

STAT 9912: Acceleration, Greed and Hedging in Optimization

Taught by Dr. Jason Altschuler
Notes written by Faraz Rahman

Contents

1 Setting the stage	2
2 Foundational Questions about GD	2
2.1 Why move in the direction $-\nabla f$	2
3 Quadratics	3
3.1 Do we converge?	3
3.2 How fast do we converge?	4
3.3 Optimal Schedules	6
3.3.1 2 Step Schedules	6
3.3.2 The General Case	7
3.4 The Chebyshev Polynomials and Accelerated Methods	8
3.5 Random Stepsizes	11
3.6 Solving the Extremal Problem via Potential Theory:	11
3.7 Understanding the best case scenarios via a two player game	12

1 Setting the stage

Mathematical optimization tackles the set of problems that can be formulated as

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) \\ & \text{subject to} \quad \mathbf{x} \in \mathcal{X} \end{aligned}$$

This course aims to study progress in first order optimizers: algorithms that aim to solve the above problem given an oracle that can compute the following given \mathbf{x}

1. $f(\mathbf{x})$
2. $\nabla f(\mathbf{x})$

The canonical algorithm, gradient descent (GD), is given by the following iterative update rule

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \alpha_t \nabla f(\mathbf{x}_t)$$

where $\{\alpha_t\}$ parameterizes different algorithms of GD.

The course will study gradient descent in various problem settings, along with how choices of $\{\alpha_t\}$ can affect convergence rates, numerical robustness, and // TODO: ADD DESC.

2 Foundational Questions about GD

- Q1) Why move in the direction $-\nabla f$
- Q2) Can GD converge? (with any step size)
- Q3) How fast? (optimal step size)

2.1 Why move in the direction $-\nabla f$

The gradient descent algorithm, as many other tools in science and engineering, can be derived from solving a linear approximation to the optimization objective.

Writing the first-order Taylor-expansion at \mathbf{x} we get

$$f(\mathbf{x} + \mathbf{v}) \approx f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle$$

If we want to use this crude optimization to (try to) move \mathbf{x} in a direction that reduces f , we can setup the following problem.

$$\begin{aligned} & \underset{\mathbf{v}}{\text{minimize}} \quad f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle \\ & \text{subject to} \quad \|\mathbf{v}\|_2^2 \leq 1 \end{aligned}$$

In English, this asks: “If we were only allowed to move one unit away from \mathbf{x} what direction should we move to decrease $f(\mathbf{x})$ the most?”.

We will find that

$$\mathbf{v}^* \propto -\nabla f(\mathbf{x})$$

// TODO: ADD SHORT DERIVATION HERE

A comment should be made that the choice of the ℓ_2 norm here is not obvious. Different norms will induce different update rules creating a general class of algorithms called methods of steepest descent.

3 Quadratics

The simplest place to start our study is in the optimization of convex quadratic functions.

We can consider any function that can be expressed as

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{H}\mathbf{x} - \mathbf{b}^\top \mathbf{x} + c$$

where $\mathbf{H} \succeq 0$ (necessary for f to be convex) and \mathbf{H} is symmetric.¹

One assumption we will make is that $m\mathbf{I} \preceq \mathbf{H} \preceq M\mathbf{I}$ where $m, M > 0$ making the f M -smooth and m -strongly convex.

This problem has a closed-form optimal solution given by

$$\mathbf{x}^* = \mathbf{H}^{-1}\mathbf{b}$$

which can be derived via the FOC. This equality will be useful for deriving convergence rates for GD on quadratics.

3.1 Do we converge?

Much can be said about GD by studying the difference to between the current solution and the optimal solution

$$\mathbf{x}_t - \mathbf{x}^*$$

over time. Intuitively we want the distance to go down... and fast!

Plugging in the update rule we can see how

$$(\mathbf{x}_{t+1} - \mathbf{x}^*) = (\mathbf{x}_t - \mathbf{x}^*) - \alpha_t \nabla f(\mathbf{x})$$

Note that

$$\begin{aligned} \nabla f(\mathbf{x}) &= \mathbf{H}\mathbf{x}_t - \mathbf{b} \\ &= \mathbf{H}\mathbf{x}_t - \mathbf{H}\mathbf{H}^{-1}\mathbf{b} \\ &= \mathbf{H}\mathbf{x}_t - \mathbf{H}\mathbf{x}^* \\ &= \mathbf{H}(\mathbf{x}_t - \mathbf{x}^*) \end{aligned}$$

¹We can assume \mathbf{H} is symmetric without loss of generality because if \mathbf{H} is not symmetric, we can replace it with $\hat{\mathbf{H}} = \frac{1}{2}(\mathbf{H} + \mathbf{H}^\top)$, which yields an equivalent function \hat{f}

where the third lines comes from the closed form solution $\mathbf{x}^* = \mathbf{H}^{-1}\mathbf{b}$.

If we plug in this term for $\nabla f(x)$ we get a recurrence

$$\begin{aligned} (\mathbf{x}_{t+1} - \mathbf{x}^*) &= (\mathbf{x}_t - \mathbf{x}^*) - \alpha_t \mathbf{H}(\mathbf{b}x_t - \mathbf{x}^*) \\ &= (\mathbf{I} - \alpha_t \mathbf{H})(\mathbf{x}_t - \mathbf{x}^*) \end{aligned}$$

This is a Linear Dynamical System (LDS). Recall that an LDS

$$\mathbf{z}_{t+1} = \mathbf{A}\mathbf{z}_t$$

will converge so long as $|\lambda_i| < 1$ for all eigenvalues λ_i of \mathbf{A} . // TODO: CHECK THIS + ADD A REFERENCE

Recall that we assumed that $m\mathbf{I} \preceq \mathbf{H} \preceq M\mathbf{I}$, so our dynamics matrix $\mathbf{A}_t = (\mathbf{I} - \alpha_t \mathbf{H})$ we know that the eigenvalues for this dynamical system sit in

$$\lambda_i \in [1 - \alpha M, 1 - \alpha m]$$

Hence, if set all steps to an adequately chosen $\alpha_t = \alpha$ the above LDS will converge to 0.

3.2 How fast do we converge?

Before we can answer this question we need to define what it means to converge “quickly”. Consider the metric

$$R_T = \frac{\|\mathbf{x}_T - \mathbf{x}^*\|}{\|\mathbf{x}_0 - \mathbf{x}^*\|}$$

which captures the fraction of the original distance remaining in our optimization. Naturally we want this as low as possible. Taking the geometric mean $(R_T)^{1/T}$ can tell us about the average contraction rate of our error per optimization step.

In our analysis we will focus on worst-case analysis over the functions f , and random inits x_0 , so if $\alpha_t = \alpha$ is constant, then it suffices to consider R_1 . We will start here

Consider the following minimax problem where we look for the algorithm (parametrized by constant step-size α) that minimizes the the worst case 1-step contraction rate R_1

$$\operatorname{argmin}_{\alpha} \max_{f, \mathbf{x}_0} R_1 := \frac{\|\mathbf{x}_1 - \mathbf{x}^*\|}{\|\mathbf{x}_0 - \mathbf{x}^*\|}$$

Using the LDS update to expand the numerator,

$$\operatorname{argmin}_{\alpha} \max_{\mathbf{H}, \mathbf{x}_0} \frac{\|(\mathbf{I} - \alpha \mathbf{H})(\mathbf{x}_0 - \mathbf{x}^*)\|}{\|\mathbf{x}_0 - \mathbf{x}^*\|}$$

Some may recognize the inner maximization as the matrix norm of $(\mathbf{I} - \alpha \mathbf{H})$. Since \mathbf{H} is symmetric this is equivalent to the maximum eigenvalue.² So we can express the minimax as

$$\operatorname{argmin}_{\alpha} \max_{\mathbf{H}} \lambda_{\max}(\mathbf{I} - \alpha \mathbf{H})$$

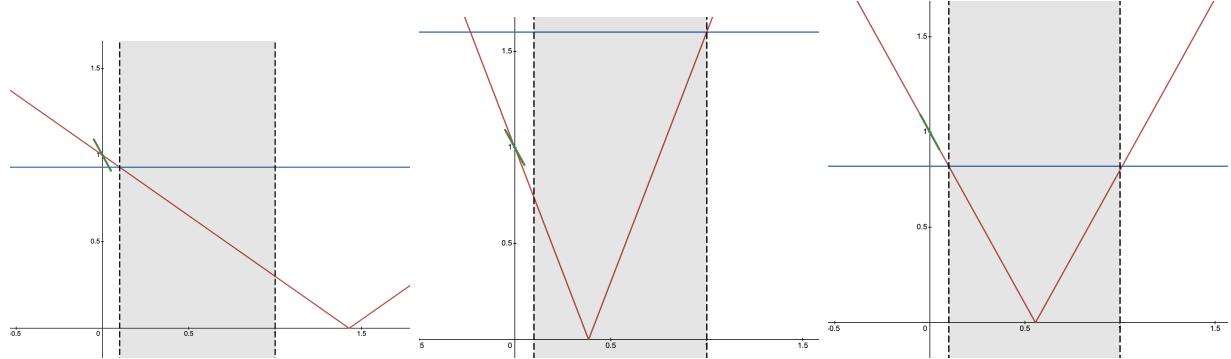
²More information on the matrix norm can be found here.

// TODO: CONSIDER ADDING SOME EXPLANATION OF SPECTRAL / MATRIX NORM? I BELIEVE THE ARGUMENT TO GET THAT ALL WE CARE ABOUT THE WORST CASE EIGENVALUE IS SEE THAT WE CAN DIAGONALIZE \mathbf{H} AS $\mathbf{Q}\Lambda\mathbf{Q}^\top$ AND SO ON.

Since we assumed $m\mathbf{I} \preceq \mathbf{H} \preceq M\mathbf{I}$ what we can further simplify to

$$\operatorname{argmin}_\alpha \max_{m \leq \lambda \leq M} |1 - \alpha\lambda|$$

There is a nice geometric interpretation to this problem. If we think $|1 - \alpha\lambda|$ as a linear function, we can think of sweeping the slope of an absolute function pinned at $(0, 1)$.



You can play around with this exact plot on Desmos here. // TODO: ADD SOME MORE DESCRIPTION FOR THE PLOTS.

Analytically, since the function is a (peice-wise) linear, we know that the maximum is at either boundary allowing us to write

$$\operatorname{argmin}_\alpha \max \{|1 - \alpha m|, |1 - \alpha M|\}$$

With arithmetic we can arrive at

$$\alpha^* = \frac{2}{M + m}$$

If we plug this into R_1 we get a **Convergence Rate**:

$$R_1 = \frac{M - m}{M + m} = 1 - \mathcal{O}\left(\frac{1}{k}\right)$$

// TODO: DEFINE THE CONDITION NUMBER AND WRITE OUT SOME MORE ARITHMETIC FOR IT
If we want **Iteration Complexity** (a.k.a. running time): Recall:

$$\text{error}_n \leq R_1^n \cdot \text{error}_0$$

Suppose we want to ask how many iterations need to get $\text{error}_n = \epsilon$

i.e.

$$R_1^n \leq \epsilon \Rightarrow n = k \log\left(\frac{1}{\epsilon}\right)$$

Note the fact: $1 - \delta = // TODO: RECALL THE LOG TRICK HERE AND WRITE IT$

3.3 Optimal Schedules

In the above section, we derived optimal rates for GD when we can pick a single fixed step-size. The natural follow-up is to ask if we can do better with non-constant step-sizes.

3.3.1 2 Step Schedules

For simplicity, we should start by computing an optimal schedule for $n = 2$ steps. The problem will take a familiar form,

$$\min_{\alpha, \beta} \max_{f, \mathbf{x}_0 \neq \mathbf{x}^*} \frac{\|\mathbf{x}_2 - \mathbf{x}^*\|}{\|\mathbf{x}_0 - \mathbf{x}^*\|}$$

where α, β are our step-sizes.

Suppose R_1^* was the optimal one-step contraction. Our hope would be that the extra flexibility of step-size would allow $\sqrt{R_2^*} < R_1^*$.

To analyze this, we can follow very similar steps to above. The two-step error recurrence can be expressed as

$$\mathbf{x}_2 - \mathbf{x}^* = (\mathbf{I} - \beta \mathbf{H})(\mathbf{I} - \alpha \mathbf{H})(\mathbf{x}_0 - \mathbf{x}^*)$$

Note that the contraction term is a matrix polynomial of degree ≤ 2 . // TODO: MAKE SURE THIS LOGIC IS CLEAR TO A FIRST TIME READER.

$$\begin{aligned} p(\mathbf{H}) &= (\mathbf{I} - \beta \mathbf{H})(\mathbf{I} - \alpha \mathbf{H}) \\ &= \mathbf{I} - (\alpha + \beta) \mathbf{H} + \alpha \beta \mathbf{H}^2 \end{aligned}$$

In this expression we see that $p(\mathbf{0}) = \mathbf{I}$. This is analogous to our plot of $y = |1 - \alpha \lambda|$ above where the function was pinned $(0, 1)$ at. Since p pinned at $(\mathbf{0}, \mathbf{I})$ has two degrees of freedom and is parameterized by two step-size choices α, β , we have a one-to-one mapping with the set $\hat{\mathcal{P}}_2 = \{p \in \mathcal{P}_2 : p(\mathbf{0}) = \mathbf{I}\}$. Squinting at $p(\mathbf{H})$ in factorized form will tell us that α, β are the inverse roots of p . Easiest way to see this is we pass a scalar x and get $p(x) = (1 - \beta x)(1 - \alpha x)$. Setting $p(x) = 0$ we get solutions $x \in \{\alpha^{-1}, \beta^{-1}\}$.

This lets us write the min of the minimax as

$$\min_{\text{poly } p \in \hat{\mathcal{P}}_2} \max_{m\mathbf{I} \preceq \mathbf{H} \preceq M\mathbf{I}} R_2 = \|p(\mathbf{H})\|$$

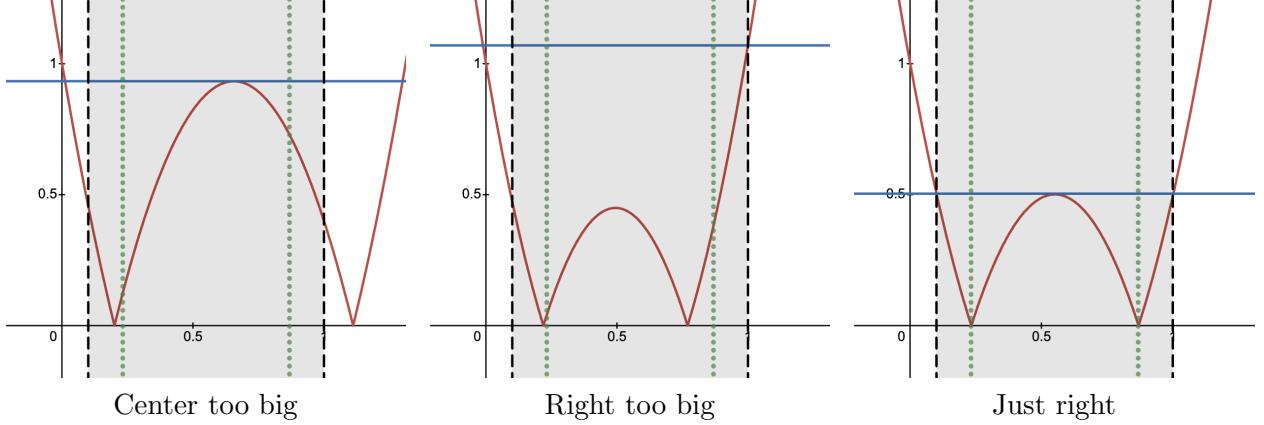
Because \mathbf{H} is symmetric, we may diagonalize as $\mathbf{H} = \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^\top$ and extract the orthogonal matrices out of the polynomial to get $p(\mathbf{H}) = \mathbf{Q} p(\boldsymbol{\Lambda}) \mathbf{Q}^\top$. Additionally, matrix-norms are invariant under orthogonal transforms so $\| \mathbf{Q} p(\boldsymbol{\Lambda}) \mathbf{Q}^\top \| = \| p(\boldsymbol{\Lambda}) \|$ meaning that our inner maximization reduces to the worst-case eigenvalue once again.

$$\min_{p \in \hat{\mathcal{P}}_2} \max_{m \leq \lambda \leq M} |p(\lambda)|$$

In fact, if we expand p , this simplified form is an exact generalization of the 1-step picture that we drew above!

$$\min_{\alpha, \beta} \max_{m \leq \lambda \leq M} |1 - (\alpha + \beta)\lambda + \alpha\beta\lambda^2|$$

We can plot the polynomial and the maximum value over $[m, M]$ and analyze the worst-case rate. An interactive plot is available [here](#).



// TODO: ADD MORE DESCRIPTION HERE.

The optimal roots (inverse step-sizes) are shown as the dotted green line at $\{\alpha^{-1}, \beta^{-1}\} = \left\{ \frac{M+m}{2} \pm \frac{M-m}{2\sqrt{2}} \right\}$

// TODO: ADD DERIVATION OF OPTIMAL HERE? IN CLASS HE JUST TOLD US WITH WAS THIS

3.3.2 The General Case

The general problem $n \geq 2$ steps with step sizes given by $\{\alpha_i\}_{i \in [n]}$ can be expressed as

$$R_n^* = \min_{\alpha, \beta} \max_{f, \mathbf{x}_0 \neq \mathbf{x}^*} \frac{\left\| \left(\prod_{i \in [n]} (\mathbf{I} - \alpha_i \mathbf{H}) \right) (\mathbf{x}_0 - \mathbf{x}^*) \right\|}{\|\mathbf{x}_0 - \mathbf{x}^*\|}$$

All of the above steps may be repeated to arrive at the form

$$\min_{p \in \mathcal{P}_n} \max_{m \leq \lambda \leq M} |p(\lambda)|$$

This problem is closely related to the Chebyshev Polynomials and the solution to the n -step problem will turn out to be exactly the n th Chebyshev Polynomials T_n (modulo some affine transforms on the input/output). The next section will dive deeper into what the Chebyshev Polynomials, but we will use them here without exposition to analyze some properties of multi-step schedules.

// TODO: ADD A ONE LINE EXPOSITION OF CHEBYSHEV POLYNOMIALS AND L

The solution to the minimax is given by

$$p_n^* = \frac{T_n(L(\lambda))}{T_n(L(0))}$$

The maximum value that T_n takes on is 1 so the value of the minimax is

$$R_n^* = \frac{1}{T_n(L(0))}$$

We get that

$$R_n = \frac{1}{T_n(L(0))} = \dots = \left(\frac{\sqrt{k} - 1}{\sqrt{k} + 1} \right)^n$$

$$\Rightarrow \# \text{ iterations to get to } R_n \leq \epsilon \text{ is } n = \Theta(\sqrt{k} \log(1/\epsilon))$$

// TODO: WRITE OUT THE JUSTIFICATION FOR THIS USING BIG N LIMIT FOR THE "K" TERM.

Philosophical Aside:

Why was the optimal steps for $n = 2$ suboptimal for $n = 1$.

It means that we have fewer choices to pick from we have to pick something "in-between" to trade off different edge cases ("hedging").

Young 1953. Acceleration Methods FnTML.

3.4 The Chebyshev Polynomials and Accelerated Methods

Some more resources (Mason & Handscomb, Riulin, Trefethen)

Various Definitions of Chebyshev Polys:

- **Explicit:**

$$T_n(z) = \frac{(z - \sqrt{z^2 - 1})^n + (z + \sqrt{z^2 - 1})^n}{2}$$

- **Trigonometric:**

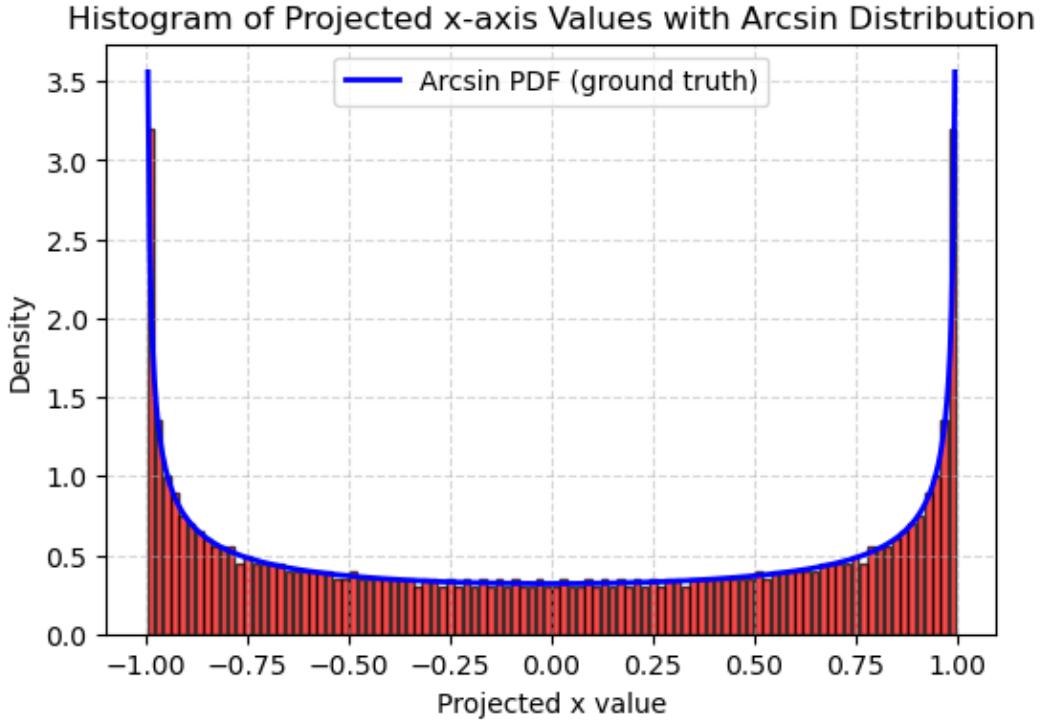
$$T_n(z) = \cos(n \cdot \arccos(z)), \quad z \in [-1, 1]$$

- **Roots:**

$$\left\{ \cos\left(\frac{2t+1}{2n}\pi\right) \right\}_{t=0 \dots n-1}$$

– $\{\frac{2t+1}{2n}\}_{t=0 \dots n-1}$ are pretty angles spread uniformly around the unit circle. So applying cos can be thought of as projecting the points from the unit circle onto the x-axis.

– Here is a nice picture for the distribution of roots (uniformly) on a circle then projected on to the x-axis. it is called the arcsine distribution and will come up in random step-size schedules



- **3-term recurrence**

$$T_{n+1} = 2zT_n(z) - T_{n-1}(z), \quad T_0(z) = 1, T_1(z) = z$$

- **Extremal:**

$$T_n(z)/2^{n-1} = \underset{p \text{ def } n, p \text{ monic}}{\operatorname{argmin}} \max_{z \in [-1,1]} |p(z)|$$

- This is why they showed up in our step-size derivation above.

A core message to takeaway here is that there are many equivalent ways to derive step-size schedules yielding accelerated descent *and* they all correspond to some way to intuit the Chebyshev polynomials. We can see this because to pick a "next" step size we need to know the "next" root of the n step polynomial.

To begin we can look at the Polyak Heavy Ball Method / Momentum and the 3-term recurrence.

To get step size for the next step we want, we can grab expand the contraction matrix polynomial using the recursive definition above. If we substitute in $\nabla f(\mathbf{x}_n)$ when possible we see that we can recover a momentum looking update rule which uses *both* the current gradient and the previous update to guide the next step: see the simplified form boxed below. The constants c_1, c_2 and c_3 are left undefined to avoid unnecessary arithmetic.

$$\begin{aligned} \mathbf{x}_{n+1} - \mathbf{x}^* &= p_{n+1}(\mathbf{H})(\mathbf{x}_0 - \mathbf{x}^*) \\ &= ((c_1 + c_2 \mathbf{H}) p_n(\mathbf{H}) - c_3 p_{n-1}(\mathbf{H})) (\mathbf{x}_0 - \mathbf{x}^*) \end{aligned}$$

POLYAK'S HEAVY BALL ALGORITHM

$$\text{Want } x_{n+1} - x^* = P_{n+1}(H) (x_0 - x^*)$$

$$= [c_1 + c_2 H] P_n(H) - c_3 P_{n-1}(H)$$

$$= \underbrace{(c_1 + c_2 H)}_{\substack{x_n - x \\ \text{---}}} P_n(H) (x_0 - x^*)$$

$$c_1(x_n - x^*) + c_2 \underbrace{H(x_n - x^*)}_{\nabla f(x)}$$

$$f(x) = \frac{1}{2} x^T H x - b^T x$$

$$\nabla f(x) = Hx - b = H(x - x^*)$$

\therefore Simplify.

But if you plug 3-term recurrence back into our $\mathbf{x}_{n+1} - \mathbf{x}^* = p_{n+1}(\mathbf{x}_0 - \mathbf{x}^*)$

Discussion of Conjugate gradients that I kind of missed, but conjugate gradients are based on orthogonal polynomials

3.5 Random Stepsizes

Suppose $\alpha_t \sim \mu$ iid. What is the optimal μ to sample from?

Informally our rate is

$$\max_{f, \mathbf{x}_0 = \mathbf{x}^*} \frac{\|\mathbf{x}_n - \mathbf{x}^*\|}{\|\mathbf{x}_0 - \mathbf{x}^*\|}$$

Suppose $\alpha_t \sim \mu$ iid. What is the optimal μ to sample from? What is the corresponding rate?

$$R(\mu) = \max_{f, \mathbf{x}_0} \lim_{n \rightarrow \infty} \mathbf{E} \left[\frac{\|\mathbf{x}_n - \mathbf{x}^*\|^{1/n}}{\|\mathbf{x}_0 - \mathbf{x}^*\|} \right]$$

Observations: rate is invariant w.r.t. ordering of the step-sizes $\{\alpha_t\}$

1. Optimal deterministic step sizes are the inverse roots of Chebyshev polynomials
2. Convergence rates depends only on $\hat{\mu}_n = \frac{1}{n} \sum \delta(\alpha_t)$ (i.e. uniform on the roots)
3. Intuitively the transform of the uniform on the unit circle is

3.6 Solving the Extremal Problem via Potential Theory:

- Step-size Problem (copy paste from above):

$$\min_{\gamma} \max_{\lambda \in [m, M]} \mathbf{E}_{\beta \sim \gamma} \log \left| 1 - \frac{\lambda}{\beta} \right|$$

- Electrostatics Problem:

Consider the following minimization of potentials on a line where the $\log 1/|\beta - \alpha|$ kernel defines the energy between particles

$$\min_{\gamma \in \mathcal{P}(m, M)} \mathbf{E}_{\beta, \lambda \sim \gamma} \log \frac{1}{|\beta - \alpha|}$$

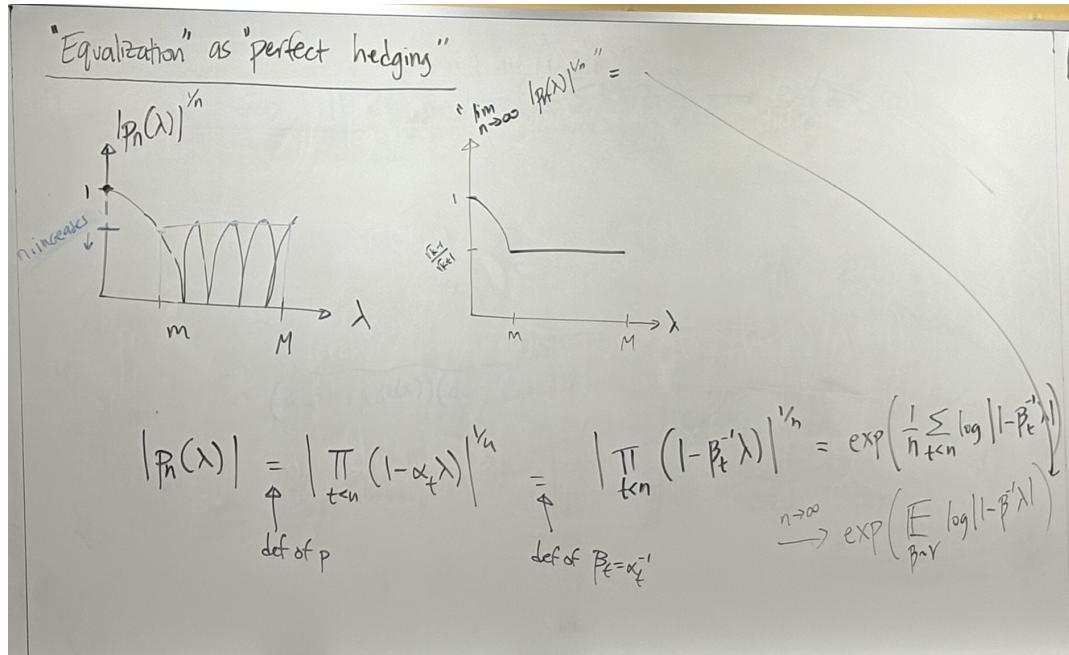
Fundamental Theorem of Potential Theory: Characterized uniquely by an equalizing property

$$\lambda \mapsto \mathbf{E} \log \left(\frac{1}{|\beta - \alpha|} \right)$$

In other words a sample from the continuum of points leads to the same log potential no matter where you are on the graph. One way to think about this intuitively is that this is a "stationary" potential since (if we allow particles to move around) there is no place where the potential is higher / lower than where at right now.

Jason provided some intuitive argument for the actual conclusion that arcsine is the optimal γ but it was more complex // **TODO: LOOK DEEPER INTO THIS.**

Connection to perfect hedging:



In the limit of Chebyshev step sizes what ends up happening is that we are "equally" hedged against everything. This means that **CHEBYSHEV BASED METHODS NEVER DO BETTER** than

$$\frac{\sqrt{k} - 1}{\sqrt{k} + 1}$$

There is some connection to Green's function that I should learn about: fundamentally, the step size optimum is given by the negative green's function evaluated on $[m, M]^c$ or something like this on the complex plane. // **TODO: LEARN MORE ABOUT THIS**

3.7 Understanding the best case scenarios via a two player game

- Min plays a step-size distribution
- Max plays a quadratic function

Game

$$Rate = \inf_{\gamma} \sup_{\lambda \in [m, M]} \mathbf{E}_{\beta \sim \gamma} \log |1 - \beta^{-1} \lambda|$$

To make the game more symmetric we can lift the function player into distribution over λ s giving

$$\inf_{\gamma \in \mathcal{P}(E)} \sup_{\rho \in \mathcal{P}(E)} \mathbf{E}_{\beta \sim \gamma, \lambda \sim \rho} \log |1 - \beta^{-1} \lambda|$$