

ETH Zurich

---

# **Advanced Graph Algorithms and Optimization**

---

Rasmus Kyng & Maximilian Probst

Spring 2021

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 submit a pull request at [https://github.com/rjkyng/agao21\\_script](https://github.com/rjkyng/agao21_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.

# Contents

<b>1</b>	<b>Course Introduction</b>	<b>7</b>
1.1	Overview . . . . .	7
1.2	Electrical Flows and Voltages - a Graph Problem from Middle School? . . .	7
1.3	Convex Optimization . . . . .	14
1.4	More Graph Optimization Problems . . . . .	16
<b>I</b>	<b>Introduction to Convex Optimization</b>	<b>19</b>
<b>2</b>	<b>Some Basic Optimization, Convex Geometry, and Linear Algebra</b>	<b>20</b>
2.1	Overview . . . . .	20
2.2	Optimization Problems . . . . .	20
2.3	A Characterization of Convex Functions . . . . .	22
2.3.1	First-order Taylor Approximation . . . . .	23
2.3.2	Directional Derivatives . . . . .	24
2.3.3	Lower Bounding Convex Functions with Affine Functions . . . . .	24
2.4	Conditions for Optimality . . . . .	26
<b>3</b>	<b>Convexity and Second Derivatives, Gradient Descent and Acceleration</b>	<b>27</b>
3.1	A Review of Linear Algebra . . . . .	27
3.2	Characterizations of Convexity and Optimality via Second Derivatives . . . .	29
3.2.1	A Necessary Condition for Local Extrema . . . . .	30
3.2.2	A sufficient condition for local extrema . . . . .	31

3.2.3	Characterization of convexity . . . . .	31
3.3	Gradient Descent - An Approach to Optimization? . . . . .	33
3.3.1	A Quantitative Bound on Changes in the Gradient . . . . .	33
3.3.2	Analyzing Gradient Descent . . . . .	34
3.4	Accelerated Gradient Descent . . . . .	36
<b>II</b>	<b>Spectral Graph Theory</b>	<b>42</b>
<b>4</b>	<b>Introduction to Spectral Graph Theory</b>	<b>43</b>
4.1	Recap: Incidence and Adjacency Matrices, the Laplacian Matrix and Electrical Energy . . . . .	43
4.2	Understanding Eigenvalues of the Laplacian . . . . .	45
4.2.1	Test Vector Bounds on $\lambda_2$ and $\lambda_n$ . . . . .	46
4.2.2	Eigenvalue Bounds Beyond Test Vectors . . . . .	47
4.2.3	The Loewner Order, aka. the Positive Semi-Definite Order . . . . .	48
4.2.4	Upper Bounding a Laplacian's $\lambda_n$ Using Degrees . . . . .	49
4.2.5	The Loewner Order and Laplacians of Graphs. . . . .	49
4.2.6	The Path Inequality . . . . .	50
4.2.7	Lower Bounding $\lambda_2$ of a Path Graph . . . . .	51
4.2.8	Laplacian Eigenvalues of the Complete Binary Tree . . . . .	52
<b>5</b>	<b>Conductance, Expanders and Cheeger's Inequality</b>	<b>54</b>
5.1	Conductance and Expanders . . . . .	55
5.2	A Lower Bound for Conductance via Eigenvalues . . . . .	56
5.3	An Upper Bound for Conductance via Eigenvalues . . . . .	58
5.4	Conclusion . . . . .	61
<b>6</b>	<b>Random Walks</b>	<b>62</b>
6.1	A Primer on Random Walks . . . . .	62
6.2	Convergence Results for Random Walks . . . . .	63

6.2.1	Making Random Walks Lazy . . . . .	64
6.2.2	Convergence of Lazy Random Walks . . . . .	65
6.2.3	The Rate of Convergence . . . . .	66
6.3	Properties of Random Walks . . . . .	67
6.3.1	Hitting Times . . . . .	68
6.3.2	Commute Time . . . . .	69
<b>7</b>	<b>Pseudo-inverses and Effective Resistance</b>	<b>71</b>
7.1	What is a (Moore-Penrose) Pseudoinverse? . . . . .	71
7.2	Electrical Flows Again . . . . .	72
7.3	Effective Resistance . . . . .	73
7.3.1	Effective Resistance is a Distance . . . . .	76
<b>8</b>	<b>Different Perspectives on Gaussian Elimination</b>	<b>77</b>
8.1	An Optimization View of Gaussian Elimination for Laplacians . . . . .	77
8.2	An Additive View of Gaussian Elimination . . . . .	80
<b>9</b>	<b>Random Matrix Concentration and Spectral Graph Sparsification</b>	<b>84</b>
9.1	Matrix Sampling and Approximation . . . . .	84
9.2	Matrix Concentration . . . . .	87
9.2.1	Matrix Functions . . . . .	88
9.2.2	Monotonicity and Operator Monotonicity . . . . .	88
9.2.3	Some Useful Facts . . . . .	90
9.2.4	Proof of Matrix Bernstein Concentration Bound . . . . .	91
9.3	Spectral Graph Sparsification . . . . .	92
<b>10</b>	<b>Solving Laplacian Linear Equations</b>	<b>99</b>
10.1	Solving Linear Equations Approximately . . . . .	99
10.2	Preconditioning and Approximate Gaussian Elimination . . . . .	100
10.3	Approximate Gaussian Elimination Algorithm . . . . .	101
10.4	Analyzing Approximate Gaussian Elimination . . . . .	104

10.4.1	Normalization, a.k.a. Isotropic Position . . . . .	105
10.4.2	Martingales . . . . .	105
10.4.3	Martingale Difference Sequence as Edge-Samples . . . . .	107
10.4.4	Stopped Martingales . . . . .	108
10.4.5	Sample Norm Control . . . . .	109
10.4.6	Random Matrix Concentration from Trace Exponentials . . . . .	111
10.4.7	Mean-Exponential Bounds from Variance Bounds . . . . .	112
10.4.8	The Overall Mean-Trace-Exponential Bound . . . . .	113

# 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 making it safe to do math.

Recall that a typical battery that you buy 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.

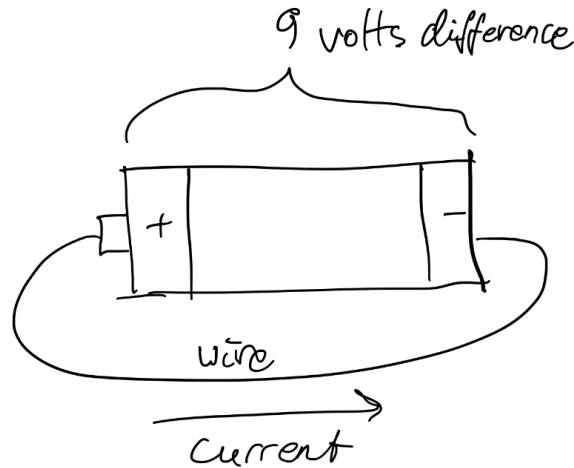


Figure 1.1: A 9 volts battery with a wire attached.

We can also imagine a kind of battery that tries to send a certain amount of current 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.

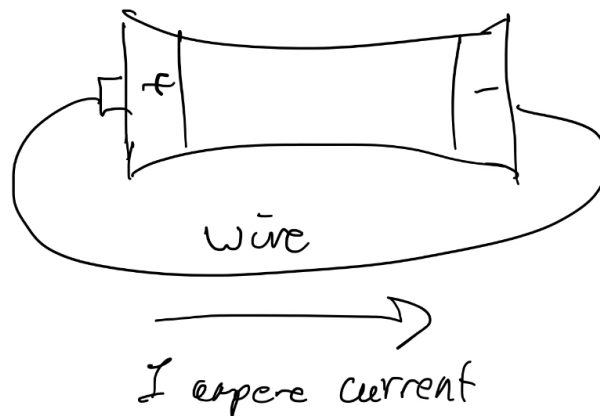


Figure 1.2: A 1 ampere battery with a wire attached.

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.

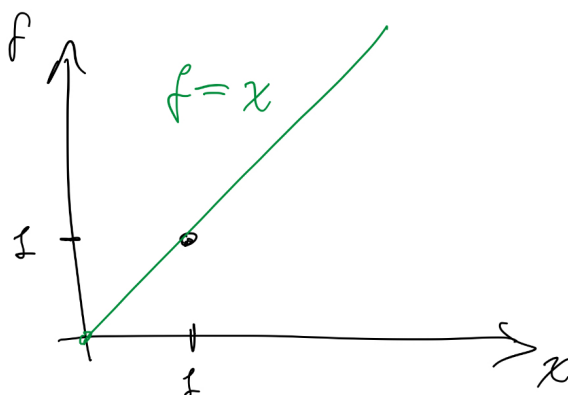


Figure 1.3: Ohm's Law for a resistor with resistance  $r = 1$ .

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

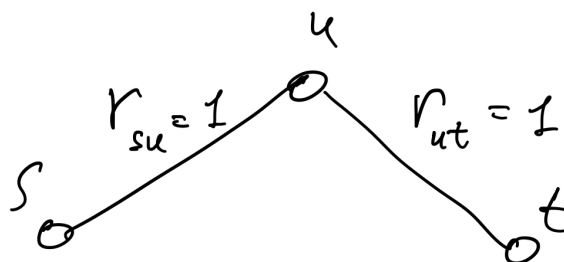


Figure 1.4: A path of two resistors.

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 Kirchoff'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  $x = 0$ .

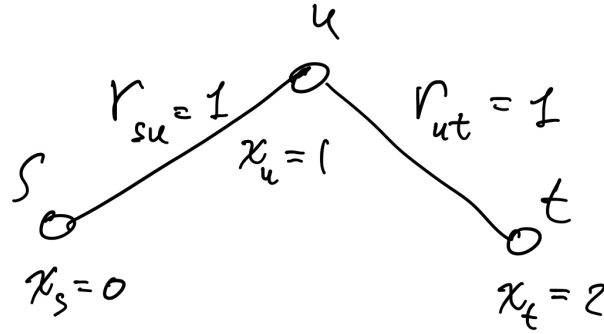


Figure 1.5: A path of two resistors.

Let us try one more example:

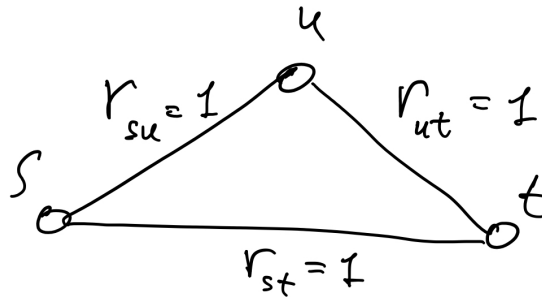


Figure 1.6: A network with three resistors.

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 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  $a = x + 0$ ,  $x_u = x + \frac{1}{3}$ ,  $x_t = x + \frac{2}{3}$ . 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 a 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  $\mathbf{x} \in \mathbb{R}^V$ . Ohm's Law says that the electrical flow induced by these voltages will be  $\mathbf{f}(u, v) = \frac{1}{r(u, v)}(\mathbf{x}(u) - \mathbf{x}(v))$ .

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

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”.

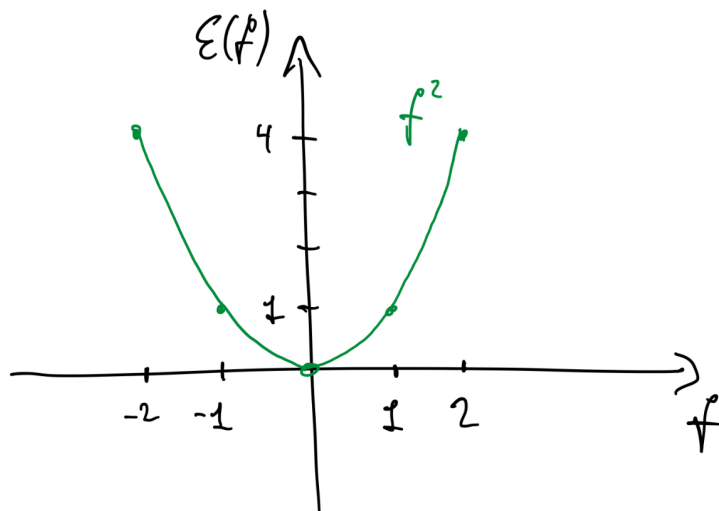


Figure 1.7: Energy has a function of flow in a resistor with resistance  $r = 1$ .

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? We can phrase this as an optimization problem:

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

We call this problem *electrical energy-minimizing flow*. It turns out, that the flow  $\mathbf{f}^*$  that minimizes the electrical energy among all flows that satisfy  $\mathbf{B}\mathbf{f} = \mathbf{d}$  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  $\mathbf{L}\mathbf{x} = \mathbf{d}$ . Consider the function  $c(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{L}\mathbf{x} - \mathbf{x}^\top \mathbf{d}$ . 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  $\mathbf{L}\mathbf{x} = \mathbf{d}$ . You will confirm this in Exercise 4.

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

$$\min_{\mathbf{x} \in \mathbb{R}^V} \frac{1}{2} \mathbf{x}^\top \mathbf{L}\mathbf{x} - \mathbf{x}^\top \mathbf{d}$$

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  $\mathbf{L}\mathbf{x} = \mathbf{d}$  are in fact minimizers of  $c$ . The answer is yes, and 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  $\mathbf{f}^* = \mathbf{R}^{-1} \mathbf{B}^\top \mathbf{x}^*$ . In Exercise 5, you will explore this relationship more.

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.

**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.

Either, 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.

---

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

One of the most basic algorithms in this framework of “guess and adjust” is called *Gradient Descent*, which we will study in two weeks. 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 developing 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  $O(n^3 m)$  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.

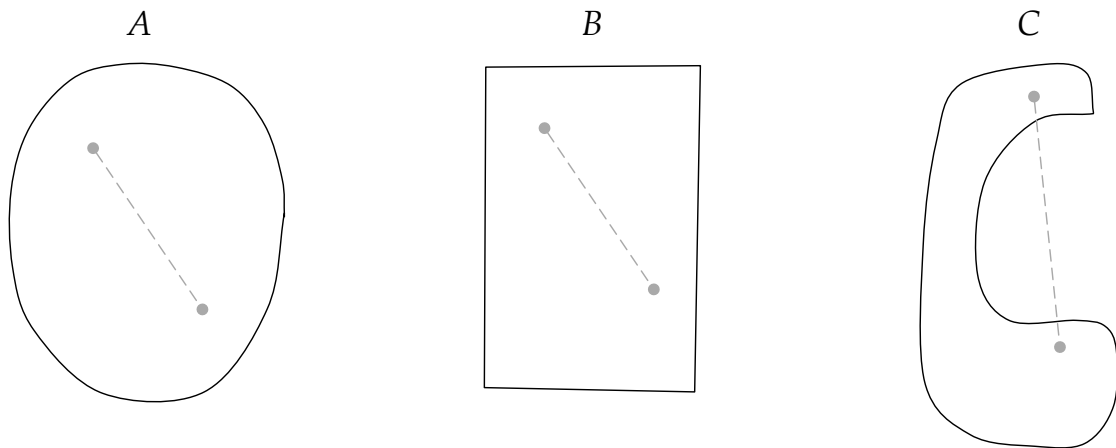


Figure 1.9: A depiction of convex and non-convex sets. The sets  $A$  and  $B$  are convex since the straight line between any two points inside them is also in the set. The set  $C$  is not convex.

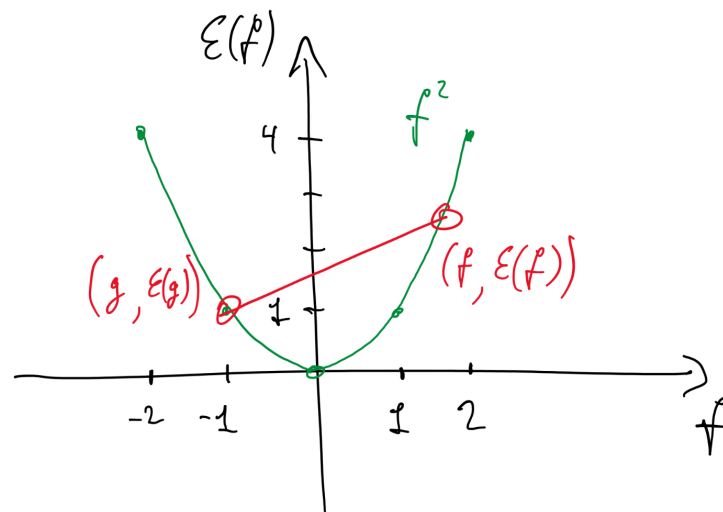


Figure 1.8: Energy has a function of flow in a resistor with resistance  $r = 1$ . The function is convex.

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 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  $\mathbf{d}$  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

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

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  $\mathbf{f}^*$  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 another algorithm, that finds an optimal flow  $\mathbf{f}^*$  *exactly*, in time<sup>5</sup>  $m^{10/7} \log^{O(1)} n$ .

---

<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.

<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.

<sup>5</sup>And there's even a paper on arXiv.org that brings this further down to  $m^{11/8} \log^{O(1)} n$ .



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 in both directions to exist between a vertex to exist, and we require that that 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 require now try to minimize  $\|\mathbf{C}^{-1}\mathbf{f}\|_\infty$ , where  $\mathbf{C} = \text{diag}_{e \in E} \mathbf{c}(e)$ . Then our problem becomes

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

For this capacitated, directed maximum flow problem, our best algorithms run in about  $O(m\sqrt{n})$  time<sup>6</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 even harder still, by simultaneously trying to route to types of flow (imagine pipes with Coke and Pepsi). Our problem now looks like

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

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.

**Random walks in a graph.** Google famously uses<sup>7</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 edge 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

<sup>6</sup>Provided the capacities are integers satisfying a condition like  $\mathbf{c} \leq n^{100}\mathbf{1}$ .

<sup>7</sup>At least they did at some point.

set by starting at some vertex  $u_0 = u$ , and then randomly picking 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>8</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>8</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:

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

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

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

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**.

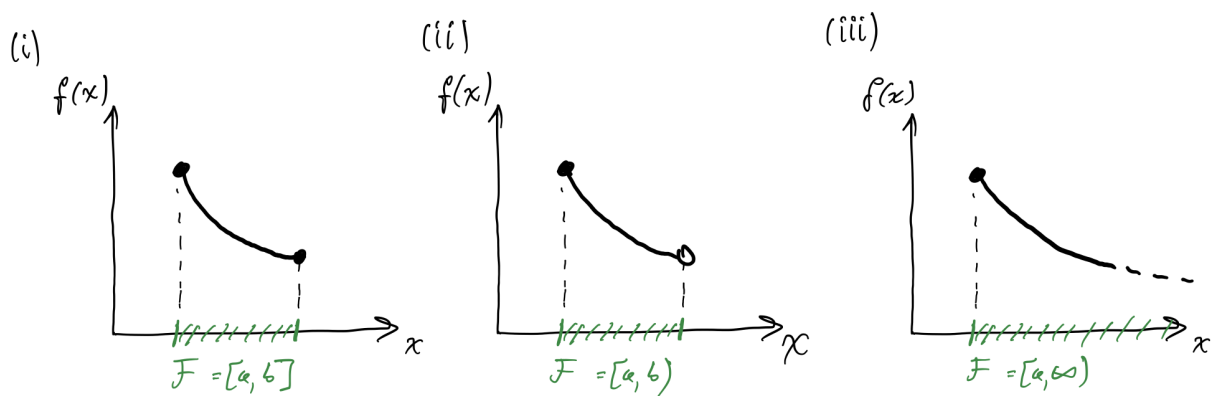


Figure 2.1

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:

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

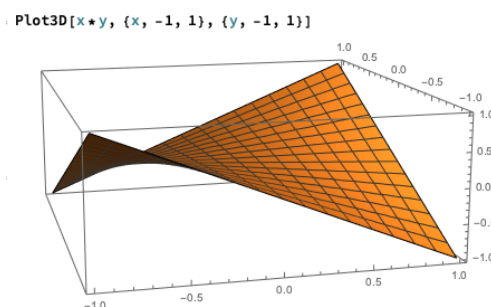


Figure 2.2: This plot shows the function  $f(x, y) = xy$ . For any fixed  $y_0$ , the function  $h(x) = f(x, y_0) = xy_0$  is linear in  $x$ , and so is a convex function in  $x$ . But is  $f$  convex?

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 functions  $f$  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  $\nabla f(\mathbf{x})$  is:

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

**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.

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

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

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). Rockafellar'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.

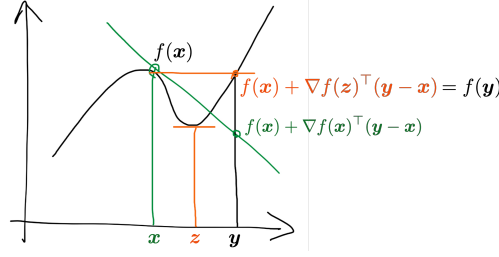


Figure 2.3: The convex function  $f(\mathbf{y})$  sits above the linear function in  $\mathbf{y}$  given by  $f(\mathbf{x}) + \nabla f(\mathbf{x})^\top(\mathbf{y} - \mathbf{x})$ .

### 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  $\|\lambda \mathbf{d}\|_2 = \lambda \|\mathbf{d}\|_2$ ):

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

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



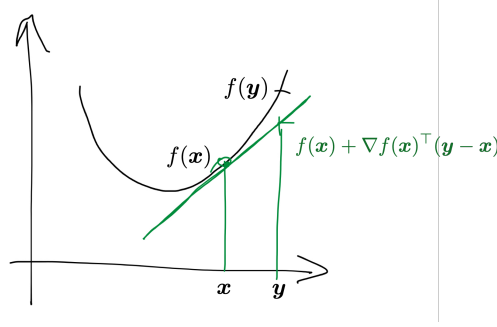


Figure 2.4: The convex function  $f(\mathbf{y})$  sits above the linear function in  $\mathbf{y}$  given by  $f(\mathbf{x}) + \nabla f(\mathbf{x})^\top(\mathbf{y} - \mathbf{x})$ .

*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  $\mathbf{d} = \mathbf{x} - \mathbf{y}$  we have that:

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

[  $\impliedby$  ] 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:

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

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

Observe that  $\mathbf{y} - \mathbf{z} = (1 - \theta)(\mathbf{y} - \mathbf{x})$  and  $\mathbf{x} - \mathbf{z} = \theta(\mathbf{y} - \mathbf{x})$ . 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. □

## 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. □

**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 3.3.5 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. □

# Chapter 3

## Convexity and Second Derivatives, Gradient Descent and Acceleration

**Notation for this chapter.** In this chapter, we sometimes consider a multivariate functions  $f$  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  $\mathbf{A}$  be 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  $\mathbf{\Lambda} \in \mathbb{R}^{n \times n}$  s.t.*

1.  $\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$ .
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.  $\mathbf{A}_{ii} = \lambda_i(\mathbf{A})$ .

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.

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

2.

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

Theorem 3.1.2 is an immediate corollary of Theorem 4.1.1, 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 4.1.1.* 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  $\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$  where  $\mathbf{\Lambda}$  is a diagonal matrix of eigenvalues of  $\mathbf{A}$ , which we take to be sorted in increasing order. Then  $\mathbf{x}^\top \mathbf{A} \mathbf{x} = \mathbf{x}^\top \mathbf{V}^\top \mathbf{\Lambda} \mathbf{V} \mathbf{x} = (\mathbf{V} \mathbf{x})^\top \mathbf{\Lambda} (\mathbf{V} \mathbf{x}) = \sum_{j=1}^i \lambda_j \mathbf{c}(j)^2 \leq \lambda_i \|\mathbf{c}\|_2^2 = \lambda_i$ . 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,

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

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:

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

**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  $\mathbf{H}_f(\mathbf{x}) \in \mathbb{R}^{n \times n}$  s.t.

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

This is equivalent to saying that for all  $\delta$

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

where by definition:

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

We say that  $f$  is *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  $\mathbf{z} \in [\mathbf{x}, \mathbf{y}]$ :

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

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 3.3.5 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 3.3.5 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 3.3.5 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 1 of this week's exercises.

**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  $\mathbf{x}_\theta - \mathbf{x} = \theta(\mathbf{y} - \mathbf{x})$ .

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

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

□

### 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  $f(\mathbf{x}_1) \leq f(\mathbf{x}_0) - \frac{\|\nabla f(\mathbf{x}_0)\|_2^2}{2\beta}$ .

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 Week 2 exercises (Exercise 3).

## 3.4 Accelerated Gradient Descent

It turns out that we can get an algorithm that converges substantially faster than Gradient Descent, using an approach known as *Accelerated Gradient Descent*, which was developed by Nesterov [Nes83]. This algorithm in turn improved on some earlier results by Nemirovski and Yudin [NY83]. The phenomenon of acceleration was perhaps first understood in the

context of quadratic functions, minimizing  $\mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{x}^\top \mathbf{b}$  when  $\mathbf{A}$  is positive definite – for this case, the Conjugate Gradient algorithm was developed independently by Hestenes and Stiefel [HS<sup>+</sup>52] (here at ETH!), and by Lanczos [Lan52]. In the past few years, providing more intuitive explanations of acceleration has been a popular research topic. This presentation is based on an analysis of Nesterov’s algorithm developed by Diakonikolas and Orecchia [DO19].

We will adopt a slightly different approach to analyzing this algorithm than what we used in the previous section for Gradient Descent.

We will use  $\mathbf{x}_0$  to denote the starting point of our algorithm, and we will produce a sequence of iterates  $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$ . At each iterate  $\mathbf{x}_i$ , we will compute the gradient  $\nabla f(\mathbf{x}_i)$ . However, the way we choose  $\mathbf{x}_{i+1}$  based on what we know so far will now be a little more involved than what we did for Gradient Descent.

To help us understand the algorithm, we are going to introduce two more sequences of iterates  $\mathbf{y}_0, \mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_k \in \mathbb{R}^n$  and  $\mathbf{v}_0, \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k \in \mathbb{R}^n$ .

The sequence of  $\mathbf{y}_i$ ’s will be constructed to help us get as low a function value as possible at  $f(\mathbf{y}_i)$ , which we will consider our current solution and the last iterate  $\mathbf{y}_k$  will be the output solution of our algorithm.

The sequence of  $\mathbf{v}_i$ ’s will be constructed to help us get a lower bound on  $f(\mathbf{x}^*)$ .

By combining the upper bound on the function value of our current solution  $f(\mathbf{y}_i)$  with a lower bound on the function value at an optimal solution  $f(\mathbf{x}^*)$ , we get an upper bound on the gap  $f(\mathbf{y}_i) - f(\mathbf{x}^*)$  between the value of our solution and the optimal one. Finally, each iterate  $\mathbf{x}_i$ , which will be where we evaluate gradient  $\nabla f(\mathbf{x}_i)$ , is chosen through a trade-off between wanting to reduce the upper bound and wanting to increase the lower bound.

**The upper bound sequence:  $\mathbf{y}_i$ ’s.** The point  $\mathbf{y}_i$  will be chosen from  $\mathbf{x}_i$  to minimize an upper bound on  $f$  based at  $\mathbf{x}_i$ . This is similar to what we did in the previous section. We let  $\mathbf{y}_i = \mathbf{x}_i + \boldsymbol{\delta}_i$  and use choose  $\boldsymbol{\delta}_i$  to minimize the upper bound  $f(\mathbf{y}_i) \leq f(\mathbf{x}_i) + \nabla f(\mathbf{x}_i)^\top \boldsymbol{\delta}_i + \frac{\beta}{2} \|\boldsymbol{\delta}_i\|^2$ , which gives us

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

We will introduce a notation for this upper bound

$$U_i = f(\mathbf{y}_i) \leq f(\mathbf{x}_i) - \frac{\|\nabla f(\mathbf{x}_i)\|_2^2}{2\beta}. \quad (3.6)$$

**Philosophizing about lower bounds<sup>1</sup>.** A crucial ingredient to establishing an upper bound on gap <sub>$i$</sub>  was a lower bound on  $f(\mathbf{x}^*)$ .

---

<sup>1</sup>YMMV. People have a lot of different opinions about how to understand acceleration, and you should take my thoughts with a grain of salt.

In our analysis of Gradient Descent, in Equation (3.4), we used the lower bound  $f(\mathbf{x}^*) \geq f(\mathbf{x}_i) - \|\nabla f(\mathbf{x}_i)\|_2 \|\mathbf{x}_i - \mathbf{x}^*\|_2$ . We can think of the Gradient Descent analysis as being based on a tension between two statements: Firstly that “a large gradient implies we quickly approach the optimum” and secondly “the function value gap to optimum cannot exceed the magnitude of the current gradient (scaled by distance to opt)”.

This analysis does not use that we have seen many different function values and gradients, and each of these can be used to construct a lower bound on the optimum value  $f(\mathbf{x}^*)$ , and, in particular, it is not clear that the last gradient provides the best bound. To do better, we will try to use lower bounds that take advantage of all the gradients we have seen.

**Definition 3.4.1.** We will adopt a new notation for inner products that sometimes is more convenient when dealing with large expressions:  $\langle \mathbf{a}, \mathbf{b} \rangle \stackrel{\text{def}}{=} \mathbf{a}^\top \mathbf{b}$ .

**The lower bound sequence:  $\mathbf{v}_i$ 's.** We can introduce weights  $a_i > 0$  for each step and combine the gradients we have observed into one lower bound based on a weighted average. Let us use  $A_i = \sum_{j \leq i} a_j$  to denote the sum of the weights. Now a general lower bound on the function value at any  $\mathbf{v} \in \mathbb{R}^n$  is :

$$f(\mathbf{v}) \geq \frac{1}{A_i} \sum_{j \leq i} a_j (f(\mathbf{x}_j) + \langle \nabla f(\mathbf{x}_j), \mathbf{v} - \mathbf{x}_j \rangle)$$

However, to use Cauchy-Schwarz on each individual term here to instantiate this bound at  $\mathbf{x}^*$  does not give us anything useful. Instead, we will employ a somewhat magical trick: we introduce a regularization term

$$\phi(\mathbf{v}) \stackrel{\text{def}}{=} \frac{\sigma}{2} \|\mathbf{v} - \mathbf{x}_0\|_2^2.$$

We will choose the value  $\sigma > 0$  later. Now we derive our lower bound  $L_i$

$$\begin{aligned} f(\mathbf{x}^*) &\geq \frac{1}{A_i} \left( \phi(\mathbf{x}^*) + \sum_{j \leq i} a_j f(\mathbf{x}_j) + \langle a_j \nabla f(\mathbf{x}_j), \mathbf{x}^* - \mathbf{x}_j \rangle \right) - \frac{\phi(\mathbf{x}^*)}{A_i} \\ &\geq \min_{\mathbf{v} \in \mathbb{R}^n} \left\{ \frac{1}{A_i} \left( \phi(\mathbf{v}) + \sum_{j \leq i} a_j f(\mathbf{x}_j) + \langle a_j \nabla f(\mathbf{x}_j), \mathbf{v} - \mathbf{x}_j \rangle \right) \right\} - \frac{\phi(\mathbf{x}^*)}{A_i} \\ &= L_i \end{aligned}$$

We will let  $\mathbf{v}_i$  be the  $\mathbf{v}$  obtaining the minimum in the optimization problem appearing in the definition of  $L_i$ , so that

$$L_i = \frac{1}{A_i} \left( \phi(\mathbf{v}_i) + \sum_{j \leq i} a_j f(\mathbf{x}_j) + \langle a_j \nabla f(\mathbf{x}_j), \mathbf{v}_i - \mathbf{x}_j \rangle \right) - \frac{\phi(\mathbf{x}^*)}{A_i}$$

**How we will measure convergence.** We have designed the upper bound  $U_i$  and the lower bound  $L_i$  such that  $\text{gap}_i = f(\mathbf{y}_i) - f(\mathbf{x}^*) \leq U_i - L_i$ .

As you will show in Exercise 3, we can prove the convergence of Gradient Descent directly by an induction that establishes  $1/\text{gap}_i \leq C \cdot i$  for some constant  $C$  depending on the Lipschitz gradient parameter  $\beta$  and the distance  $\|\mathbf{x}_0 - \mathbf{x}^*\|_2$ .

To analyze Accelerated Gradient Descent, we will adopt a similar, but slightly different strategy, namely trying to show that  $(U_i - L_i)r(i)$  is non-increasing for some positive “rate function”  $r(i)$ . Ideally  $r(i)$  should grow quickly, which would imply that  $\text{gap}_i$  quickly gets small. We will also need to show that  $(U_0 - L_0)r(0) \leq C$  for some constant  $C$  again depending on  $\beta$  and  $\|\mathbf{x}_0 - \mathbf{x}^*\|_2$ . Then, we’ll be able to conclude that

$$\text{gap}_i \cdot r(i) \leq (U_i - L_i)r(i) \leq (U_{i-1} - L_{i-1})r(i-1) \leq \dots \leq (U_0 - L_0)r(0) \leq C,$$

and hence  $\text{gap}_i \leq C/r(i)$ .

This framework is fairly general. We could have also used it to analyze Gradient Descent, and it works for many other optimization algorithms too.

We are going to choose our rate function  $r(i)$  to be exactly  $A_i$ , which of course is no accident! As we will see, this interacts nicely with our lower bound  $L_i$ . Hence, our goals are to

1. provide an upper bound on  $A_0(U_0 - L_0)$ ,
2. and show that  $A_{i+1}(U_{i+1} - L_{i+1}) \leq A_i(U_i - L_i)$ ,

**Establishing the convergence rate.** Let’s start by looking at the change in the upper bound scaled by our rate function:

$$A_{i+1}U_{i+1} - A_iU_i = A_{i+1}(f(\mathbf{y}_{i+1}) - f(\mathbf{x}_{i+1})) - A_i(f(\mathbf{y}_i) - f(\mathbf{x}_{i+1})) + (A_{i+1} - A_i)f(\mathbf{x}_{i+1}) \quad (3.7)$$

$$\begin{aligned} &\leq A_{i+1} \left( -\frac{\|\nabla f(\mathbf{x}_{i+1})\|_2^2}{2\beta} \right) && \text{First term controlled by Equation (3.6).} \\ &\quad - A_i \langle \nabla f(\mathbf{x}_{i+1}), \mathbf{y}_i - \mathbf{x}_{i+1} \rangle && \text{Second term bounded by Theorem 3.3.5.} \\ &\quad + a_{i+1}f(\mathbf{x}_{i+1}) && \text{Third term uses } a_{i+1} = A_{i+1} - A_i. \end{aligned}$$

The solution  $\mathbf{v}_i$  to the minimization in the lower bound  $L_i$  turns out to be relatively simple to characterize. By using derivatives to find the optimum, we first analyze the initial value of the lower bound  $L_0$ .

**Claim 3.4.2.**

1.  $\mathbf{v}_0 = \mathbf{x}_0 - \frac{a_0}{\sigma} \nabla f(\mathbf{x}_0)$
2.  $L_0 = f(\mathbf{x}_0) - \frac{a_0}{2\sigma} \|\nabla f(\mathbf{x}_0)\|_2^2 - \frac{\sigma}{2a_0} \|\mathbf{x}^* - \mathbf{x}_0\|_2^2.$

You will prove Claim 3.4.2 in the exercises for Week 3. Noting  $A_0 = a_0$ , we see from Equation (3.6) and Part 2 of Claim 3.4.2, that

$$A_0(U_0 - L_0) \leq \left( \frac{a_0^2}{2\sigma} - \frac{a_0}{2\beta} \right) \|\nabla f(\mathbf{x}_0)\|_2^2 + \frac{\sigma}{2} \|\mathbf{x}^* - \mathbf{x}_0\|_2^2 \quad (3.8)$$

It will be convenient to introduce notation for the rescaled lower bound  $A_i L_i$  *without* optimizing over  $\mathbf{v}$ .

$$m_i(\mathbf{v}) = \phi(\mathbf{v}) - \phi(\mathbf{x}^*) + \sum_{j \leq i} a_j f(\mathbf{x}_j) + \langle a_j \nabla f(\mathbf{x}_j), \mathbf{v} - \mathbf{x}_j \rangle$$

Thus  $A_i L_i - A_{i+1} L_{i+1} = m_i(\mathbf{v}_i) - m_{i+1}(\mathbf{v})$ . Now, it is not too hard to show the following relationships.

**Claim 3.4.3.**

1.  $m_i(\mathbf{v}) = m_i(\mathbf{v}_i) + \frac{\sigma}{2} \|\mathbf{v} - \mathbf{v}_i\|_2^2$
2.  $m_{i+1}(\mathbf{v}) = m_i(\mathbf{v}) + a_{i+1} f(\mathbf{x}_{i+1}) + \langle a_{i+1} \nabla f(\mathbf{x}_{i+1}), \mathbf{v} - \mathbf{x}_{i+1} \rangle$
3.  $\mathbf{v}_{i+1} = \mathbf{v}_i - \frac{a_{i+1}}{\sigma} \nabla f(\mathbf{x}_{i+1})$

And again, you will prove Claim 3.4.3 in the exercises for Week 3. *Hint for Part 1: note that  $m_i(\mathbf{v})$  is a quadratic function, minimized at  $\mathbf{v}_i$  and its Hessian equals  $\sigma I$  at all  $\mathbf{v}$ .*

Given Claim 3.4.3, we see that

$$\begin{aligned} A_i L_i - A_{i+1} L_{i+1} &= m_i(\mathbf{v}_i) - m_{i+1}(\mathbf{v}_{i+1}) \end{aligned} \quad (3.9)$$

$$= -a_{i+1} f(\mathbf{x}_{i+1}) - \langle a_{i+1} \nabla f(\mathbf{x}_{i+1}), \mathbf{v}_{i+1} - \mathbf{x}_{i+1} \rangle - \frac{\sigma}{2} \|\mathbf{v}_{i+1} - \mathbf{v}_i\|_2^2 \quad (3.10)$$

$$= -a_{i+1} f(\mathbf{x}_{i+1}) - \langle a_{i+1} \nabla f(\mathbf{x}_{i+1}), \mathbf{v}_i - \mathbf{x}_{i+1} \rangle + \frac{a_{i+1}^2}{2\sigma} \|\nabla f(\mathbf{x}_{i+1})\|_2^2 \quad (3.11)$$

This means that by combining Equation (3.7) and (3.11) we get

$$\begin{aligned} A_{i+1}(U_{i+1} - L_{i+1}) - A_i(U_i - L_i) &\leq \left( \frac{-A_{i+1}}{2\beta} + \frac{a_{i+1}^2}{2\sigma} \right) \|\nabla f(\mathbf{x}_{i+1})\|_2^2 \\ &\quad + \langle \nabla f(\mathbf{x}_{i+1}), A_{i+1} \mathbf{x}_{i+1} - a_{i+1} \mathbf{v}_i - A_i \mathbf{y}_i \rangle. \end{aligned}$$



Now, this means that  $A_{i+1}(U_{i+1} - L_{i+1}) - A_i(U_i - L_i) \leq 0$  if

$$A_{i+1}\mathbf{x}_{i+1} - a_{i+1}\mathbf{v}_i - A_i\mathbf{y}_i = \mathbf{0} \text{ and } A_{i+1}/\beta \geq a_{i+1}^2/\sigma$$

We can get this by letting  $\mathbf{x}_{i+1} = \frac{A_i\mathbf{y}_i + a_{i+1}\mathbf{v}_i}{A_{i+1}}$ , and  $\sigma = \beta$  and  $a_i = \frac{i+1}{2}$ , which implies that  $A_i = \frac{(i+1)(i+2)}{4} > a_i^2$ .

By Equation (3.8), these parameter choices also imply that

$$A_0(U_0 - L_0) \leq \frac{\beta}{2} \|\mathbf{x}^* - \mathbf{x}_0\|_2^2.$$

Finally, by induction, we get  $A_i(U_i - L_i) \leq \frac{\beta}{2} \|\mathbf{x}^* - \mathbf{x}_0\|_2^2$ . Dividing through by  $A_i$  and using  $\text{gap}_i \leq U_i - L_i$  results in the following theorem.

**Theorem 3.4.4.** *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 Accelerated Gradient Descent algorithm given by*

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

*ensures that the  $k$ th iterate satisfies*

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

# 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

In Chapter 1, we looked at undirected graphs and we introduce the incidence matrix and the Laplacian of the graph. Let us recall these.

We consider an undirected weighted graph  $G = (V, E, \mathbf{w})$ , with  $n = |V|$  vertices and  $m = |E|$  edges, where  $\mathbf{w} \in \mathbb{R}_+^E$  assigns positive weight for every edge. Let's assume  $G$  is connected.

To introduce the *edge-vertex incidence matrix* of the graph, 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.

Let  $\mathbf{W} \in \mathbb{R}^{E \times E}$  be the diagonal matrix given by  $\mathbf{W} = \text{diag}(\mathbf{w})$ , i.e  $\mathbf{W}(e, e) = \mathbf{w}(e)$ . 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  $w(e) = 1/R(e, e)$ .

The weighted adjacency matrix  $\mathbf{A} \in \mathbb{R}^{V \times V}$  of a graph is given by

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

Note that we treat the edges as undirected here, so  $\mathbf{A}^\top = \mathbf{A}$ . The weighted degree of a vertex is defined as  $d(v) = \sum_{\{u, v\} \in E} w(u, v)$ . Again we treat the edges as undirected. Let  $\mathbf{D} = \text{diag}(\mathbf{d})$  be the diagonal matrix in  $\mathbb{R}^{V \times V}$  with weighted degrees on the diagonal.

In Problem Set 1, you showed that  $\mathbf{L} = \mathbf{D} - \mathbf{A}$ , and that for  $\mathbf{x} \in \mathbb{R}^V$ ,

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

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 electrical voltages  $\mathbf{x}$  will induce an electrical flow  $\mathbf{f} = \mathbf{R}^{-1} \mathbf{B}^\top \mathbf{x}$ . We defined the electrical energy of a flow  $\mathbf{f} \in \mathbb{R}^E$  to be

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

And, from Ohm's Law, we can then see that

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

Hence, define the electrical energy associated with a set of voltages to be

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

**The Courant-Fisher Theorem.** Let us also recall the Courant-Fischer theorem, which we proved in Chapter 3 (Theorem 4.1.1).

**Theorem 4.1.1** (The Courant-Fisher 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.

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

2.

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

In fact, from our proof of the Courant-Fischer theorem in Chapter 3, we can also extract a slightly different statement:

**Theorem 4.1.2** (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.

$$\lambda_i = \min_{\substack{\mathbf{x} \perp \mathbf{x}_1, \dots, \mathbf{x}_{i-1} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$$

2.

$$\lambda_i = \max_{\substack{\mathbf{x} \perp \mathbf{x}_{i+1}, \dots, \mathbf{x}_n \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$$

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 off with 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*

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

**Eigenvalues of the Laplacian of a 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  $\lambda_2 = \lambda_3 = \dots = \lambda_n = n$ .

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.

#### 4.2.1 Test Vector Bounds on $\lambda_2$ and $\lambda_n$

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  $\frac{\mathbf{x}^\top \mathbf{L} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$ , by instead replacing  $\mathbf{x}$  with the component orthogonal to  $\mathbf{x}$ , which we denoted by  $\mathbf{y}$ .

We can conclude from Equation (4.2) that for *any* vector  $\mathbf{y} \perp \mathbf{1}$ ,

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

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  $\lambda_2(\mathbf{L})$ .

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

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

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

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

This means 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  $\mathbf{x} \in \mathbb{R}^V$  be given by  $\mathbf{x}(1) = -1$ , and  $\mathbf{x}(n) = 1$  and  $\mathbf{x}(i) = 0$  for  $i \neq 0, 1$ .

Here we get

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

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.

## 4.2.2 Eigenvalue Bounds Beyond Test Vectors

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  $n$ .

We then looked at eigenvalues of  $P_n$ , the unit weight path on  $n$  vertices, and we showed using *test vector* bounds that

$$\lambda_2(\mathbf{L}_{P_n}) \leq \frac{12}{n^2} \text{ and } 1 \leq \lambda_n(\mathbf{L}_{P_n}). \quad (4.3)$$

Ideally we would like to prove an almost matching upper bound on  $\lambda_2$  and an almost matching lower bound on  $\lambda_n$ , but it is not clear how to get that from the Courant-Fischer theorem.

To get there, we start we need to introduce some more tools.

### 4.2.3 The Loewner Order, aka. the Positive Semi-Definite Order

We'll 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  $\mathbf{B} \succeq \mathbf{A}$ .

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 4.2.2.** *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.

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

□

Note that the converse of Claim 4.2.2 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 4.2.3.** 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.

#### 4.2.4 Upper Bounding a Laplacian's $\lambda_n$ 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 4.2.2

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

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

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

#### 4.2.5 The Loewner Order and Laplacians of Graphs.

It's sometimes convenient to overload the for 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  $H = (V, F)$  is a subgraph of  $G$ , then

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

We can see this from the Laplacian quadratic form:

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

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.

## 4.2.6 The Path Inequality

Now, we'll see a general tool for comparing two graphs  $G$  and  $H$  to prove an 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 define 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 4.2.5** (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  $\mathbf{a}^\top \mathbf{b} \leq \|\mathbf{a}\|_2 \|\mathbf{b}\|_2$ :

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

□

#### 4.2.7 Lower Bounding $\lambda_2$ of a Path Graph

We will now use Lemma 4.2.5 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 4.2.5:

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

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 consirably tighter by being more careful about the sums.

## 4.2.8 Laplacian Eigenvalues of the Complete Binary Tree

Let's do the same analysis with the complete binary tree with unit weight edges on  $n = 2^{d+1} - 1$  vertices, which we denote by  $T_d$ .

$T_d$  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.

**The largest eigenvalue.** We'll start by above bounding  $\lambda_n(\mathbf{L}_{T_d})$  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 4.2.4,  $\lambda_n(\mathbf{L}) \leq 6$ . So we can bound the largest eigenvalue above and below by constant.

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

**Lemma 4.2.6.** *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.  $\square$

**$\lambda_2$  in a tree.** Since a complete binary tree  $T_d$  has diameter  $2d \leq 2 \log_2(n)$ , by Lemma 4.2.6,  $\lambda_2(\mathbf{L}_{T_d}) \geq \frac{1}{2n \log_2(n)}$ .

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

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

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

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

# Chapter 5

## 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).

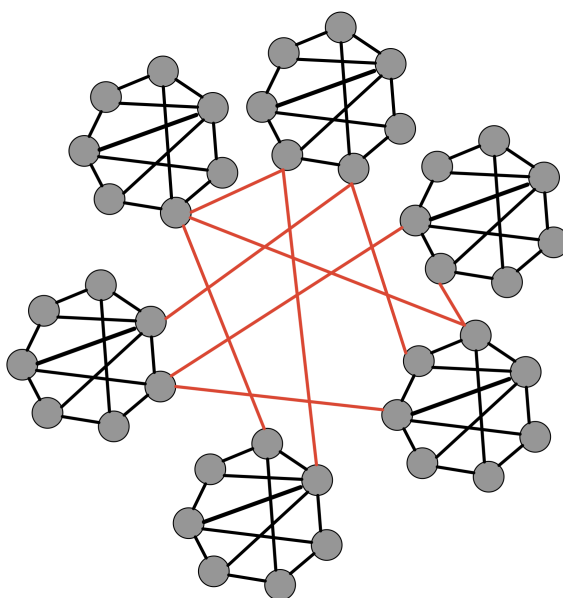


Figure 5.1: After removing the red edges (of which there are few in relation to the total number of edges), each connected component in  $G$  is "well-connected".

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

Building on the last two lectures, 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!

## 5.1 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  $\phi(S) \leq \sqrt{2\phi}$ .

In the graded homework, we ask you to make the procedure  $\text{CERTIFYORCUT}(G, \phi)$  explicit, and then to show how to use it to compute a  $\phi$ -expander decomposition.

## 5.2 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

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

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

$$\phi(G) = \min_{\substack{\emptyset \subset S \subset V, \\ \text{vol}(S) \leq \text{vol}(V)/2}} \frac{\mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S}{\mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S} \quad (5.1)$$

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  $\mathbf{y} = \mathbf{D}^{1/2} \mathbf{1}$ , to get by Courant-Fischer (see Theorem 4.1.2) that

$$\lambda_1(\mathbf{N}) = \min_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \mathbf{N} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} \leq \frac{\mathbf{y}^\top \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2} \mathbf{y}}{\mathbf{y}^\top \mathbf{D} \mathbf{y}} = \frac{\mathbf{1}^\top \mathbf{L} \mathbf{1}}{\mathbf{y}^\top \mathbf{D} \mathbf{y}} = 0 \quad (5.2)$$

because  $\mathbf{D}^{-1/2} \mathbf{D}^{1/2} = \mathbf{I}$  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  $\mathbf{N}$ :

$$\lambda_2(\mathbf{N}) = \min_{\substack{\mathbf{x} \perp \mathbf{D}^{1/2} \mathbf{1} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}^\top \mathbf{N} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \min_{\substack{\mathbf{z} \perp \mathbf{d} \\ \mathbf{z} \neq \mathbf{0}}} \frac{\mathbf{z}^\top \mathbf{D}^{1/2} \mathbf{N} \mathbf{D}^{1/2} \mathbf{z}}{\mathbf{z}^\top \mathbf{D}^{1/2} \mathbf{D}^{1/2} \mathbf{z}} = \min_{\substack{\mathbf{z} \perp \mathbf{d} \\ \mathbf{z} \neq \mathbf{0}}} \frac{\mathbf{z}^\top \mathbf{L} \mathbf{z}}{\mathbf{z}^\top \mathbf{D} \mathbf{z}}. \quad (5.3)$$

**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 5.3 is a relaxation of the minimization problem 5.1 defining  $\phi(G)$ .

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

**Theorem 5.2.1** (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  $\text{vol}(S) \leq \text{vol}(V)/2$ .

□

### 5.3 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 5.3.1** (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  $\emptyset \subset S \subset V$ , such that

$$\frac{\mathbf{1}_S^\top \mathbf{L} \mathbf{1}_S}{\mathbf{1}_S^\top \mathbf{D} \mathbf{1}_S} \leq \sqrt{2 \cdot \frac{\mathbf{z}^\top \mathbf{L} \mathbf{z}}{\mathbf{z}^\top \mathbf{D} \mathbf{z}}}. \quad (5.4)$$

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

- we *renumber* 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

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

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

- 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 5.4 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 distribution with probability density function

$$p(t) = \begin{cases} 2 \cdot |t| & t \in [\mathbf{z}_{sc}(1), \mathbf{z}_{sc}(n)], \\ 0 & \text{otherwise.} \end{cases} \quad (5.5)$$

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

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 5.3.2.** 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  $\mathbf{z}_{sc}(i) \leq \mathbf{z}_{sc}(j)$ ,

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

Distinguishing by cases, we get

$$\text{sgn}(i) \cdot \mathbf{z}_{sc}(i)^2 - \text{sgn}(j) \cdot \mathbf{z}_{sc}(j)^2 = \begin{cases} |\mathbf{z}_{sc}(i)^2 - \mathbf{z}_{sc}(j)^2| & \text{sgn}(i) = \text{sgn}(j), \\ \mathbf{z}_{sc}(i)^2 + \mathbf{z}_{sc}(j)^2 & \text{otherwise.} \end{cases}$$

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

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

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

$$\mathbb{E}_\tau[|E(S_\tau, V \setminus S_\tau)|] \leq \sqrt{\mathbf{z}_{sc}^\top \mathbf{L} \mathbf{z}_{sc} \cdot 2 \mathbf{z}_{sc}^\top \mathbf{D} \mathbf{z}_{sc}} = \sqrt{2 \cdot \frac{\mathbf{z}_{sc}^\top \mathbf{L} \mathbf{z}_{sc}}{\mathbf{z}_{sc}^\top \mathbf{D} \mathbf{z}_{sc}}} \cdot \mathbf{z}_{sc}^\top \mathbf{D} \mathbf{z}_{sc} \quad (5.6)$$

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 5.6 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 5.3.2 only ensures our claim in expectation, this is already sufficient to conclude that there exists some set  $S$  that satisfies the same guarantees deterministically. 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.  $\square$

## 5.4 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 5.4.1** (Cheeger's Inequality). *We have  $\frac{\lambda_2(\mathbf{N})}{2} \leq \phi(G) \leq \sqrt{2 \cdot \lambda_2(\mathbf{N})}$ .*

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 6

## 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  $G = (V, E, w)$ .

### 6.1 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

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

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  $w(v, v_t)$  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

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

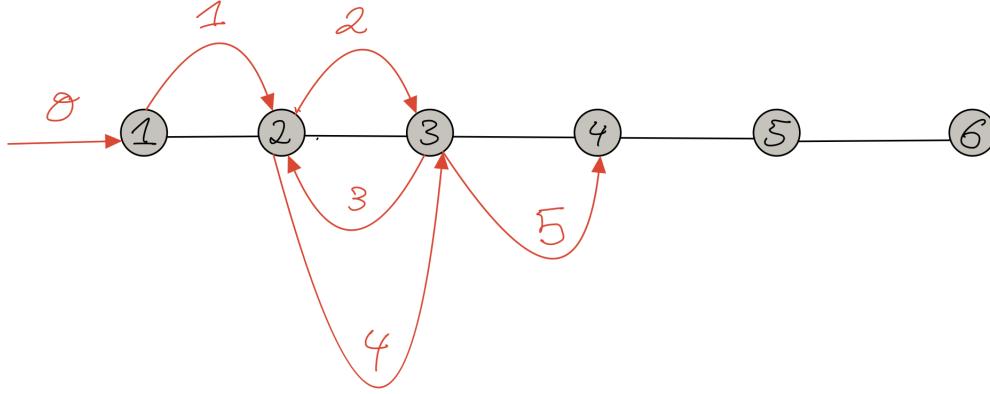


Figure 6.1: A (possibly random) walk where the red edges indicate the edges that the particle moves along. Here the walk visits the vertices  $v_0 = 1, v_1 = 2, v_2 = 3, v_3 = 2, v_4 = 3, v_5 = 4$ .

Thus,  $\mathbf{W}_{uv} = \mathbb{P}[v_{t+1} = v \mid v_t = u]$  (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  $\mathbf{p}_t = \mathbf{W}^t\mathbf{p}_0$ .

## 6.2 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 6.2.1.** *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

$$\mathbf{W}\pi = \mathbf{A}\mathbf{D}^{-1} \cdot \frac{\mathbf{d}}{\mathbf{1}^\top \mathbf{d}} = \frac{\mathbf{A}\mathbf{1}}{\mathbf{1}^\top \mathbf{d}} = \frac{\mathbf{d}}{\mathbf{1}^\top \mathbf{d}} = \pi.$$

□

For many graphs one can show that for  $t \rightarrow \infty$ , we have that  $\mathbf{p}_t \rightarrow \pi$ , i.e. that independent of the starting distribution  $\mathbf{p}_0$ , 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  $\mathbf{p}_0$  being 1 at one vertex and 0 at the other.

### 6.2.1 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

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

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

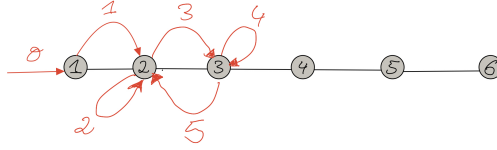


Figure 6.2: A lazy random walk where the red edges indicate the edges that the particle moves along. Here the lazy walk visits the vertices  $v_0 = 1, v_1 = 2, v_2 = 2, v_3 = 3, v_4 = 3, v_5 = 2$ .

**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 5.2 where we added a normalization such that  $\psi_1^\top \psi_1 = 1$ ).

**Lemma 6.2.2.** *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 6.2.3.** *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 6.2.2, the corollary follows. □

## 6.2.2 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  $\mathbf{N}$

$$\forall i, \psi_i^\top \mathbf{D}^{-1/2} \mathbf{p}_0 = \alpha_i \iff \mathbf{D}^{-1/2} \mathbf{p}_0 = \sum_{i=1}^n \alpha_i \psi_i \iff \mathbf{p}_0 = \sum_{i=1}^n \alpha_i \mathbf{D}^{1/2} \psi_i. \quad (6.1)$$

The above equivalences are best understood if you start from the middle. To get to the left side, you need to observe that adding  $\psi_i^\top$  to both sides cancels all terms  $\psi_j$  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

$$\mathbf{p}_t = \sum_{i=1}^n \alpha_i \left(1 - \frac{\nu_i}{2}\right)^t \mathbf{D}^{1/2} \boldsymbol{\psi}_i = \alpha_1 \mathbf{D}^{1/2} \boldsymbol{\psi}_1 + \sum_{i=2}^n \alpha_i \left(1 - \frac{\nu_i}{2}\right)^t \mathbf{D}^{1/2} \boldsymbol{\psi}_i. \quad (6.2)$$

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 6.1).

**Theorem 6.2.4.** *For any connected graph  $G$ , we have that the lazy random walk converges to the stationary distribution of  $G$ .*

### 6.2.3 The Rate of Convergence

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

**Theorem 6.2.5.** *For any  $\mathbf{p}_0$ , at any time step  $t$ , we have for  $\mathbf{p}_t = \tilde{\mathbf{W}}^t \mathbf{p}_0$  that*

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

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 6.2.6.** *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*

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

From Equation 6.2, we obtain that

$$\mathbf{p}_t(b) - \boldsymbol{\pi}(b) = \mathbf{1}_b^\top (\mathbf{p}_t - \boldsymbol{\pi}) = \mathbf{1}_b^\top \left( \sum_{i=2}^n \alpha_i \left(1 - \frac{\nu_i}{2}\right)^t \mathbf{D}^{1/2} \boldsymbol{\psi}_i \right) \quad (6.3)$$

$$= \sum_{i=2}^n \alpha_i \left(1 - \frac{\nu_i}{2}\right)^t \mathbf{1}_b^\top \mathbf{D}^{1/2} \boldsymbol{\psi}_i \leq \left(1 - \frac{\nu_2}{2}\right)^t \cdot \sum_{i=2}^n \alpha_i \mathbf{1}_b^\top \mathbf{D}^{1/2} \boldsymbol{\psi}_i \quad (6.4)$$

Taking the absolute value on both sides, we obtain that

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

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 6.1,  $\sum_{i=2}^n \alpha_i^2 = \sum_{i=2}^n \left(\boldsymbol{\psi}_i^\top \mathbf{D}^{-1/2} \mathbf{p}_0\right)^2 \leq \|\mathbf{D}^{-1/2} \mathbf{p}_0\|_2^2 = \|\mathbf{D}^{-1/2} \mathbf{1}_a\|_2^2 = 1/d_a$ .
- Finally, we show that  $\sum_{i=2}^n \left(\mathbf{1}_b^\top \mathbf{D}^{1/2} \boldsymbol{\psi}_i\right)^2 \leq \sum_{i=1}^n \left(\mathbf{1}_b^\top \mathbf{D}^{1/2} \boldsymbol{\psi}_i\right)^2 = \|\mathbf{D}^{1/2} \mathbf{1}_b\|_2^2 = d_b$  (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

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

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

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

## 6.3 Properties of Random Walks

We now shift our focus away from convergence of random walks and consider some interesting properties of random walks. 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  $\mathbf{L}\mathbf{x} = \mathbf{b}$  can solve a host of problems in random walks.

### 6.3.1 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  $H_{a,s} = \min\{t \mid \mathbf{v}_t = s\}$ . 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  $\mathbf{L}\mathbf{h} = \mathbf{b}$ . 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.  $\mathbf{h}_s = 0$ . For  $a \neq s$ , 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

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

Using that  $\mathbf{h}_a = \mathbf{1}_a^\top \mathbf{h} = \mathbf{1}_a^\top \mathbf{I}\mathbf{h}$ , 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

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

where we have an extra degree of freedom in formulating a constraint  $1 - \alpha = \mathbf{1}_s^\top (\mathbf{I} - \mathbf{W}^\top) \mathbf{h}$  (a degree of freedom that we indeed need in a second).

**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  $\mathbf{h}_t = 0$ , we can reconstruct  $\mathbf{h}$  from  $\mathbf{x}$  straight-forwardly by subtracting  $\mathbf{x}_t\mathbf{1}$ .

**Theorem 6.3.1.** *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  $\mathbf{h} = \mathbf{x} - \mathbf{x}_t\mathbf{1}$ . 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 solve 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  $s$ :  $\mathbb{E}[H_{a,s}] = \mathbf{h}_a = \mathbf{x}_a - \mathbf{x}_s$ .

### 6.3.2 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  $C_{a,b} = H_{a,b} + H_{b,a}$ .

**More of the Same?** While superficially, commute and hitting times might seem similar, in many ways, commute times are more natural to consider than hitting times. Key properties of commute times that do not hold for hitting times:

- they are symmetric, and
- one can define a metric using the commute times, and
- arguing about additional edge insertions into the graph  $G$  is known to monotonically decrease  $C_{a,b}$  (for all  $a, b \in V$ ). For hitting times, it can be truly hard to argue about such insertions.

**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  $\mathbf{L}\mathbf{x} = \mathbf{b}_b$  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  $\mathbf{L}\mathbf{y} = \mathbf{b}_a$  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  $z$ . Note that, if we revert the flow by negating  $\mathbf{y}$ , then we can still compute the hitting time by taking  $\mathbb{E}[H_{z,a}] = -(-\mathbf{y}_z - (-\mathbf{y}_a)) = -(\mathbf{y}_a - \mathbf{y}_z)$ .

Thus, inducing voltages  $\mathbf{x} - \mathbf{y}$  on the graph  $G$ , we now have by linearity that  $\mathbb{E}[C_{a,b}] = \mathbb{E}[H_{a,b} + H_{b,a}] = |\mathbf{x}_a - \mathbf{x}_b| + |\mathbf{y}_a - \mathbf{y}_b|$ . But these voltages are also induced by  $\mathbf{L}(\mathbf{x} - \mathbf{y}) = \mathbf{b}_b - \mathbf{b}_a = \mathbf{d}^\top \mathbf{1}(\mathbf{1}_a - \mathbf{1}_b)$  (again by linearity). That is the flow that routes  $2m$  units of flow from  $b$  to  $a$ .

**Theorem 6.3.2.** *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  $\mathbf{L}\mathbf{z} = 2m \cdot (\mathbf{1}_b - \mathbf{1}_a)$ .*

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  $\mathbb{E}[C_{a,b}] = 2m \cdot R_{\text{eff}}(a, b)$ . This is a classic result.

# Chapter 7

## Pseudo-inverses and Effective Resistance

### 7.1 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)  $(\mathbf{L}^+)^\top = \mathbf{L}^+$  (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,

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

where  $\mathbf{V} = [\mathbf{v}_1 \ \cdots \ \mathbf{v}_n]$  and  $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_n\}$ . 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.

## 7.2 Electrical Flows Again

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

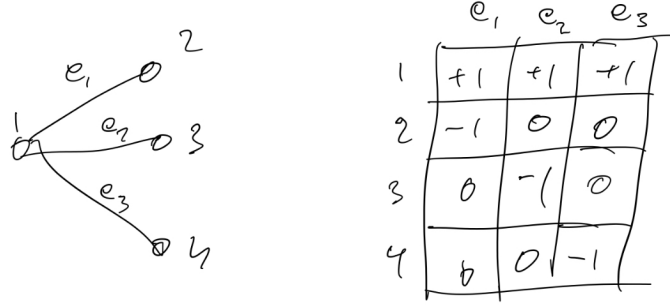


Figure 7.1: An example of a graph and its incidence matrix  $\mathbf{B}$ .

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  $\mathbf{r}(e)$  = 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 7.2.1.**

$$\tilde{\mathbf{f}} = \arg \min_{\mathbf{B}\mathbf{f}=\mathbf{d}} \mathbf{f}^\top \mathbf{R}\mathbf{f} \text{ where } \mathbf{f}^\top \mathbf{R}\mathbf{f} = \sum_e \mathbf{r}(e)\mathbf{f}(e)^2, \quad (7.1)$$

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

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

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  $\mathbf{x} \in \mathbb{R}^V$ ,

$$\frac{1}{2}\mathbf{f}^\top \mathbf{R}\mathbf{f} \geq \mathbf{d}^\top \mathbf{x} - \frac{1}{2}\mathbf{x}^\top \mathbf{L}\mathbf{x}. \quad (7.2)$$



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,

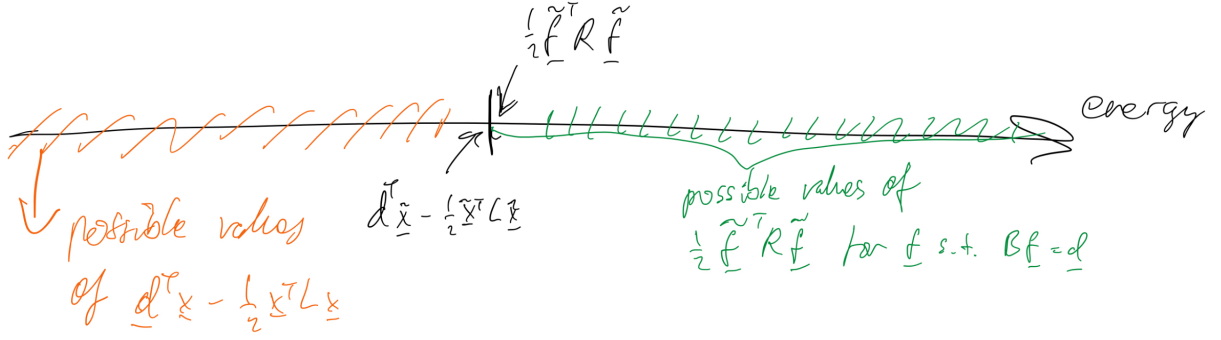
$$\frac{1}{2} \tilde{\mathbf{f}}^\top \mathbf{R} \tilde{\mathbf{f}} = \mathbf{d}^\top \tilde{\mathbf{x}} - \frac{1}{2} \tilde{\mathbf{x}}^\top \mathbf{L} \tilde{\mathbf{x}}. \quad (7.3)$$

By combining Equation (7.2) and Equation (7.3), 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 (7.1) holds.

The drawing below shows how the quantities line up:



□

## 7.3 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

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

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 7.3.1.** 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  $\tilde{\mathbf{f}} = \mathbf{R}^{-1} \mathbf{B}^\top \tilde{\mathbf{x}}$ , 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 7.3.2.** 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  $\mathbf{L}^{+1/2}$  be the square root of  $\mathbf{L}^+$ . So

$$R_{\text{eff}}(a, b) = (\mathbf{e}_b - \mathbf{e}_a)^\top \mathbf{L}^+ (\mathbf{e}_b - \mathbf{e}_a) = \|\mathbf{L}^{+1/2} (\mathbf{e}_b - \mathbf{e}_a)\|^2.$$

**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 7.2: 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  $1 = \mathbf{f}(i) = \frac{\Delta(i)}{\mathbf{r}(i)}$ , 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  $\mathbf{r}(1), \mathbf{r}(2), \dots, \mathbf{r}(k)$ .

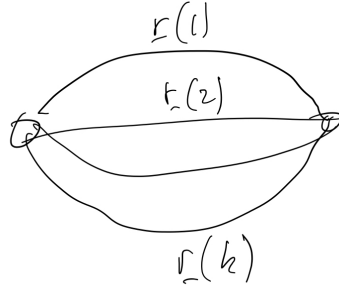


Figure 7.3: A graph on just two vertices with  $k$  parallel multiedges.

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

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

### 7.3.1 Effective Resistance is a Distance

**Definition 7.3.3.** 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 7.3.4.**  $R_{\text{eff}}$  is a distance.

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

**Claim 7.3.5.** 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.

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

Rearranging  $\tilde{\mathbf{x}}(c) = \frac{\sum_{(u,c)} w(u, c) \tilde{\mathbf{x}}(u)}{\sum_{(u,c)} w(u, c)}$ . 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 7.3.3 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 7.3.5 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 8

## Different Perspectives on Gaussian Elimination

### 8.1 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:

$$\mathbf{L} = \begin{pmatrix} W & -\mathbf{a}^\top \\ -\mathbf{a} & \text{diag}(\mathbf{a}) + \mathbf{L}_{-1} \end{pmatrix}. \quad (8.1)$$

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

$$\begin{pmatrix} W & -\mathbf{a}^\top \\ -\mathbf{a} & \text{diag}(\mathbf{a}) \end{pmatrix} \quad (8.2)$$

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

$$y = \frac{1}{W} (b + \mathbf{a}^\top \mathbf{z}). \quad (8.3)$$

Observe that

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

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 8.1.1.**

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 8.1.1 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 (8.3) to get the value for  $y$ .

**Remark 8.1.2.** 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.

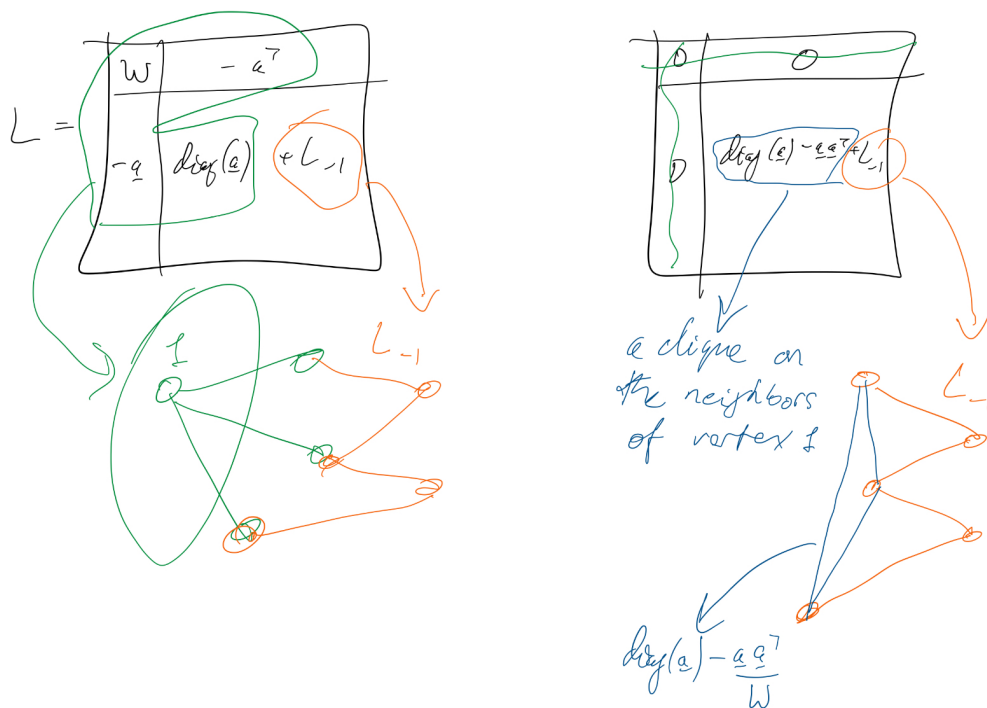
*Proof of Claim 8.1.1.* To establish the first part, we note that  $\mathbf{1}^\top \mathbf{d}' = \mathbf{1}^\top \mathbf{c} + b \frac{1}{W} \mathbf{a}^\top \mathbf{1} = \mathbf{1}^\top \mathbf{c} + b = \mathbf{1}^\top \mathbf{d} = 0$ . To establish the second part, we notice that  $\mathbf{L}_{-1}$  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 8.1.3.** 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 8.1.3 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  $\mathbf{S}$ .  $\square$



## 8.2 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  $\mathbf{L} = \mathbf{L}\mathbf{L}^\top$ , where  $\mathbf{L} \in \mathbb{R}^{n \times n}$  is a lower triangular matrix, i.e.  $\mathbf{L}(i, j) = 0$  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 8.2.1.** *Given an invertible square lower triangular matrix  $\mathbf{L}$ , we can solve the linear equation  $\mathbf{L}\mathbf{y} = \mathbf{b}$  in time  $O(\text{nnz}(\mathbf{L}))$ . Similarly, given an upper triangular matrix  $\mathbf{U}$ , we can solve linear equations  $\mathbf{U}\mathbf{z} = \mathbf{c}$  in time  $O(\text{nnz}(\mathbf{U}))$ .*

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 8.2.2.** Strictly speaking, the lemma requires us to have access an adjacency list representation of  $\mathbf{L}$  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  $\mathbf{M} = \mathbf{L}\mathbf{L}^\top$ , we can now solve linear equations in  $\mathbf{M}$  in time  $O(\text{nnz}(\mathbf{L}))$ .

**Remark 8.2.3.** 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.  $\mathbf{P}\mathbf{L}\mathbf{P}^\top$  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  $\mathbf{L} = \mathbf{L}\mathbf{L}^\top$ , where  $\mathbf{L}$  is not invertible? For graph Laplacians we have a simple characterization 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 8.2.4.** *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  $\mathbf{L} = \mathbf{L}\mathbf{L}^\top$  that we produce will have the property that all diagonal entries of  $\mathbf{L}$  are strictly non-zero, except that  $\mathbf{L}(n, n) = 0$ . From let us  $\hat{\mathbf{L}}$  as the matrix whose entries agree with  $\mathbf{L}$ , except that  $\hat{\mathbf{L}}(n, n) = 1$ . Let  $\mathbf{D}$  be the diagonal matrix with  $\mathbf{D}(i, i) = 1$  for  $i < n$  and  $\mathbf{D}(n, n) = 0$ . Then  $\mathbf{L}\mathbf{L}^\top = \hat{\mathbf{L}}\mathbf{D}\hat{\mathbf{L}}^\top$ , and  $\hat{\mathbf{L}}$  is invertible, and  $\mathbf{D}^+ = \mathbf{D}$ . 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  $\mathbf{L}^+ = \Pi_{\mathbf{L}}(\hat{\mathbf{L}}^\top)^{-1}\mathbf{D}\hat{\mathbf{L}}^{-1}\Pi_{\mathbf{L}}$ , and this matrix can be applied in time  $O(\text{nnz}(\mathbf{L}))$ .

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

**Theorem 8.2.5** (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

$$\mathbf{L}^{(i-1)} = \mathbf{0}_{n \times n}. \quad (8.4)$$

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 8.2.6.** *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  $i < n - 1$ , 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 (8.4) holds.  $\square$

*Sketch of proof of Claim 8.2.6.* We will focus on the first elimination, as the remaining are similar. Adopting the same notation as in Equation (8.1), 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 8.1.1 that  $\mathbf{L}^{(1)}(\{2, \dots, n\}, \{2, \dots, n\})$  is the Laplacian of a connected graph. This proves Claim 8.2.6 for the case  $i = 1$ . An induction following the same pattern can be used to prove the claim for all  $i < n$ .  $\square$

# Chapter 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  $\epsilon > 0$*

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

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  $\mathbf{X} \in \mathbb{R}$  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.*

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

We can rearrange this to get the desired statement. □

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

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

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

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

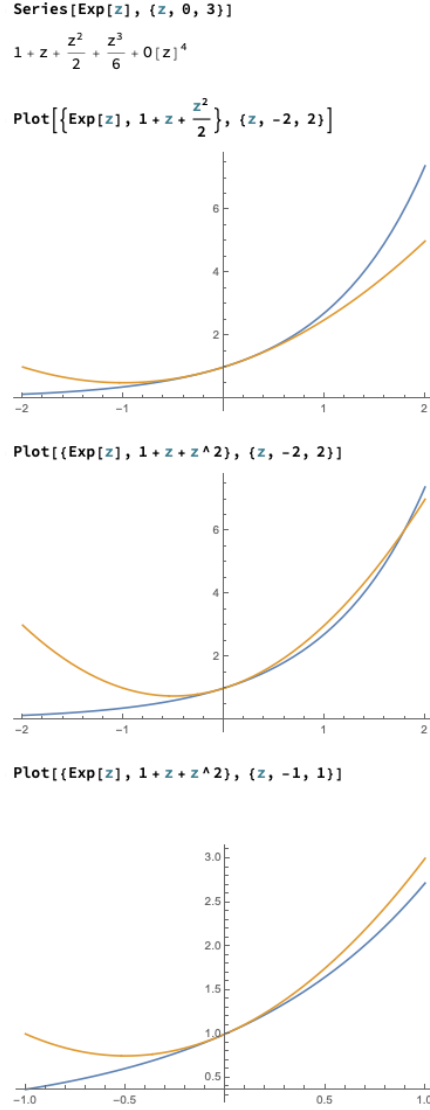


Figure 9.1: Plotting  $\exp(z)$  compared to  $1 + z + z^2$ .

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  $\epsilon > 0$*

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

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 9.2.1.** *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,  $\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$  where  $\mathbf{\Lambda} = \text{diag}_i\{\lambda_i(\mathbf{A})\}$  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.*

$$\text{Tr}(\mathbf{A}) = \text{Tr}(\mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top) = \text{Tr}\left(\mathbf{\Lambda} \underbrace{\mathbf{V}^\top \mathbf{V}}_{\mathbf{I}}\right) = \text{Tr}(\mathbf{\Lambda}) = \sum_i \lambda_i.$$

□

## 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  $\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$ , let

$$f(\mathbf{A}) = \mathbf{V} \text{diag}_i\{f(\lambda_i)\} \mathbf{V}^\top.$$

**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  $f(\mathbf{A}) = \exp(\mathbf{A}) = \mathbf{V} \exp(\mathbf{\Lambda}) \mathbf{V}^\top$  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 10.4.11.* 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.*

$$\begin{aligned} \mathbf{I} + \mathbf{A} + \mathbf{A}^2 - \exp(\mathbf{A}) &= \mathbf{V} \mathbf{I} \mathbf{V}^\top + \mathbf{V} \mathbf{A} \mathbf{V}^\top + \mathbf{V} \mathbf{A}^2 \mathbf{V}^\top - \mathbf{V} \exp(\mathbf{A}) \mathbf{V}^\top \\ &= \mathbf{V} (\mathbf{I} + \mathbf{A} + \mathbf{A}^2 - \exp(\mathbf{A})) \mathbf{V}^\top \\ &= \mathbf{V} \text{diag}\{1 + \lambda_i + \lambda_i^2 - \exp(\lambda_i)\} \mathbf{V}^\top \end{aligned}$$

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.*

$$\begin{aligned} \mathbf{A} - \log(\mathbf{I} + \mathbf{A}) &= \mathbf{V} \mathbf{A} \mathbf{V}^\top - \mathbf{V} \log(\mathbf{A} + \mathbf{I}) \mathbf{V}^\top \\ &= \mathbf{V} (\mathbf{A} - \log(\mathbf{A} + \mathbf{I})) \mathbf{V}^\top \\ &= \mathbf{V} \text{diag}\{\lambda_i - \log(1 + \lambda_i)\} \mathbf{V}^\top \end{aligned}$$

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 9.2.1, 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 9.2.1.* 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  $0 < \theta \leq 1/R$ ,

$$\begin{aligned}
\mathbb{E} \operatorname{Tr}(\exp(\theta \mathbf{X})) &= \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_{k-1}} \mathbb{E}_{\mathbf{X}_k} \operatorname{Tr} \exp \left( \underbrace{\theta \mathbf{X}_{<k}}_{\mathbf{H}} + \underbrace{\theta \mathbf{X}_k}_{=\log \exp(\theta \mathbf{X}_k)} \right), \quad \{\mathbf{X}_i\} \text{ are independent} \\
&\leq \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_{k-1}} \operatorname{Tr} \exp \left( \theta \mathbf{X}_{<k} + \log \mathbb{E} \exp(\theta \mathbf{X}_k) \right), \quad \text{by 10.4.9 and 9.2.16} \\
&\leq \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_{k-1}} \operatorname{Tr} \exp \left( \theta \mathbf{X}_{<k} + \log \mathbb{E} [\mathbf{I} + \theta \mathbf{X}_k + \theta^2 \mathbf{X}_k^2] \right), \quad \text{by 10.4.14, 9.2.7, and 10.4.1} \\
&\leq \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_{k-1}} \operatorname{Tr} \exp \left( \theta \mathbf{X}_{<k} + \theta^2 \mathbb{E} \mathbf{X}_k^2 \right), \quad \text{by 10.4.12 and 9.2.7} \\
&= \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_{k-2}} \mathbb{E}_{\mathbf{X}_{k-1}} \operatorname{Tr} \exp \left( \underbrace{\theta^2 \mathbb{E} \mathbf{X}_k^2 + \theta \mathbf{X}_{<k-1}}_{\mathbf{H}} + \theta \mathbf{X}_{k-1} \right), \\
&\quad \vdots \\
&\leq \operatorname{Tr} \exp \left( \theta^2 \sum_i \mathbb{E} [\mathbf{X}_i^2] \right), \\
&\leq \operatorname{Tr} \exp (\theta^2 \sigma^2 \mathbf{I}), \quad \text{by 9.2.7 and } \sum_i \mathbb{E} [\mathbf{X}_i^2] \preceq \sigma^2 \mathbf{I} \\
&= n \cdot \exp(\theta^2 \sigma^2).
\end{aligned}$$

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

$$\mathbf{A} \approx_\epsilon \mathbf{B} \text{ if and only if } \frac{1}{1+\epsilon} \mathbf{A} \leq \mathbf{B} \leq (1+\epsilon) \mathbf{A}.$$

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  $\tilde{G}$ .

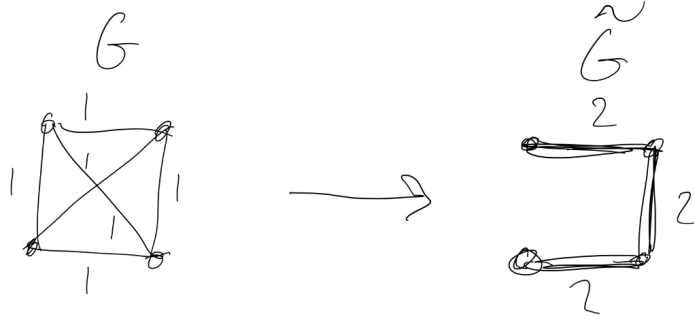


Figure 9.2: A graph  $G$  and a spectral sparsifier  $\tilde{G}$ , satisfying  $\mathbf{L}_G \approx_\epsilon \mathbf{L}_{\tilde{G}}$  for  $\epsilon = 2.42$ .

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.  $\square$

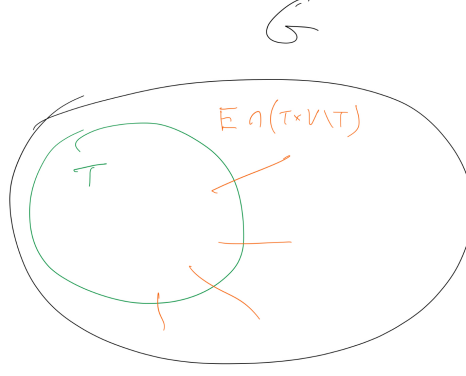


Figure 9.3: The cut  $c_G(T)$  in  $G$ .

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  $\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top = \sum_i \lambda_i \mathbf{v}_i \mathbf{v}_i^\top$  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  $\mathbf{\Pi}_L = \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$ .

Now that we understand the projection  $\mathbf{\Pi}_L$ , it is not hard to show the following claim.

**Claim 9.3.8.**

1.  $\mathbf{\Pi}_L \approx_\epsilon \mathbf{L}^{+1/2} \tilde{\mathbf{L}} \mathbf{L}^{+1/2}$  implies  $\mathbf{L} \approx_\epsilon \tilde{\mathbf{L}}$ .
2. For  $\epsilon \leq 1$ , we have that  $\left\| \mathbf{\Pi}_L - \mathbf{L}^{+1/2} \tilde{\mathbf{L}} \mathbf{L}^{+1/2} \right\| \leq \epsilon/2$  implies  $\mathbf{\Pi}_L \approx_\epsilon \mathbf{L}^{+1/2} \tilde{\mathbf{L}} \mathbf{L}^{+1/2}$ .

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 2, we assume  $\mathbf{\Pi}_L \approx_\epsilon \mathbf{L}^{+1/2} \tilde{\mathbf{L}} \mathbf{L}^{+1/2}$ . 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  $\mathbf{v} = \mathbf{L}^{+1/2} \mathbf{z}$  for some  $\mathbf{z}$ , as  $\ker(\mathbf{L})^\perp = \text{im}(\mathbf{L}) = \text{im}(\mathbf{L}^{+1/2})$ . Hence

$$\mathbf{v}^\top \tilde{\mathbf{L}} \mathbf{v} = \mathbf{z}^\top \mathbf{L}^{+1/2} \tilde{\mathbf{L}} \mathbf{L}^{+1/2} \mathbf{z} \geq \frac{1}{1+\epsilon} \mathbf{z}^\top \mathbf{L}^{+1/2} \mathbf{L} \mathbf{L}^{+1/2} \mathbf{z} = \frac{1}{1+\epsilon} \mathbf{v}^\top \mathbf{L} \mathbf{v}$$

and similarly

$$\mathbf{v}^\top \tilde{\mathbf{L}} \mathbf{v} = \mathbf{z}^\top \mathbf{L}^{+1/2} \tilde{\mathbf{L}} \mathbf{L}^{+1/2} \mathbf{z} \leq (1+\epsilon) \mathbf{z}^\top \mathbf{L}^{+1/2} \mathbf{L} \mathbf{L}^{+1/2} \mathbf{z} = (1+\epsilon) \mathbf{v}^\top \mathbf{L} \mathbf{v}.$$

Thus we have established  $\mathbf{L} \approx_\epsilon \tilde{\mathbf{L}}$ .

To prove Part 2, we assume  $\left\| \mathbf{\Pi}_L - \mathbf{L}^{+1/2} \tilde{\mathbf{L}} \mathbf{L}^{+1/2} \right\| \leq \epsilon/2$ . This is equivalent to

$$-\frac{\epsilon}{2} \mathbf{I} \preceq \mathbf{L}^{+1/2} \tilde{\mathbf{L}} \mathbf{L}^{+1/2} - \mathbf{\Pi}_L \preceq \frac{\epsilon}{2} \mathbf{I}$$

But since

$$\mathbf{1}^\top (\mathbf{L}^{+1/2} \tilde{\mathbf{L}} \mathbf{L}^{+1/2} - \mathbf{\Pi}_L) \mathbf{1} = 0,$$

we can in fact sharpen this to

$$-\frac{\epsilon}{2} \mathbf{\Pi}_L \preceq \mathbf{L}^{+1/2} \tilde{\mathbf{L}} \mathbf{L}^{+1/2} - \mathbf{\Pi}_L \preceq \frac{\epsilon}{2} \mathbf{\Pi}_L.$$

Rearranging, we then conclude

$$(1 - \frac{\epsilon}{2}) \mathbf{\Pi}_L \preceq \mathbf{L}^{+1/2} \tilde{\mathbf{L}} \mathbf{L}^{+1/2} \preceq (1 + \frac{\epsilon}{2}) \mathbf{\Pi}_L.$$

Finally, we note that  $1/(1+\epsilon) \leq (1 - \frac{\epsilon}{2})$  to reach our conclusion,  $\mathbf{\Pi}_L \approx_\epsilon \mathbf{L}^{+1/2} \tilde{\mathbf{L}} \mathbf{L}^{+1/2}$ .  $\square$

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

$$\left\| \mathbf{\Pi}_L - \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2} \right\| \leq \epsilon/2. \quad (9.3)$$

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

$$\mathbf{Y}_e = \begin{cases} \frac{\mathbf{w}(e)}{p_e} \mathbf{b}_e \mathbf{b}_e^\top & \text{with probability } p_e \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

This way,  $\tilde{\mathbf{L}} = \sum_e \mathbf{Y}_e$ . Note that  $\mathbb{E}[\mathbf{Y}_e] = p_e \frac{\mathbf{w}(e)}{p_e} \mathbf{b}_e \mathbf{b}_e^\top = \mathbf{w}(e) \mathbf{b}_e \mathbf{b}_e^\top$ , and so

$$\mathbb{E}[\tilde{\mathbf{L}}] = \sum_e \mathbb{E}[\mathbf{Y}_e] = \mathbf{L}.$$

By linearity of  $\Phi$ ,

$$\mathbb{E}[\Phi(\tilde{\mathbf{L}})] = \Phi(\mathbb{E}[\tilde{\mathbf{L}}]) = \mathbf{\Pi}_L.$$

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

$$p_e = \min(\alpha \|\Phi(\mathbf{w}(e) \mathbf{b}_e \mathbf{b}_e^\top)\|, 1)$$

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) holds.

This means that we can apply Theorem 9.2.1 to our  $\mathbf{X} = \sum_e \mathbf{X}_e$ , with  $R = \frac{1}{\alpha}$  and  $\sigma^2 = \frac{1}{\alpha}$ , to get

$$\Pr \left[ \left\| \mathbf{\Pi}_L - \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2} \right\| \geq \epsilon/2 \right] \leq 2n \exp \left( \frac{-0.25\epsilon^2}{(\epsilon + 4)/\alpha} \right)$$

Since  $0 < \epsilon < 1$ , this means that if  $\alpha = 40\epsilon^{-2} \log(n/\delta)$ , then

$$\Pr \left[ \left\| \mathbf{\Pi}_L - \mathbf{L}^{+/2} \tilde{\mathbf{L}} \mathbf{L}^{+/2} \right\| \geq \epsilon/2 \right] \leq \frac{2n\delta^2}{n^2} \leq \delta/2.$$

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

$$\mathbb{E} [|\tilde{E}|] = \sum_e p_e \leq \alpha \sum_e \mathbf{w}(e) \left\| \mathbf{L}^{+/2} \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/2} \right\|$$

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

$$\begin{aligned} \sum_e \mathbf{w}(e) \left\| \mathbf{L}^{+/2} \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/2} \right\| &= \sum_e \mathbf{w}(e) \text{Tr} \left( \mathbf{L}^{+/2} \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/2} \right) \\ &= \text{Tr} \left( \mathbf{L}^{+/2} \left( \sum_e \mathbf{w}(e) \mathbf{b}_e \mathbf{b}_e^\top \right) \mathbf{L}^{+/2} \right) \\ &= \text{Tr}(\mathbf{\Pi}_L) = n - 1. \end{aligned}$$

Thus with our choice of  $\alpha$ ,

$$\mathbb{E} [|\tilde{E}|] \leq 40\epsilon^{-2} \log(n/\delta)n.$$

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

$$\|\Phi(\mathbf{w}(e) \mathbf{b}_e \mathbf{b}_e^\top)\| = \mathbf{w}(e) \left\| \mathbf{L}^{+/2} \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/2} \right\| \leq \mathbf{w}(e) \left\| \mathbf{L}^{+/2} \mathbf{b}_e \right\|_2^2.$$

Recall that in Chapter 7, we saw that the effective between vertex  $v$  and vertex  $u$  is given by  $\left\| \mathbf{L}^{+/2}(\mathbf{e}_u - \mathbf{e}_v) \right\|_2^2$ , 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  $\mathbf{L}\mathbf{L}^\top \approx_K \mathbf{M}$ , where  $\mathbf{L}$  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  $\mathbf{M} \approx_K \mathbf{L}\mathbf{L}^\top$ , 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  $\mathbf{L}^{-1}\mathbf{z}$  or  $(\mathbf{L}^\top)^{-1}\mathbf{z}$  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  $\mathbf{M} \approx_K \mathbf{L}\mathbf{D}\mathbf{L}^\top$ , where  $\mathbf{L}$  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  $\mathbf{L}^{-1}\mathbf{z}$  and  $(\mathbf{L}^\top)^{-1}\mathbf{z}$  and  $\mathbf{D}^+\mathbf{z}$  given a vector  $\mathbf{z}$ .

- $T_{proj}$  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.  $\mathcal{L}\mathcal{L}^\top = \mathbf{L}$

Let  $\mathbf{S}_0 = \mathbf{L}$ ;

**for**  $i = 1$  **to**  $i = n - 1$  **do**

$$\left[ \begin{array}{l} \mathbf{l}_i = \frac{1}{\sqrt{\mathbf{S}_{i-1}(i, i)}} \mathbf{S}_{i-1}(:, i); \\ \mathbf{S}_i = \mathbf{S}_{i-1} - \mathbf{l}_i \mathbf{l}_i^\top. \end{array} \right.$$

$\mathbf{l}_n = \mathbf{0}_{n \times 1}$ ;

**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.

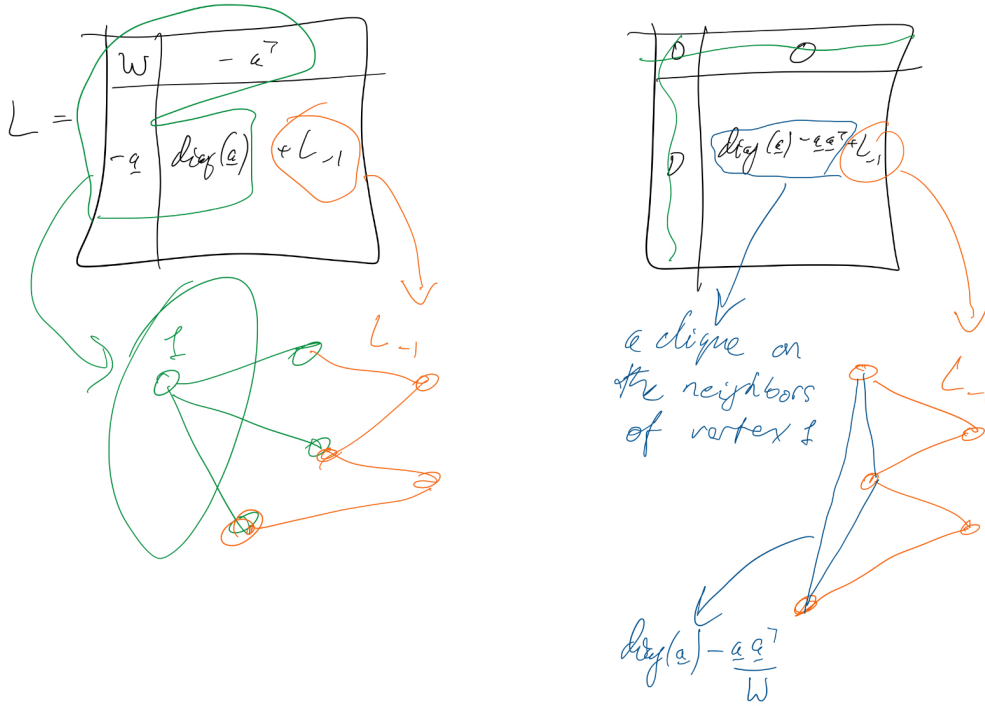


Figure 10.1: Gaussian Elimination:  $\text{CLIQUE}(1, \mathbf{L}) = \text{STAR}(1, \mathbf{L}) - \frac{1}{\mathbf{L}(1,1)} \mathbf{L}(:, 1) \mathbf{L}(:, 1)^\top$ .

In Chapter 8, we proved that  $\text{CLIQUE}(v, \mathbf{S})$  is a graph Laplacian – it follows from the proof of Claim 8.1.1 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  $\text{CLIQUE}(v, \mathbf{S})$

$\mathbf{Y}_v \leftarrow \mathbf{0}_{n \times n}$ ;

**foreach** *Multiedge*  $e = (v, i)$  *from*  $v$  *to* *a neighbor*  $i$  **do**

Randomly pick a neighbor  $j$  of  $v$  with probability  $\frac{\mathbf{w}(j,v)}{\mathbf{w}_v}$ ;

If  $i \neq j$ , let  $\mathbf{Y}_v \leftarrow \mathbf{Y}_v + \frac{\mathbf{w}(i,v)\mathbf{w}(j,v)}{\mathbf{w}(i,v) + \mathbf{w}(j,v)} \mathbf{b}_{i,j} \mathbf{b}_{i,j}^\top$ ;

**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  $O(\deg_{\mathbf{S}}(v))$  time to initialize a data structure used for sampling. Overall, this means the total time for  $O(\deg_{\mathbf{S}}(v))$  samples is still  $O(\deg_{\mathbf{S}}(v))$ .

**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  $i \neq j$ .

$$\begin{aligned} \mathbf{C}(i, j) &= -\frac{\mathbf{w}(i, v)\mathbf{w}(j, v)}{\mathbf{w}_v}, \\ \mathbb{E}_{\mathbf{Y}_v(i,j)} &= -\frac{\mathbf{w}(i, v)\mathbf{w}(j, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} \left( \frac{\mathbf{w}(j, v)}{\mathbf{w}_v} + \frac{\mathbf{w}(i, v)}{\mathbf{w}_v} \right) = -\frac{\mathbf{w}(i, v)\mathbf{w}(j, v)}{\mathbf{w}_v} = \mathbf{C}(i, j). \end{aligned}$$

□

**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

Let  $\mathbf{S}_0 = \mathbf{L}$ ;

Generate a random permutation  $\pi$  on  $[n]$ ;

**for**  $i = 1$  **to**  $i = n - 1$  **do**

$\mathbf{l}_i = \frac{1}{\sqrt{\mathbf{S}_{i-1}(\pi(i), \pi(i))}} \mathbf{S}_{i-1}(:, \pi(i));$   
 $\mathbf{S}_i = \mathbf{S}_{i-1} - \text{STAR}(\pi(i), \mathbf{S}_{i-1}) + \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1})$

$\mathbf{l}_n = \mathbf{0}_{n \times 1}$ ;

**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

$$\left\| \mathbf{L}^{+/2} (\mathcal{L}\mathcal{L}^\top - \mathbf{L}) \mathbf{L}^{+/2} \right\| \leq 0.5 \quad (10.2)$$

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  $\mathbf{L}^{+/2}$ , we reintroduce the “normalizing map”  $\Phi : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$  defined by

$$\Phi(\mathbf{A}) = \mathbf{L}^{+/2} \mathbf{A} \mathbf{L}^{+/2}.$$

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

$$\mathbb{E}[Z_i \mid Z_0, \dots, Z_{i-1}] = Z_{i-1}. \quad (10.3)$$

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

$$\mathbb{E}[X_i \mid Z_0, \dots, Z_{i-1}] = 0.$$

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

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

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 martingale 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

$$\mathbb{E}[\mathbf{Z}_i \mid \mathbf{Z}_0, \dots, \mathbf{Z}_{i-1}] = \mathbf{Z}_{i-1}. \quad (10.4)$$

**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*

$$\mathbb{E}[\mathbf{L}_i \mid \text{all random variables before } \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1})] = \mathbf{L}_{i-1}.$$

*Proof.* Let's only consider  $i = 1$  here as other cases are similar.

$$\mathbf{L}_0 = \mathbf{L} = \mathbf{l}_1 \mathbf{l}_1^\top + \text{CLIQUE}(v, \mathbf{L}) + \mathbf{L}_{-1}$$

$$\begin{aligned} \mathbf{L}_1 &= \mathbf{l}_1 \mathbf{l}_1^\top + \text{CLIQUESAMPLE}(v, \mathbf{L}) + \mathbf{L}_{-1} \\ \mathbb{E}[\mathbf{L}_1 \mid \pi(1)] &= \mathbf{l}_1 \mathbf{l}_1^\top + \mathbb{E}[\text{CLIQUESAMPLE}(v, \mathbf{L}) \mid \pi(1)] + \mathbf{L}_{-1} \\ &= \mathbf{l}_1 \mathbf{l}_1^\top + \text{CLIQUE}(v, \mathbf{L}) + \mathbf{L}_{-1} \\ &= \mathbf{L}_0 \end{aligned}$$

where we used Lemma 10.3.3 to get  $\mathbb{E}[\text{CLIQUESAMPLE}(v, \mathbf{L}) \mid \pi(1)] = \text{CLIQUE}(v, \mathbf{L})$ .  $\square$

**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

$$\begin{aligned} \mathbf{L}_i - \mathbf{L}_{i-1} &= \mathbf{l}_i \mathbf{l}_i^\top + \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) - \text{STAR}(\pi(i), \mathbf{S}_{i-1}) \\ &= \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) - \text{CLIQUE}(\pi(i), \mathbf{S}_{i-1}) \\ &= \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) - \mathbb{E}[\text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) \mid \text{preceding samples}] \\ &\quad \text{by Lemma 10.3.3.} \end{aligned} \tag{10.5}$$

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  $\mathbf{L}_{j^*} \not\preceq 1.5\mathbf{L}$ .

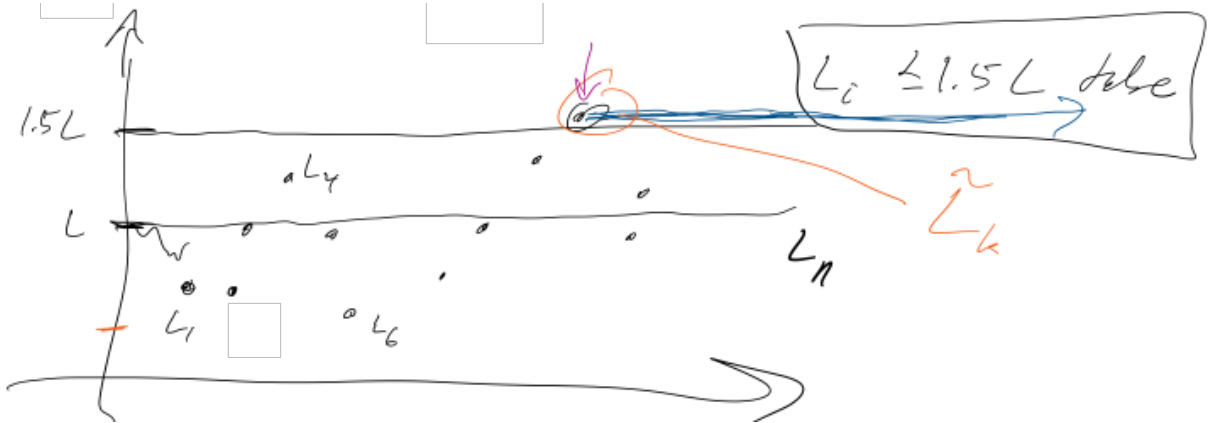


Figure 10.2: Gaussian Elimination :  $\text{CLIQUE}(1, \mathbf{L}) = \text{STAR}(1, \mathbf{L}) - \frac{1}{\mathbf{L}(1,1)} \mathbf{L}(:, 1) \mathbf{L}(:, 1)^\top$ .

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.  $\left\| \mathbf{L}^{+1/2}(\tilde{\mathbf{L}}_i - \mathbf{L}) \mathbf{L}^{+1/2} \right\| \leq 0.5$  implies  $\left\| \mathbf{L}^{+1/2}(\mathbf{L}_i - \mathbf{L}) \mathbf{L}^{+1/2} \right\| \leq 0.5$

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  $\|\mathbf{L}^{+/2}(\tilde{\mathbf{L}}_i - \mathbf{L})\mathbf{L}^{+/2}\| \leq 0.5$  is false, and the implication is vacuously true. If the, on the other hand, if 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 goin to be to show that  $\|\mathbf{L}^{+/2}(\tilde{\mathbf{L}}_i - \mathbf{L})\mathbf{L}^{+/2}\| \leq 0.5$  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

$$\tilde{\mathbf{X}}_i = \begin{cases} \mathbf{X}_i & \text{if for all } j < i \text{ we have } \mathbf{L}_j \preceq 1.5\mathbf{L} \\ \mathbf{0} & \text{otherwise} \end{cases} \quad (10.8)$$

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,*

$$\left\| \mathbf{L}^{+/2} \mathbf{w}_{\mathbf{S}}(e) \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/2} \right\| \leq R,$$

---

<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  $\text{CLIQUE\_SAMPLE}(v, \mathbf{S})$  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  $\mathbf{w} = \mathbf{w}_{\mathbf{S}}$  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 7. 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  $e \in \text{STAR}(v, \mathbf{S}_i)$ ,*

$$\left\| \Phi(\mathbf{w}_{\mathbf{S}_i}(e) \mathbf{b}_e \mathbf{b}_e^\top) \right\| \leq R.$$

*then all  $e \in \text{STAR}(\pi(i), \mathbf{S}_i)$ ,*

$$\left\| \Phi(\mathbf{Y}_{\pi(i), e}) \right\| \leq R.$$

**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

$$\left\| \mathbf{L}^{+/2} \mathbf{w}(\hat{e}) \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/2} \right\| \leq 1.$$

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*

$$\left\| \Phi(\mathbf{Y}_{\pi(i),e}) \right\| \leq 1/K.$$

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 Subsection, 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

$$\begin{aligned} \mathbb{E} \left[ \exp (\theta \tilde{\mathbf{X}}_{i,e}) \mid \text{preceding samples} \right] &\preceq \mathbf{I} + \mathbb{E} \left[ (\theta \tilde{\mathbf{X}}_{i,e})^2 \mid \text{preceding samples} \right] \\ &\preceq \mathbf{I} + \frac{1}{K} \theta^2 \cdot \mathbb{E} \left[ \Phi (\mathbf{Y}_{\pi(i),e}) \mid \text{preceding samples} \right] \end{aligned} \quad (10.11)$$



### 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  $0 < \theta \leq 0.5/\sqrt{K}$

$$\begin{aligned}
& \mathbb{E} \operatorname{Tr} \left( \exp \left( \theta \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i \right) \right) \tag{10.12} \\
&= \mathbb{E}_{(<\hat{n})} \mathbb{E}_{\pi(\hat{n})} \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},1}} \cdots \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}-1}} \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}}} \operatorname{Tr} \exp \left( \underbrace{\sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \sum_{e=1}^{d_{\hat{n}}-1} \theta \tilde{\mathbf{X}}_{\hat{n},e}}_{\mathbf{H}} + \theta \tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}} \right) \\
&\quad \tilde{\mathbf{X}}_{\hat{n},1}, \dots, \tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}} \text{ are independent conditional on } (<\hat{n}), \pi(\hat{n}) \\
&\leq \mathbb{E}_{(<\hat{n})} \mathbb{E}_{\pi(\hat{n})} \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},1}} \cdots \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}-1}} \operatorname{Tr} \exp \left( \sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \sum_{e=1}^{d_{\hat{n}}-1} \theta \tilde{\mathbf{X}}_{\hat{n},e} + \frac{1}{K} \theta^2 \cdot \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}}} \Phi(\mathbf{Y}_{\pi(\hat{n}),d_{\hat{n}}}) \right) \\
&\quad \text{By Equation (10.11) and Corollary 10.4.13 .} \\
&\vdots \quad \text{Repeat for each multi-edge sample } \tilde{\mathbf{X}}_{\hat{n},1}, \dots, \tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}-1} \\
&\leq \mathbb{E}_{(<\hat{n})} \mathbb{E}_{\pi(\hat{n})} \operatorname{Tr} \exp \left( \sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \sum_{e=1}^{d_{\hat{n}}} \frac{1}{K} \theta^2 \cdot \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},e}} \Phi(\mathbf{Y}_{\pi(\hat{n}),e}) \right) \\
&= \mathbb{E}_{(<\hat{n})} \mathbb{E}_{\pi(\hat{n})} \operatorname{Tr} \exp \left( \sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \frac{1}{K} \theta^2 \operatorname{CLIQUE}(\pi(\hat{n}), \mathbf{S}_{\hat{n}-1}) \right)
\end{aligned}$$

To further bound 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 its surprising strength.

We can treat  $\frac{1}{K} \theta^2 \operatorname{CLIQUE}(\pi(\hat{n}), \mathbf{S}_{\hat{n}-1})$  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  $\|\mathbf{A}\| \leq R$ .

*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  $\|\frac{1}{K}\theta^2\Phi(\text{CLIQUE}(\pi(i), \mathbf{S}_{i-1}))\| \leq 1$ , 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

$$\begin{aligned} &\mathbb{E} \text{Tr} \left( \exp\left(\theta \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i\right) \right) \\ &\leq \mathbb{E}_{(<\hat{n})} \mathbb{E}_{\pi(\hat{n})} \text{Tr} \exp \left( \sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \frac{1}{K} \theta^2 \text{CLIQUE}(\pi(\hat{n}), \mathbf{S}_{\hat{n}-1}) \right) \\ &\leq \mathbb{E}_{(<\hat{n})} \text{Tr} \exp \left( \sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \frac{6\theta^2}{K(n+1-i)} \mathbf{I} \right) \end{aligned}$$

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):

$$0.5\mathbf{L} \preceq \mathbf{L}\mathbf{L}^\top \preceq 1.5\mathbf{L}.$$

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  $\mathbf{L}$ ). 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.

# Bibliography

- [DO19] Jelena Diakonikolas and Lorenzo Orecchia. Conjugate gradients and accelerated methods unified: The approximate duality gap view. *arXiv preprint arXiv:1907.00289*, 2019.
- [HS<sup>+</sup>52] Magnus R Hestenes, Eduard Stiefel, et al. Methods of conjugate gradients for solving linear systems. *Journal of research of the National Bureau of Standards*, 49(6):409–436, 1952.
- [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.
- [Lan52] Cornelius Lanczos. Solution of systems of linear equations by minimized iterations. *J. Res. Nat. Bur. Standards*, 49(1):33–53, 1952.
- [Nes83] Y. E. Nesterov. A method for solving the convex programming problem with convergence rate  $o(1/k^2)$ . *Dokl. Akad. Nauk SSSR*, 269:543–547, 1983.
- [NY83] A Nemirovski and D Yudin. Information-based complexity of mathematical programming. *Izvestia AN SSSR, Ser. Tekhnicheskaya Kibernetika (the journal is translated to English as Engineering Cybernetics. Soviet J. Computer & Systems Sci.)*, 1, 1983.
- [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.