Algorithm 18:** NEGATIVESSSP( $G, s$ )

---

```

1  $W_{\max} \leftarrow \|\mathbf{w}\|_\infty$ .
2  $\phi_0 \leftarrow \mathbf{0}$ .
3 for  $i = 1, 2, \dots, \eta = \lceil \log_2(W_{\max}n^2) \rceil$  do
4    $\phi_i \leftarrow \text{PRICEFUNCTIONONNORMALIZEDGRAPH} \left( (G \cdot 2^i / \lceil \log_2(W_{\max}) \rceil)_{\phi_{i-1}} \right)$ .
5    $\phi_{\text{almost}} \leftarrow \frac{\lceil \log_2(W_{\max}) \rceil}{\lceil \log_2(W_{\max}n^2) \rceil} \cdot \phi_\eta$ .
6 Obtain a price function  $\phi$  from  $\phi_{\text{almost}}$  by rounding.
7 return DIJKSTRA( $G, s, \phi$ ).

```

---

Here, we will not give proof that this algorithm indeed yields Theorem 18.2.1, but simply use its result. For proof, also that  $\phi$  can be computed efficiently in the second-to-last line of the algorithm, we refer you to the original paper.

### 18.3 The (Normalized) Negative-Weight SSSP Problem

Recall, we can now assume that the most negative weight in  $G$  is  $-2$ , and we want to produce a "price" function  $\phi$  such that the most negative weight in  $G_\phi$  is  $-1$ . For the rest of this section, we focus on giving an implementation of PRICEFUNCTIONONNORMALIZEDGRAPH() with runtime  $\tilde{O}(m\sqrt{m})$ .

**Constant-Degree Assumption.** It will be convenient to assume that  $G$  has constant-degree. This is w.l.o.g. since we can take an arbitrary graph  $G$  and transform it into a

metric-preserving graph  $G'$  with  $O(m)$  vertices and edges by replacing each vertex  $v$  with  $d$  edges incident, by a zero-weight cycle with  $d$  vertices where each edge on the cycle receives one of the edges originally incident to  $v$ .

**Notation.** We introduce the notation  $G^{+i}$  to denote the graph  $(V, E, \mathbf{w} + i \cdot \mathbf{1})$ , i.e. the graph  $G$  after adding  $i$  to each edge weight. We will use

- $G^{+0}$  which is the initial graph;
- $G^{+1}$  which is the graph for which  $\phi$  will be a real price function; and
- $G^{+2}$  which has the property that every edge weight is non-negative.

**The Algorithm.** To obtain Single-Source Shortest Paths from a dedicated source  $s$ , we use the following algorithm. The algorithm first decomposes  $G$  into LDDs where it ignores edge weights and uses weights 1 instead. It then computes a price function for each cluster  $X_i$  in the LDD. Finally, combining these price functions in the most straight-forward into a global function  $\phi$  (that is not a price function but as we will see is close to a price function), we can invoke our heuristic for Bellman-Ford on the graph  $G_\phi$ , and then return the computed distances after adjusting them for the price function.

---

**Algorithm 19:** PRICEFUNCTIONONNORMALIZEDGRAPH( $G = (V, E, \mathbf{w})$ )

---

- 1  $(\mathcal{X} = \{X_1, X_2, \dots, X_k\}, E_{del}) \leftarrow \text{COMPUTELDD}(G^{+2}, \sqrt{n})$ .  
// Find "price" function  $\phi_1$  on clusters.
  - 2 **for**  $i = 1, 2, \dots, k$  **do**  $\phi_1^i \leftarrow \text{COMPUTEPRICEFUNCTION}(G^{+1}[X_i])$ . ;
  - 3  $\phi_1 \leftarrow [\phi_1^1, \phi_1^2, \dots, \phi_1^k]^\top$ .  
// Adapt "price" function to  $\phi_2$  which also works for DAG edges.
  - 4 Let  $\pi$  be a topological order of  $X_1, X_2, \dots, X_k$  in  $(G/\mathcal{X}) \setminus E_{del}$ .
  - 5 **for**  $i = 1, 2, \dots, k$  **do**  $\phi_2^i \leftarrow \phi_1^i + (2\|\phi_1\|_\infty + 1) \cdot \pi(X_1) \cdot \mathbf{1}_{X_1}$ . ;
  - 6  $\phi_2 \leftarrow [\phi_2^1, \phi_2^2, \dots, \phi_2^k]^\top$ .  
// Obtain "price" function  $\phi_3$  that works for all of  $G^{+1}$ .
  - 7  $\phi_3 \leftarrow \text{COMPUTEPRICEFUNCTION}(G_{\phi_2}^{+1})$ .
  - 8 **return**  $\phi_2 + \phi_3$ .
- 

**Analysis.** From the algorithm, it is immediately clear that  $\phi_2 + \phi_3$  is a price function for  $G^{+1}$  and thus satisfies our claim. It remains to bound the runtime.

**Claim 18.3.1.** *For each cluster  $X_i \in \mathcal{X}$ , we have that  $\text{COMPUTEPRICEFUNCTION}(G^{+1}[X_i])$  runs in time  $\tilde{O}(n\sqrt{n})$ .*

*Proof.* Recall that  $\text{COMPUTEPRICEFUNKTION}(G^{+1}[X_i])$  internally adds a dummy source  $d$  and edges of weight 0 from  $d$  to everyone in  $G^{+1}[X_i]$ , and then computes the price function in time  $\tilde{O}(\sum_v \eta(d, v))$  (where we use that  $G$  is assumed constant-degree).

Let us now show that  $\eta(d, v) = \tilde{O}(\sqrt{n})$  for all  $v$  which then establishes the claim. We note that  $\eta(\cdot, \cdot)$  is here defined w.r.t.  $G^{+1}[X]$  (plus the vertex  $d$ ). Assume for the sake of contradiction that  $\eta(d, v) \geq 8\sqrt{n} \log n + 1$ . Since there is an edge  $(d, v)$  of weight 0 for all  $v \in X_i$  in this graph, we have that the weight of the shortest  $d$  to  $v$  path  $P$  must be of non-positive weight in  $G^{+1}$ . Let  $u$  be the first vertex on this shortest path after  $d$ . Then, the shortest  $u$  to  $v$  path  $P'$  is exactly the path  $P$  omitting the first vertex, and thus also of weight at most 0. Note that  $P'$  is contained in  $G[X_i]$  (without  $d$ ) and consists of at least  $8\sqrt{n} \log n$  edges.

Finally, note that in  $G^{+0}$  each edge carries weight one less than in  $G^{+1}$ . So the path  $P'$  is of weight at most  $-8\sqrt{n} \log n$  in  $G = G^{+0}$ , thus  $\text{dist}_G(u, v) \leq -8\sqrt{n} \log n$ . But at the same time, we know from the LDD decomposition (see Theorem 17.2.2) that since  $u$  and  $v$  are in the same cluster w.r.t.  $G^{+2}$  that  $\text{dist}_G(v, u) < \text{dist}_{G^{+2}}(v, u) \leq 8\sqrt{n} \log n$ .

But this implies a negative cycle in  $G$  which yields a contradiction.  $\square$

To analyze the runtime of the call to  $\text{COMPUTEPRICEFUNKTION}(G_{\phi_2}^{+1})$  in Line 7, we establish that  $\phi_2$  indeed is a price function for  $G^{+1} \setminus E_{\text{del}}$ .

**Claim 18.3.2.**  $\phi_2$  is a price function for  $G^{+1} \setminus E_{\text{del}}$ .

*Proof.* Since price functions are invariant to shifts by 1, it is not hard to verify that for every edge  $e$  in  $G[X_i]$ , we have that  $\mathbf{w}_{\phi_2}(e) = \mathbf{w}_{\phi_1}(e) \geq -1$  and thus in  $G^{+1}$ , the function is a price function restricted to edges in clusters.

Let us now consider any edge  $e = (u, v)$  in  $G \setminus E_{\text{del}}$  that is not in a cluster. Let  $u \in X_i$  and  $v \in X_j$  for  $i \neq j$ . We have by definition, that  $\pi(X_i) + 1 \leq \pi(X_j)$ . But  $|\phi_1(v) - \phi_1(u)| \leq 2\|\phi_1\|_\infty$  and in  $G$  the smallest possible edge weight is  $-2$  by assumption. Thus,

$$\begin{aligned} \mathbf{w}_{\phi_2}(e) &= \mathbf{w}(e) + \phi_1(u) - \phi_1(v) + (2\|\phi_1\|_\infty + 1)(\pi(X_j) - \pi(X_i)) \\ &\geq \mathbf{w}(e) + \phi_1(u) - \phi_1(v) + 2\|\phi_1\|_\infty + 1 \\ &\geq \mathbf{w}(e) + 1 \\ &\geq -1. \end{aligned}$$

Again, this confirms that  $\phi_2$  is also a price function on the edges that are between clusters (and not in  $E_{\text{del}}$ ) with respect to  $G^{+1}$ .  $\square$

**Claim 18.3.3.** The runtime of the call to  $\text{COMPUTEPRICEFUNKTION}(G_{\phi_2}^{+1})$  in Line 7 is  $\tilde{O}(n\sqrt{n})$  in expectation.

*Proof.* Again recall that  $\text{COMPUTEPRICEFUNKTION}(G_{\phi_2}^{+1})$  internally adds a dummy source  $d$  and edges of weight 0 from  $d$  to everyone in  $G_{\phi_2}^{+1}$ , and then computes the price function

in time  $\tilde{O}(\sum_v \eta(d, v))$  here  $\eta(\cdot, \cdot)$  is with respect to  $\phi_2$  (where we use that  $G$  is assumed constant-degree).

To finish the proof it remains to show that for every vertex  $v \in V$ ,  $\mathbb{E}[\eta(d, v)] = O(\sqrt{n} \log n)$ . To see this, fix any shortest path  $P$  from  $d$  to  $v$  in the graph  $G_{\phi_2}^{+1}$  (plus the vertex  $d$ ). Let  $P'$  be again the shortest  $u$  to  $v$  path where  $u$  is the vertex that appears on  $P$  after  $d$ . Again, we have that  $P'$  has weight at most 0 in  $G^{+1}$ , since  $P$  is of weight at most 0 since it is a shortest  $d$  to  $v$  path and a zero-weight edge  $(d, v)$  exists.

Further note that  $P'$  is also a shortest  $u$  to  $v$  path in  $G^{+1}$ . Finally, note that every edge that is negative in  $G_{\phi_2}$  is in  $E_{del}$ . Let  $N$  be the set of negative edges on  $P'$  in  $G^{+1}$ . We now want to analyze  $E[|N \cap E_{del}|]$  which yields an upper bound on  $\mathbb{E}[\eta(d, v)]$ .

But to this end, we need to use three different facts:

1. since every edge in  $N$  is a negative edge in  $G^{+1}$ , its weight in  $G^{+2}$  (where we computed the LDDs) is at most 1, and
2. each edge  $e$  is added to  $E_{del}$  with probability at most  $O(w^{+2}(e)/\sqrt{n})$  by the LDD decomposition theorem (see Theorem 17.2.2), and
3. there are at most  $n - 1$  edges on any shortest path, and thus also on  $P'$  and thus  $|N| \leq n - 1$ .

Combining these facts, we have that  $E[|N \cap E_{del}|] = \tilde{O}((n - 1)/\sqrt{n}) = \tilde{O}(\sqrt{n})$ . The claim follows.  $\square$

## 18.4 A Near-Linear-Time Algorithm

Using the same framework as above one can in fact obtain runtime  $\tilde{O}(m)$  for implementing `PRICEFUNCTIONONNORMALIZEDGRAPH()`. The only missing trick is that one can also scale off the maximum number of negative edges on any shortest path. To get an idea, observe that if we would know beforehand that this number is at most  $\sqrt{n}$ , then invoking the LDD algorithm with parameter  $n^{1/4}$  in-place of  $\sqrt{n}$  will yield an  $\tilde{O}(m^{5/4})$  algorithm. The key claim to look at is Theorem 18.3.3 where we then can use that any shortest path has at most  $\sqrt{n}$  negative edges and so its total weight in  $G^{+2}$  is at most  $\sqrt{n}$ , and therefore, we now add at most  $n^{1/4}$  edges to  $E_{del}$  in expectation.

# Chapter 19

## Maximum Matching

In this chapter, we will learn about classical algorithms for computing maximum matching in graphs. In this problem, the input is an undirected graph  $G = (V, E, \mathbf{w})$  on  $n$  vertices and  $m$  edges with integer edge weights  $\mathbf{w} : E \rightarrow \{1, \dots, W\}$ . A matching is a set of vertex-disjoint edges, and our goal is to find a matching  $M$  in  $G$  whose total weight  $\mathbf{w}(M) = \sum_{e \in M} \mathbf{w}(e)$  is maximized.

### 19.1 Bipartite Matching

Let us first study maximum weight matching in bipartite graphs. Assume  $G = (V = L \cup R, E, \mathbf{w})$  is bipartite with vertices in  $L$  and  $R$  having no edges in between them, respectively. The basic strategy is augmenting paths.

#### 19.1.1 Unit Edge Weights

To begin with, let us first consider unit edge weights (that is,  $\mathbf{w} = 1$ ), we could maintain a matching  $M$  and repeatedly finding augmenting paths to increase its size until we reach optimum. More specifically, given any matching  $M$ , we can apply breath-first search from free vertices and maintain a tree of vertices reachable from free vertices in  $L$  via simple alternating paths. In the end, if some free vertices in  $R$  are reachable via simple alternating paths, then we can recover an augmenting path. Let us formalize this scheme with more details.

**Definition 19.1.1** (augmenting path). Given a matching  $M$  of  $G$ , a vertex  $v$  is called *free* (with respect to  $M$ ) if  $v \notin V(M)$ , and *matched* if  $v \in V(M)$ .

Any simple path  $\langle u_1, u_2, \dots, u_l \rangle$  is called an *alternating path* if it begins with an unmatched vertex and whose edges belong alternately to the  $M$  and not to  $M$ , or more formally:

1.  $u_1$  is free.

2. All edges  $(u_{2i-1}, u_{2i})$  are unmatched ( $(u_{2i-1}, u_{2i}) \notin M$ ),  $\forall 1 \leq i \leq l/2$ .
3. All edges  $(u_{2i}, u_{2i+1})$  are matched ( $(u_{2i}, u_{2i+1}) \in M$ ),  $\forall 1 \leq i \leq (l-1)/2$ .

An alternating path  $\langle u_1, u_2, \dots, u_l \rangle$  is an *augmenting path* if  $u_l$  is free.

If we could find an augmenting path  $\langle u_1, u_2, \dots, u_l \rangle$  of  $M$ , then the edge set

$$M' = M \cup \{(u_1, u_2), (u_3, u_4), \dots, (u_{l-1}, u_l)\} \setminus \{(u_2, u_3), (u_4, u_5), \dots, (u_{l-2}, u_{l-1})\}$$

would be a matching of size  $|M| + 1$ . We argue that augmenting paths always exist whenever  $M$  is not maximum.

**Lemma 19.1.2.** *If  $M$  is not a maximum matching, then there exists an augmenting path with respect to  $M$ .*

*Proof.* Let  $M^*$  be the maximum matching of  $G$ . Then, in the graph  $(V, M \cup M^*)$ , all vertices have degree at most 2, and hence this graph is a union of simple cycles and paths. In each cycle  $C$ ,  $M \cap C$  and  $M^* \cap C$  have equal size, and in each path  $P$ , the sizes of  $M \cap P$  and  $M^* \cap P$  differ by at most 1. Since  $|M^*| > |M|$ , there must exist a path  $P$  in  $(V, M \cup M^*)$  where  $|M^* \cap P| = |M \cap P| + 1$  which would be an augmenting path.  $\square$

In bipartite graphs, any augmenting path starting at a free vertex in  $L$  should terminate at a free vertex in  $R$ . To efficiently find an augmenting path, we reduce this task to a reachability problem on a directed graph. More specifically, build a directed graph  $\vec{G}$  as following.

1. Start with  $G$ , and for edges  $(u, v) \notin M$ , assign a direction from  $L$  to  $R$ , and for edges  $(u, v) \in M$ , assign a direction from  $R$  to  $L$ .
2. Add an auxiliary source and terminal  $s, t$  to  $\vec{G}$ . For each free vertex  $u \in L$ , add a directed edge  $(s, u)$  to  $\vec{G}$ ; symmetrically, for each free vertex  $v \in R$ , add a directed edge  $(v, t)$  to  $\vec{G}$ .

**Lemma 19.1.3.** *Augmenting paths exist if  $s$  can reach  $t$  in  $\vec{G}$ .*

*Proof.* Let  $\langle s, u_1, u_2, \dots, u_l, t \rangle$  be a simple directed path from  $s$  to  $t$  in  $\vec{G}$ , then by definition  $\langle u_1, u_2, \dots, u_l \rangle$  is an augmenting path.  $\square$

The whole algorithm for unweighted bipartite matching can be summarized as the following pseudo-code Algorithm 20.

---

**Algorithm 20:** maximum matching in bipartite graph  $G = (L \cup R, E)$ 


---

```

1 initialize  $M \leftarrow \emptyset$ ;
2 repeat
3   build directed graph  $\vec{G}$ ;
4   if  $s$  can reach  $t$  in  $\vec{G}$  then
5     let  $\langle s, u_1, u_2, \dots, u_l, t \rangle$  be a simple directed path from  $s$  to  $t$  in  $\vec{G}$ ;
6     update
7      $M \leftarrow M \cup \{(u_1, u_2), (u_3, u_4), \dots, (u_{l-1}, u_l)\} \setminus \{(u_2, u_3), (u_4, u_5), \dots, (u_{l-2}, u_{l-1})\}$ ;
8   else
9   break;
9 return  $M$ ;

```

---

**Alternate description of augmenting path searching.** As will be useful for future extensions to non-bipartite graphs, we present here an alternate explanation of how augmenting paths are computed in bipartite graphs. For this task, we could grow trees of alternating paths from free vertices in  $L$  a breath-first tree manner until we reach some free vertices in  $R$ . More specifically, for any free vertex  $u \in L$ , we maintain a rooted breath-first search tree  $T_u$ , such that edges on the odd levels are unmatched, and edges on the even levels are matched.

To grow the search trees  $\{T_u \mid u \in L\}$  in a breath-first search manner, we iterate over the levels  $l = 1, 2, \dots$ . In the  $l$ -th iteration, go over all free vertices  $u \in L$  and all the leaf vertices  $v \in V(T_u)$ . For each edge  $(v, w) \notin M$  such that  $w$  is not contained in any search trees  $T_{u'}$  yet, add edge  $(v, w)$  to tree  $T_u$ . Furthermore, if  $w$  is matched with another vertex  $z$  in  $M$ , then add the matched edge  $(w, z)$  to  $T_u$  as well.

During the growth of all the search trees  $\{T_u\}$ , upon the first time that a free vertex  $v$  appears in some search tree  $T_u$ , we can retrieve an augmenting path and stop the breath-first search. Check Figure 19.1 for an illustration.

### 19.1.2 General Edge Weights

However, the above scheme does not work directly with edge weights. To handle general edge weights, we use a technique called the primal-dual method. The bipartite maximum weight matching problem can be formalized as an integer linear program:

$$\begin{aligned}
& \max \sum_{e \in E} \mathbf{w}(e) \cdot x(e) \\
\text{s.t. } & x(e) \in \{0, 1\} \quad \forall e \in E \\
& \sum_{e \ni u} x(e) \leq 1 \quad \forall u \in V
\end{aligned}$$

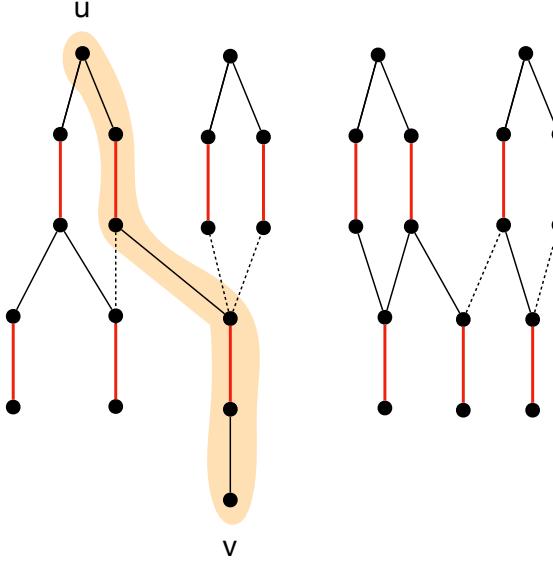


Figure 19.1: This picture shows an example of breath-first trees of alternating paths. The vertices on the top level are free vertices in  $L$ . All solid edges are tree edges, and red edges are matching edges in  $M$ , and the dashed edges are non-tree edges in  $E \setminus M$ . The orange path highlights an augmenting path found by the search tree.

Relax the constraint  $x(e) \in \{0, 1\}$  to  $x(e) \in [0, 1]$  and write the dual program:

$$\begin{aligned} & \min \sum_{u \in V} y(u) \\ \text{s.t. } & y(u) + y(v) \geq \mathbf{w}(e) \quad \forall e = (u, v) \in E \\ & y(u) \geq 0 \quad \forall u \in V \end{aligned}$$

The basic strategy is to repeatedly augment a feasible primal integral solution  $\{x(e) \mid e \in E\}$ , which is actually a matching  $M$ , as well as a feasible dual solution  $\{y(u) \mid u \in V\}$  and in the end reach the condition that:

$$\sum_{e \in E} \mathbf{w}(e) \cdot x(e) = \sum_{u \in V} y(u)$$

By the weak duality of linear programs, we conclude that the our matching  $M$  should be optimal at the end.

To describe the algorithm, we rely on the notation of tightness.

**Definition 19.1.4.** An edge  $e = (u, v)$  is **tight** if  $y(u) + y(v) = \mathbf{w}(e)$ .

At the beginning of the algorithm, set  $M = \emptyset$ ,  $y(u) = W, \forall u \in L$  and  $y(v) = 0, \forall v \in R$ , so  $y$  would be a feasible solution for the dual program. During the augmentation steps, define  $G_y = (L \cup R, E_y)$  to be the subgraph of  $G$  consisting of all tight edges. Throughout the algorithm, we will maintain the following invariant.

**Invariant 19.1.5.** During the algorithm, all edges in  $M$  are tight. Furthermore, for any  $u \in L \setminus V(M)$ , we have  $y(u) = \min_{v \in L} \{y(v)\}$ .

In each round, the algorithm augment  $M$  in  $G_y$  as long as it could, and then adjust the vertex labels  $y$  so that more edges become tight. The algorithm terminates when  $y(v) = 0$  for all free vertex  $v$ .

---

**Algorithm 21:** maximum weight matching in graph  $G = (L \cup R, E, \mathbf{w})$

---

```

1 initialize  $M \leftarrow \emptyset$ ,  $y(u) \leftarrow W$ ,  $y(v) \leftarrow 0$ ,  $\forall u \in L, v \in R$ ;
2 while  $y(u) > 0$  for some  $u \notin V(M)$  do
3   define  $G_y = (L \cup R, E_y)$  to be the subgraph of all tight edges;
4   repeatedly augment  $M$  in graph  $G_y$ , until there is no more augmenting paths in the
      same way as Algorithm 20;
5   /* adjust the dual variables  $y$  to make more tight edges */  

6   assign directions to edges in  $G_y$ ;
7   compute  $Z \subseteq V$  which is the set of vertices reachable from all free vertices in  $L$  in
      directed graph  $\vec{G}_y$ ;
8   define  $F = \{(u, v) \in E \setminus E_y \mid u \in L \cap Z, v \in R \setminus Z\}$  and
       $\Delta = \min \left\{ \min_{(u,v) \in F} \{y(u) + y(v) - \mathbf{w}(u, v), \min_{u \in L} \{y(u)\}\} \right\}$ ;
9   update  $y(u) \leftarrow y(u) - \Delta$ ,  $y(v) \leftarrow y(v) + \Delta$ ,  $\forall u \in L \cap Z, v \in R \cap Z$ ;
9 return  $M$ ;

```

---

The advantage of finding augmenting paths in graph  $G_y$  is that it preserves Invariant 19.1.5 because by definition all edges in  $G_y$  are tight. The main question here is how to adjust the dual variables to generate more tight edges.

To do this, we will implement dual adjustment through a reachability problem in graph  $G_y$ . For any edge  $(u, v) \in (L \times R) \cap E_y$ , if  $(u, v) \in M$ , assign an edge direction from  $v$  to  $u$ ; if  $(u, v) \notin M$ , assign an edge direction from  $u$  to  $v$ . Call this directed graph  $\vec{G}_y$ . Compute a vertex subset  $Z \subseteq V$  which is the set of vertices reachable from all free vertices in  $L$ . Define an edge set  $F \subseteq E \setminus E_y$  as well as a value  $\Delta$ :

$$F = \{(u, v) \in E \setminus E_y \mid u \in L \cap Z, v \in R \setminus Z\}$$

$$\Delta = \min \left\{ \min_{(u,v) \in F} \{y(u) + y(v) - \mathbf{w}(u, v), \min_{u \in L} \{y(u)\}\} \right\}$$

Then, for each  $u \in L \cap Z$  decrease  $y(u) \leftarrow y(u) - \Delta$ , and for each  $v \in R \cap Z$  increase  $y(v) \leftarrow y(v) + \Delta$ . After this dual adjustment, more tight edges will appear. Check Figure 19.2 for an illustration.

The pseudo-code Algorithm 21 summarizes this algorithm.

**Observation 19.1.6.** Invariant 19.1.5 is preserved by dual adjustments, and  $E_y$  is only monotonically growing.

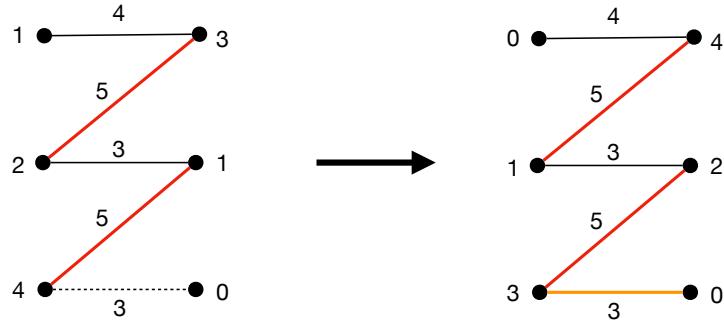


Figure 19.2: In the example, solid edges are tight, and red edges are matched. Before the dual adjustment, the dashed edge on the left-hand side is not tight. After the dual adjustment, the orange edge is a new tight edge.

**Lemma 19.1.7.** *The runtime of the algorithm is  $O(m^2)$ .*

*Proof.* Each iteration can be implemented in  $O(m)$  which mainly computes reachability in the directed version of  $G_y$ . In each iteration of the while-loop, if  $\Delta > \min_{u \in L} \{y(u)\}$ , then  $\Delta$  would be equal to some  $y(u_*) + y(v_*) - \mathbf{w}(u_*, v_*)$ ,  $(u_*, v_*) \in F$ , after adjusting the dual variables,  $G_y$  would now include a new edge  $(u_*, v_*)$ . Otherwise, if  $\Delta = \min_{u \in L} \{y(u)\}$ , all free vertices may now have zero labels, so the algorithm terminates. Overall, there can be at most  $m$  iterations, so the runtime is  $O(m^2)$ . Note that a tight edge  $(u, v)$  never becomes non-tight again because the decrease of  $y(u)$  cancels out with the increase of  $y(v)$  due to the way we adjust the vertex labels.  $\square$

**Lemma 19.1.8.** *When the algorithm terminates,  $M$  is a maximum weight matching.*

*Proof.* Let  $M^*$  be the a maximum weight matching. Let us show that  $\mathbf{w}(M) \geq \mathbf{w}(M^*)$ . In fact, since all edges in  $M$  are tight and free vertices have zero labels, we have:

$$\begin{aligned} \mathbf{w}(M) &= \sum_{u \in V} y(u) \\ &= \sum_{(u,v) \in M^*} (y(u) + y(v)) + \sum_{u \notin V(M^*)} y(u) \\ &\geq \sum_{(u,v) \in M^*} \mathbf{w}(u, v) = \mathbf{w}(M^*) \quad \square \end{aligned}$$

### 19.1.3 Faster Implementation

Next, we show how to optimize the runtime from  $O(m^2)$  to  $\tilde{O}(mn)$ . Notice that  $M$  could be augmented by at most  $n$  times, while the graph  $G_y$  may grow for at most  $m$  times. So the bottleneck of the algorithm is growing  $G_y$ , and most of the  $m$  while-loop iterations are

about growing  $G_y$ , not augmenting  $M$ . So, to improve the runtime, the idea is to compress multiple consecutive while-loop iterations which only enlarge  $G_y$  into a single iteration, such that the next iteration can augment  $M$ .

To implement this idea, the main observation is that the concatenation of multiple while-loop iterations is actually a shortest path computation from left free vertices to right free vertices. More specifically, for each edge  $(u, v) \in E$ , assign edge length  $\mu(u, v) = y(u) + y(v) - \mathbf{w}(u, v) \geq 0$ , and orient all matching edges from  $R$  to  $L$ , and non-matching edges from  $L$  to  $R$ . Then, add a terminal  $s$  which is connected to all free vertices in  $L$  (with zero weight edges). Call this directed weighted graph  $\vec{H}_y = (V, E, \mu)$ . After that, compute the single-source shortest paths from  $s$  in  $\vec{H}_y$ .

Let  $t \in R \cap V(M)$  be the free vertex which minimizes  $d(s, t, \vec{H}_y)$ , and define:

$$\Delta = \min_{u \in L} \left\{ d(s, t, \vec{H}_y), y(u) \right\}$$

Then, for each  $u \in L$  such that  $d(s, u, \vec{H}_y) \leq \Delta$ , update the vertex label:

$$y(u) \leftarrow y(u) - \Delta + d(s, u, \vec{H}_y)$$

and for each  $v \in R$  such that  $d(s, v, \vec{H}_y) \leq \Delta$ , update the vertex label:

$$y(v) \leftarrow y(v) + \Delta - d(s, v, \vec{H}_y)$$

This algorithm is summarized in Algorithm 22.

---

**Algorithm 22:** faster maximum weight matching in graph  $G = (L \cup R, E, \mathbf{w})$ 


---

```

1 initialize  $M \leftarrow \emptyset$ ,  $y(u) \leftarrow W$ ,  $y(v) \leftarrow 0$ ,  $\forall u \in L, v \in R$ ;
2 while  $y(u) > 0$  for some  $u \notin V(M)$  do
3   repeatedly augment  $M$  in graph  $G_y$ , until there is no more augmenting paths;
   /* adjust the dual variables  $y$  to make more tight edges */ 
4   build directed weighted graph  $\vec{H}_y$ ;
5   let  $t \in R \cap V(M)$  be the free vertex which minimizes  $d(s, t, \vec{H}_y)$ ;
6   compute  $\Delta \leftarrow \min_{u \in L} \left\{ d(s, t, \vec{H}_y), y(u) \right\}$ ;
7   for  $u \in L$  such that  $d(s, u, \vec{H}_y) \leq \Delta$  do
8     update the vertex label  $y(u) \leftarrow y(u) - \Delta + d(s, u, \vec{H}_y)$ ;
9   for  $v \in R$  such that  $d(s, v, \vec{H}_y) \leq \Delta$  do
10    update the vertex label  $y(v) \leftarrow y(v) + \Delta - d(s, v, \vec{H}_y)$ ;
11 return  $M$ ;

```

---

First, let us verify that our dual adjustment scheme preserves feasibility of the dual solution as well as Invariant 19.1.5.

**Lemma 19.1.9.** *In each iteration of the while-loop, dual feasibility and Invariant 19.1.5 are both preserved by our dual adjustment scheme.*

*Proof.* Let use  $y$  and  $y'$  to denote the dual variables before and after the adjustment, respectively. First, we verify that all  $y'$  values are non-negative. In fact, for any  $v \in R$  such that  $d(s, v, \overrightarrow{H}_y) \leq \Delta$ , we have  $y'(v) = y(v) + \Delta - d(s, v, \overrightarrow{H}_y) \geq 0$ ; for any  $u \in L$ , by definition of  $\Delta$ , we know that  $\Delta \leq y(u)$ , and therefore  $y'(u) = y(u) - \Delta + d(s, u, \overrightarrow{H}_y) \geq 0$ .

Now, consider any edge  $(u, v) \in E \cap (L \times R)$ . Previously before the change, we had  $y(u) + y(v) - \mathbf{w}(u, v) = \mu(u, v) \geq 0$ . Since  $y'(v)$  does not decrease, we only need to worry about the case when  $y'(u) < y(u)$ . There are two cases.

- $(u, v) \notin M$ . In this case, noticing that we always have  $y'(v) \geq y(v) + \Delta - d(s, v, \overrightarrow{H}_y)$ , therefore

$$\begin{aligned} y'(u) + y'(v) &\geq y(u) + y(v) + d(s, u, \overrightarrow{H}_y) - d(s, v, \overrightarrow{H}_y) \\ &= \mathbf{w}(u, v) + \mu(u, v) + d(s, u, \overrightarrow{H}_y) - d(s, v, \overrightarrow{H}_y) \\ &\geq \mathbf{w}(u, v) \end{aligned}$$

- $(u, v) \in M$ . In this case,  $(u, v)$  is the only in-edge for  $u$ , and so  $d(s, u, \overrightarrow{H}_y) = d(s, v, \overrightarrow{H}_y)$ . Therefore,

$$\begin{aligned} y'(u) + y'(v) &= y(u) + y(v) + d(s, u, \overrightarrow{H}_y) - d(s, v, \overrightarrow{H}_y) \\ &= \mathbf{w}(u, v) + d(s, u, \overrightarrow{H}_y) - d(s, v, \overrightarrow{H}_y) \\ &= \mathbf{w}(u, v) \end{aligned}$$

This equality also proves that all matching edges in  $M$  remain tight after the dual adjustment.

Finally, for Invariant 19.1.5, notice that  $d(s, u, \overrightarrow{H}_y) = 0$  for all free vertices  $u \in L$ , the vertex labels of free vertices on  $L$  have decreased by the same amount  $\Delta$ . So Invariant 19.1.5 also holds.  $\square$

Next, let us argue that after the dual adjustment, a new augmenting path shows up in  $G_y$ .

**Lemma 19.1.10.** *In each iteration of the while-loop, after the dual adjustment, either a new augmenting path in  $G_y$  shows up, or the algorithm terminates.*

*Proof.* If  $\Delta = \min_{u \in L} \{y(u)\}$ , then after the dual adjustment, all the free vertices have zero vertex labels, so the algorithm would terminate. Otherwise,  $\Delta = d(s, t, \overrightarrow{H}_y)$ .

Let  $P$  be the shortest path from  $s$  to  $t$  in graph  $\overrightarrow{H}_y$ . We claim that all edges on  $P$  become tight after the dual adjustment, which would imply that  $G_y$  may now contain a new augmenting

path. In fact, for any edge  $(u, v) \in P$ , if  $u \in L, v \in R$ , then

$$\begin{aligned} y'(u) + y'(v) &= y(u) + y(v) + d(s, u, \overrightarrow{H_y}) - d(s, v, \overrightarrow{H_y}) \\ &= \mathbf{w}(u, v) + \mu(u, v) + d(s, u, \overrightarrow{H_y}) - d(s, v, \overrightarrow{H_y}) \\ &= \mathbf{w}(u, v) \end{aligned}$$

The last inequality holds as  $P$  is the shortest path and thus  $\mu(u, v) + d(s, u, \overrightarrow{H_y}) = d(s, v, \overrightarrow{H_y})$ . Otherwise, if  $u \in R, v \in L$ , then  $(u, v) \in M$ , and therefore

$$\begin{aligned} y'(u) + y'(v) &= y(u) + y(v) + d(s, u, \overrightarrow{H_y}) - d(s, v, \overrightarrow{H_y}) \\ &= \mathbf{w}(u, v) + d(s, u, \overrightarrow{H_y}) - d(s, v, \overrightarrow{H_y}) \\ &= \mathbf{w}(u, v) \end{aligned}$$

This concludes the proof.  $\square$

By the above statement, we know that the while-loop contains at most  $n$  iterations, and so the total runtime becomes  $\tilde{O}(mn)$  if we use Dijkstra's algorithm to compute single-source shortest path in each iteration.

## 19.2 General Graph Matching

### 19.2.1 Unit Edge Weights

Let us now turn to maximum matching in general graphs. The main technical difficulty with general graphs is that we cannot formulate augmenting paths as a reachability problem (even in unit-weight graphs) as we did in bipartite graphs due to odd cycles formed by alternating paths. This technical issue was addressed by Edmonds [Edm65]. The key notion of Edmonds' algorithm is called *blossom*. Roughly speaking, given a matching  $M$ , a blossom is an odd-length alternating cycle whose constituents are vertices or smaller blossoms. More formally, a blossom is specified by a tuple  $(B, E_B, \beta(B))$ , where  $B \subseteq V$  is a vertex subset,  $E_B \subseteq E$  is an edge subset,  $\beta(B) \in B$  is a special vertex called the *base*. Blossoms follow an inductive definition below.

**Definition 19.2.1** ([Edm65]). A single vertex  $v$  forms a trivial blossom as a tuple  $(\{v\}, \emptyset, v)$ . Inductively, let  $B_1, B_2, \dots, B_{2l+1}$  be an odd sequence of blossoms. Suppose there exists a sequence of edges  $e_1, e_2, \dots, e_{2l+1}$  such that  $e_i \in B_i \times B_{i+1}$  ( $B_{2l+2} = B_1$ ). The tuple  $(\bigcup_{i=1}^{2l+1} B_i, \bigcup_{i=1}^{2l+1} E_{B_i}, \beta(B_1))$  is identified as a blossom if  $e_i \in M$  iff  $i$  is even. Check Figure 19.3 for an illustration.

We can show that for any blossom  $(B, E_B, \beta(B))$ ,  $E_B \cap M$  matches all vertices in  $B$  but  $\beta(B)$ . The key property of blossoms is the following statement.

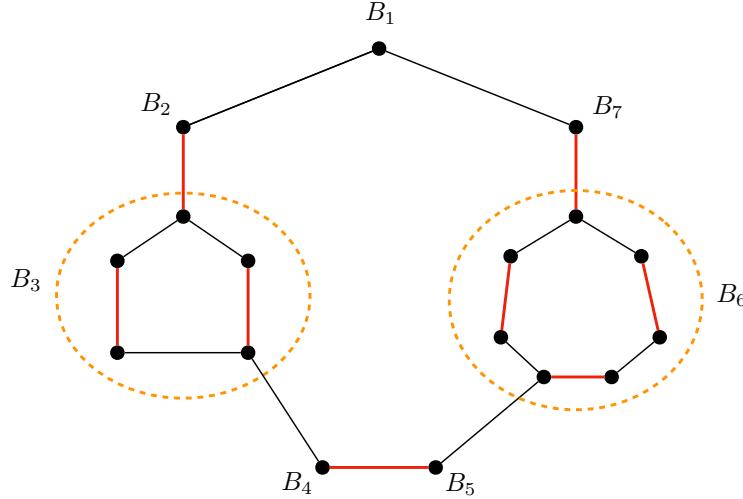


Figure 19.3: This picture shows an example of a blossom, where the red edges are matching edges, and the black ones are non-matching edges. The two subgraphs in the dashed orange cycles are two smaller blossoms.

**Observation 19.2.2.** *If the base of a blossom is free, then all vertices in the blossom are also reachable from this free base vertex via simple alternating paths.*

This property implies that we can shrink a blossom into a single node if this blossom is reachable from free vertices. Informally speaking, to search for augmenting paths from free vertices, we can apply the same breath-first search through alternating paths. If we encounter some odd cycles during the process, then this odd cycle must form a blossom. In this case, we can shrink the entire blossom into a single node, and then continue on with the breath-first search procedure. As the graph cannot shrink indefinitely, we will find an augmenting path at the end in polynomial time.

To formalize this idea, we will maintain a laminar family of blossoms  $\Omega \subseteq 2^V$ . For each free vertex  $u \in V$ , we will maintain a tree  $T_u$  rooted at  $u$  in the contracted graph  $\mathcal{G} = G/\Omega$  where each the blossom is contracted into a single node.

**Definition 19.2.3** (out/in labels). For each node in tree  $T_u$ , if its distance to the root is even, label this node as *out*, and otherwise label it as *in*. All vertices contracted within blossoms are labeled out or in accordingly.

The trees  $\{T_u \mid u \notin V(M)\}$  will satisfy the following properties.

- All the rooted trees  $\{T_u \mid u \notin V(M)\}$  will be vertex-disjoint.
- Each root-to-leaf path in any such tree  $T_u$  will be an alternating path in  $\mathcal{G}$ , and each leaf node in  $T_u$  is labeled out.

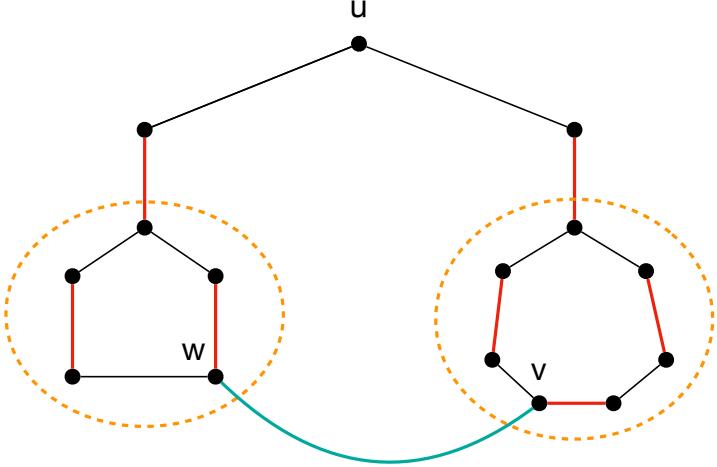


Figure 19.4: In this picture, all the red and black edges are tree edges of  $T_u$ , and blue edge  $(v, w)$  connects two out vertices. In this example, we can contract the entire tree  $T_u$  into a single node.

To grow all the trees  $\{T_u \mid u \notin V(M)\}$ , each time we select an arbitrary free vertex  $u$  as well as an out vertex  $v \in V$  contracted in a leaf node in  $B_v \in V(T_u)$ . Enumerate all edges  $(v, w) \in E$  and consider the following possibilities.

- **Blossom contraction.** Vertex  $w$  is residing in a node  $B_w \in V(T_u)$ . If  $w$  is an out vertex, then all the nodes on the tree path between  $B_v, B_w$  would make a blossom in  $G$ ; call this blossom  $B$ , add  $B$  to  $\Omega$ , and contract  $B$  in  $\mathcal{G}$ . Check Figure 19.4 for an illustration.
- **Tree extension.** Vertex  $w$  is currently not in any of the trees  $T_{u'}, u' \notin V(M)$ . In this case,  $w$  cannot be free, and let  $(w, z) \in M$  be its matched edge. We argue that  $z$  also cannot be in any existing trees  $T_{u'}, u' \notin V(M)$ , since otherwise  $w$  would be added to some trees already. Then, let us add edges both  $(v, w), (w, z)$  to tree  $T_u$ .
- **Augmentation.** Vertex  $w$  is in another tree  $T_{u'}, u' \notin V(M) \setminus \{u\}$ . Suppose  $w$  is an out vertex, then we have found an augmenting path between  $u, u'$ . In this case, flip the augmenting path to increase  $M$  and start over. Check Figure 19.5 for an illustration.

The above algorithm is summarized as pseudo-code Algorithm 23.

**Lemma 19.2.4.** *When  $M$  is not a maximum matching, then the algorithm always finds an augmenting path.*

*Proof sketch.* Take the maximum matching  $M^*$  and show that in any augmenting path in  $M \Delta M^*$  there should be an edge connecting two out vertices in two different trees  $T_u, T_{u'}$ .  $\square$

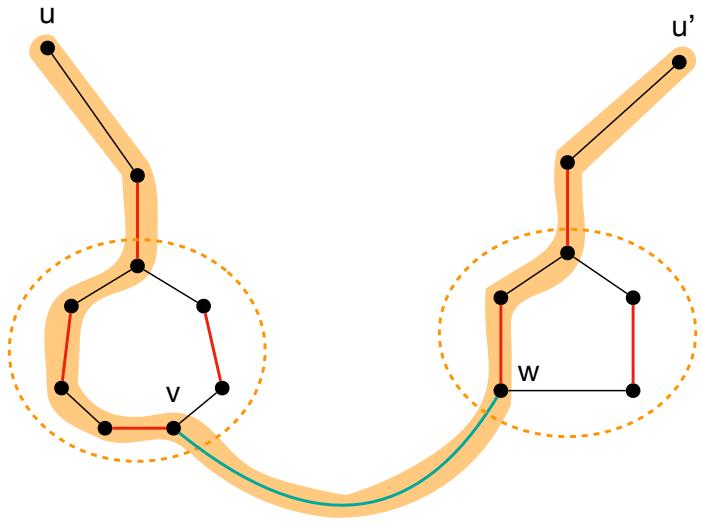


Figure 19.5: In this picture, all the red and black edges are tree edges of  $T_u, T_{u'}$ , and the blue edge  $(v, w)$  connects two out vertices in  $T_u, T_{u'}$ . Then, we can find an augmenting path between  $u, u'$  which is highlighted as the orange path.

---

**Algorithm 23:** maximum matching in general graph  $G = (V, E)$

---

```

1 initialize  $M \leftarrow \emptyset$ ;
2 repeat
3   initialize  $\Omega \leftarrow \emptyset, T_u \leftarrow \{u\}, \forall u \notin V(M)$ ;
4   repeat
5     for  $u \notin V(M)$ , out vertex  $v$  contracted in a leaf node  $B_v \in V(T_u)$ , and
        $(v, w) \in E$  do
6       if  $w$  is an out vertex then
7         contract all the nodes on the tree path between  $B_v, B_w$  as a blossom;
8       else if  $w$  is not in any tree then
9         let  $(w, z) \in M$  be a matching edge, and add  $(v, w), (w, z)$  to  $T_u$ ;
10      else if  $w$  is in another tree node  $T_{u'}$  and  $w$  is an out vertex then
11        flip the augmenting path between  $u, u'$  to increase  $|M|$ ;
12        break;
13      if  $M$  was not augmented in this iteration then
14        break;
15 return  $M$ ;

```

---

### 19.2.2 General Edge Weights

To deal with edge weights, we will combine the previous primal-dual approach with the above blossom approach. The primal linear program for maximum weight matching will have extra constraints for every odd vertex set.

$$\begin{aligned}
& \max \sum_{e \in E} \mathbf{w}(e) \cdot x(e) \\
\text{s.t. } & x(e) \in [0, 1] && \forall e \in E \\
& \sum_{e \ni u} x(e) \leq 1 && \forall u \in V \\
& \sum_{u, v \in B, (u, v) \in E} \mathbf{w}(u, v) \leq \frac{|B| - 1}{2} && \forall B \subseteq V, |B| \text{ is odd}
\end{aligned}$$

As for the dual program, besides the vertex labels  $y : V \rightarrow \mathbb{R}$ , we will introduce a new set of variables  $z : 2^V \rightarrow \mathbb{R}$  defined on all odd vertex subsets. The dual linear program is formalized below:

$$\begin{aligned}
& \min \sum_{u \in V} y(u) + \sum_{B \subseteq V, |B| \geq 3 \text{ is odd}} z(B) \cdot \frac{|B| - 1}{2} \\
\text{s.t. } & yz(u, v) = y(u) + y(v) + \sum_{B \supset \{u, v\}} z(B) \geq \mathbf{w}(e) && \forall e = (u, v) \in E \\
& y(u) \geq 0 && \forall u \in V \\
& z(B) \geq 0 && \forall B \subseteq V, |B| \text{ is odd}
\end{aligned}$$

The algorithm maintains a matching  $M$  as well as vertex labels  $y$ , and a laminar family of blossoms  $\Omega \subseteq 2^V$  as well as their dual values  $z(B), B \in \Omega$ ; the dual values of all other odd sets will be zero. Let  $\mathcal{G}$  be the contracted graph where every maximal blossom in  $\Omega$  is contracted into a single node, and then we only keep the tight edges of  $E$  under dual variables  $y, z$  in  $\mathcal{G}$ .

The following properties will be ensured during the algorithm.

**Invariant 19.2.5.** *For any free vertex  $u$ ,  $y(u) = \min_{v \in V} \{y(v)\}$ . For each blossom edge  $e \in E_B$  for some  $B \in \Omega$ ,  $e$  is tight, namely  $yz(e) = \mathbf{w}(e)$ . Also, for any maximal blossom  $B \in \Omega$  (a tree root in the laminar structure of  $\Omega$ ), we have  $z(B) > 0$ .*

At the beginning of the algorithm, initialize  $M = \emptyset, y = W, \Omega = \emptyset$ . More specifically, given any  $M$  and duals  $y, z$ , begin with empty search trees at each free vertex, and initialize  $\Omega \leftarrow \emptyset$ . For any free vertex  $u \notin V(M)$ , we will maintain a tree  $T_u$  in  $\mathcal{G}$  rooted at the contracted node containing  $u$  such that each root-to-leaf path is an alternating path. As with the unit weight setting, we also label the nodes and vertices as out or in.

If there is currently no augmenting path in  $\mathcal{G}$ , to generate an augmenting path consisting of tight edges, we will repeatedly adjust the dual variables  $y, z$  and extend the search trees. More specifically, given a value  $\Delta \geq 0$ , we perform the following dual adjustments:

- $y(v) \leftarrow y(v) - \Delta$  for all out-vertices, and  $y(v) \leftarrow y(v) + \Delta$  for all in-vertices;
- $z(B) \leftarrow z(B) + 2\Delta$  for all maximal out-blossoms, and  $z(B) \leftarrow z(B) - 2\Delta$  for all maximal in-blossoms.

**Observation 19.2.6.** *The dual adjustments always preserve Invariant 19.2.5 as well as dual feasibility.*

We will find the smallest value of  $\Delta \geq 0$  such that after the above dual adjustments, one of the following events occur.

- **Termination.** All the free vertices have zero vertex labels. In this case, the algorithm terminates and returns  $M$  as a maximum weight matching.
- **Blossom dissolution.** The dual value  $z(B)$  has reached 0 for some maximal in-blossoms; this operation does not appear in the unit-weight setting. In this case, remove  $B$  from  $\Omega$ . Suppose blossom  $B$  was contained in search tree  $T_u$  in  $\mathcal{G}$ , the sub-blossoms of  $B$  are connected as an odd alternating cycle  $B_1, B_2, \dots, B_{2k+1}$ . As  $B$  was an in-blossom, its parent in  $B$  must be connected with  $B$  via a non-matching edge with one endpoint residing in a sub-blossom  $B_i$ . Without loss of generality, assume  $i$  is an odd integer. Then, add the entire alternating path  $B_1, B_2, \dots, B_i$  to search tree  $T_u$ , and discard  $B_{i+1}, B_{i+2}, \dots, B_{2k+1}$  from  $T_u$ . After that, label  $B_1, B_2, \dots, B_i$  as in or out-blossoms according to their distance parities to root of  $T_u$ . Check Figure 19.6 for an illustration.
- **New tight edges.** Some new tight edges show up. We argue that any new tight edge must connect two out-blossoms in the contracted graph  $\mathcal{G}$ , or extend some search trees from an out-blossom to an unlabeled vertex, because in-vertices are only raising the vertex labels. Here are three possibilities which are similar to the unit-weight setting.
  - If a new tight edge connects to some unlabeled blossoms, then extend the search tree via alternating path as much as possible, and label all the blossoms according to their distance parities.
  - If a new tight edge  $e$  connects two out-blossoms from two different search trees, then we can find a new augmenting path.
  - If these two out-blossoms  $B_1, B_2$  come from the same search tree, then the unique path between  $B_1, B_2$  in the search tree together with  $e$  form an odd cycle of alternating path, which creates a new blossom; in this case, we contract this new blossom and add it to  $\Omega$ .

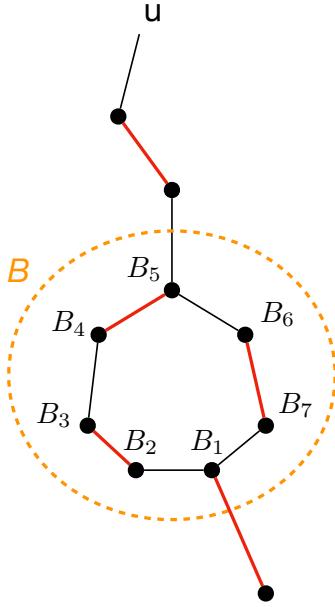


Figure 19.6: In this example, after dissolving blossom  $B$  which is drawn as the orange cycle, we keep  $B_1, B_2, B_3, B_4, B_5$  in the tree  $T_u$ , and remove  $B_6, B_7$  from  $T_u$ . Since  $B_6, B_7$  are no longer in any search trees rooted at free vertices, we would remove their in/out labels.

- **Augmentation.** After a sequence of blossom dissolution and new tight edges, there will be a new augmentation at some point. This happens only when we find a tight edge connecting two different search trees  $T_1, T_2$ . In this case, let us unlabel all vertices contracted in  $T_1, T_2$ , and flip the corresponding augmenting path.

To find the proper value of  $\Delta$  such that one of the above events occur, we can set  $\Delta$  to be the minimum among the following terms:

$$\min_{s \notin V(M)} \{y(t)\}$$

$$\frac{1}{2} \cdot \min_{B \in \Omega, \text{in blossom}} \{z(B)\}$$

$$\frac{1}{2} \cdot \min_{\substack{s, t \text{ are out and in different trees}, (s, t) \text{ is not tight}}} \{\mathbf{w}(s, t) - y(s) - y(t)\}$$

$$\min_{\substack{s \text{ is out}, t \text{ is unlabeled}}} \{\mathbf{w}(s, t) - y(s) - y(t)\}$$

The whole algorithm is summarized Algorithm 24.

---

**Algorithm 24:** maximum weight matching in graph  $G = (V, E, \mathbf{w})$ 


---

```

1 initialize  $M \leftarrow \emptyset$ ,  $y(*) \leftarrow W$ ,  $z(*) \leftarrow 0$ ,  $\Omega \leftarrow \emptyset$ ;
2 while  $\min_{u \notin V}\{y(u)\} > 0$  do
3   repeat
4     set  $\Delta$  to be the minimum among the four terms we discussed;
5      $y(v) \leftarrow y(v) - \Delta$  for all out-vertices, and  $y(v) \leftarrow y(v) + \Delta$  for all in-vertices;
6      $z(B) \leftarrow z(B) + 2\Delta$  for all maximal out-blossoms, and  $z(B) \leftarrow z(B) - 2\Delta$  for all
      maximal in-blossoms;
7     if free vertices have zero vertex labels then
8       break;
9     else if some in-blossoms  $B$  has zero dual then
10      dissolve  $B$  and add sub-blossoms to the search tree;
11    else
12      /* new tight edges appear */  
*/  

13      extend search trees, or find augmentations, or contract out-blossoms;
14    if there exists an augmentation then
15      suppose this is due to a new tight edge connecting search trees  $T_1, T_2$ , then
        unlabel all blossoms in  $T_1, T_2$ , and flip the alternating path;
16
17 return  $M$ ;

```

---

**Lemma 19.2.7.** *Between every two consecutive matching augmentation, there are at most  $O(n)$  events of blossom dissolution and new tight edges.*

*Proof.* In general, since there are at most  $n$  in-blossoms, and blossom dissolution only destroys in-blossoms, so there can be at most  $n$  such events; actually, in our algorithm described above, we never create in-blossoms, so there is none of them.

When new tight edges emerge, either we extend a search tree by two more vertices, or contract an out-blossom. Such events could also happen at most  $n$  times. This concludes the proof.  $\square$

**Corollary 19.2.8.** *The total runtime of the algorithm is  $O(mn^2)$ .*

*Proof.* Between any two consecutive matching augmentations, there are  $O(n)$  events of new tight edges or blossom dissolution. As each event takes  $O(m)$  time by scanning all edges, the total time becomes  $O(mn^2)$ .  $\square$

**Lemma 19.2.9.** *When the algorithm terminates, it returns a maximum weight matching.*

*Proof.* Let  $M^*$  be an arbitrary matching. Then, by dual feasibility and that all free vertices

have zero duals, we have:

$$\begin{aligned}
\mathbf{w}(M^*) &= \sum_{(u,v) \in M^*} \mathbf{w}(u,v) \\
&\leq \sum_{(u,v) \in M^*} yz(u,v) \\
&\leq \sum_{u \in V(M^*)} y(u) + \sum_{B \subseteq V, |B| \geq 3 \text{ is odd}} z(B) \cdot \frac{|B|-1}{2} \\
&\leq \sum_{u \in V} y(u) + \sum_{B \subseteq V, |B| \geq 3 \text{ is odd}} z(B) \cdot \frac{|B|-1}{2} \\
&= \sum_{u \in V(M)} y(u) + \sum_{B \subseteq V, |B| \geq 3 \text{ is odd}} z(B) \cdot \frac{|B|-1}{2} = \mathbf{w}(M)
\end{aligned}$$

The last equality uses the fact that there are  $\frac{|B|-1}{2}$  matched edges in every blossoms  $B \in \Omega$ . This immediately concludes the proof.  $\square$

**Remark 19.2.10.** It is possible to compress every  $O(n)$  iterations into a single iteration in  $O(m + n \log n)$ , like what we did for bipartite graphs. This improvement is rather involved and is detailed in [Gab18].

# Bibliography

- [BV04] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge university press, 2004.
- [DS08] Samuel I. Daitch and Daniel A. Spielman. Faster approximate lossy generalized flow via interior point algorithms. In *Proceedings of the Fortieth Annual ACM Symposium on Theory of Computing*, STOC '08, pages 451–460, New York, NY, USA, May 2008. Association for Computing Machinery.
- [Edm65] Jack Edmonds. Paths, trees, and flowers. *Canadian Journal of mathematics*, 17:449–467, 1965.
- [Gab18] Harold N Gabow. Data structures for weighted matching and extensions to b-matching and f-factors. *ACM Transactions on Algorithms (TALG)*, 14(3):1–80, 2018.
- [KS16] R. Kyng and S. 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, 2016.
- [LS20a] Yang P. Liu and Aaron Sidford. Faster Divergence Maximization for Faster Maximum Flow. *arXiv:2003.08929 [cs, math]*, April 2020.
- [LS20b] Yang P. Liu and Aaron Sidford. Faster energy maximization for faster maximum flow. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pages 803–814, 2020.
- [Mad13] Aleksander Madry. Navigating central path with electrical flows: From flows to matchings, and back. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 253–262. IEEE, 2013.
- [Spi19] Daniel A Spielman. Spectral and Algebraic Graph Theory, 2019.
- [SS11] Daniel A Spielman and Nikhil Srivastava. Graph sparsification by effective resistances. *SIAM Journal on Computing*, 40(6):1913–1926, 2011.
- [ST04] Daniel A. Spielman and Shang-Hua Teng. Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems. In *Proceedings of the*

*Thirty-Sixth Annual ACM Symposium on Theory of Computing*, STOC '04, page 81–90, New York, NY, USA, 2004. Association for Computing Machinery.

- [T<sup>+</sup>15] Joel A Tropp et al. An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230, 2015.
- [Tro19] Joel A Tropp. Matrix concentration & computational linear algebra. 2019.