

STAT 9912: Acceleration, Greed and Hedging in Optimization

Taught by Dr. Jason Altschuler
Notes written by Faraz Radhman

Contents

1	Setting the Stage	2
1.1	Optimization and Gradient Descent	2
1.2	Why $-\nabla f$?	2
2	Quadratics	3
2.1	Do we converge?	3
2.2	How fast do we converge?	4
2.3	Optimal Schedules	5
2.3.1	Warmup: Spectral Annihilation	6
2.3.2	Optimal 2-Step Schedules	7
2.3.3	The General Case	8
2.4	The Chebyshev Polynomials and Accelerated Methods	10
2.4.1	Polyak's Heavy Ball	11
2.4.2	Random Stepsizes via Potential Theory	12
2.5	Understanding the best case scenarios via a two player game	14
3	Nesterov's Acceleration and Curvature	14

1 Setting the Stage

1.1 Optimization and Gradient Descent

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

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

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 to solve this problem is Gradient Descent (GD) and is described by the following update

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

where $\{\alpha_t\}$ parameterizes different "schedules" for GD.

Here are some questions to start discussion:

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

1.2 Why $-\nabla f$?

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.

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

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

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

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.

2 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 $\mathbf{H} \succeq 0$ ¹ (necessary for f to be convex) and \mathbf{H} is symmetric.²

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.

2.1 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

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

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

¹ $\mathbf{H} \succeq 0 \Leftrightarrow \mathbf{x}^\top \mathbf{H} \mathbf{x} \geq 0, \quad \forall \mathbf{x} \in \mathbb{R}^m \Leftrightarrow \mathbf{H}$ is positive semidefinite (PSD) by defn. $\Leftrightarrow \mathbf{x}^\top \mathbf{H} \mathbf{x}$ is convex.

²We 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}

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

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

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} .

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.

2.2 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.³ So we can express the minimax as

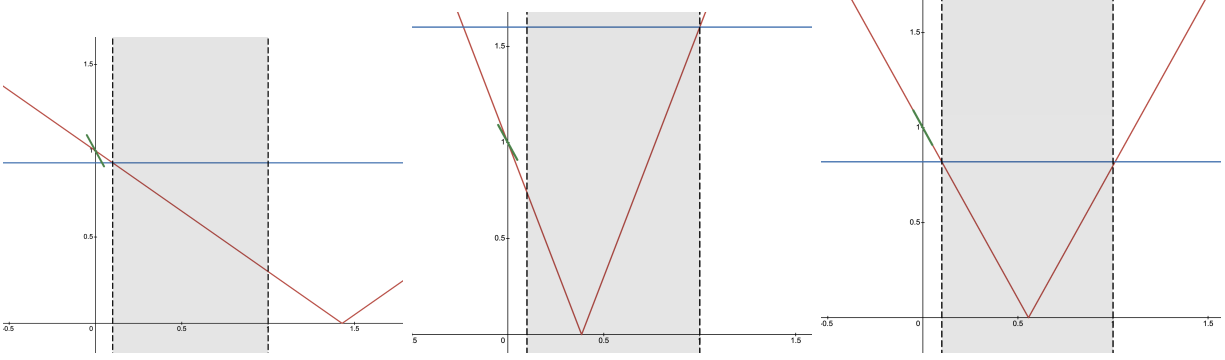
$$\operatorname{argmin}_{\alpha} \max_{\mathbf{H}} \lambda_{\max}(\mathbf{I} - \alpha \mathbf{H})$$

³More information on the matrix norm can be found here.

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](#).

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**:

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

If we want **Iteration Complexity** (a.k.a. running time): Recall:

$$error_n \leq R_1^n \cdot error_0$$

Suppose we want to ask how many iterations need to get $error_n = \epsilon$

i.e.

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

Note the fact: $1 - \delta =$

2.3 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.

2.3.1 Warmup: Spectral Annihilation⁴

We spent some time deriving the optimal constant step-size for worst-case quadratic problems with a given minimum and maximum eigenvalue. Before continuing, let's build understanding about what we can do if we have access to the entire eigen-spectrum of our problem.

In particular, let's start with the 1d-case with

$$f(x) = \frac{1}{2}Ax^2$$

The gradient $\nabla f(x) = Ax$, so after one step of gradient descent with step-size α , the next iterate is

$$x^+ = x - \alpha Ax$$

You'll notice that simply setting $\alpha = \frac{1}{A}$ yields

$$x^+ = x - \frac{A}{A}x = 0 = x^*$$

Allowing us to solve in one iteration. Now consider that we have a 2d problem.

Consider diagonal $A \in \mathbb{R}^{2 \times 2}$ and

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{A}\mathbf{x} = \frac{1}{2}(A_{11}x_1^2 + A_{22}x_2^2)$$

The gradient $\nabla f(\mathbf{x}) = \mathbf{A}\mathbf{x}$, so after one step of gradient descent with step-size α , the next iterate is

$$x_1^+ = x_1 - \alpha A_{11}x_1$$

$$x_2^+ = x_2 - \alpha A_{22}x_2$$

Like the 1d case, we can simply select the diagonal elements (eigen-values for general matrices) to be our reciprocal step-sizes to arrive knock the corresponding entry right to the optimum $\mathbf{x}^* = \mathbf{0}$. Let $\alpha_1 = \frac{1}{A_{11}}$, $\alpha_2 = \frac{1}{A_{22}}$, then after one step we knock out x_1

$$x_1^{++} = x_1^+ - \frac{A_{11}}{A_{11}}x_1^+ = 0$$

$$x_2^{++} = x_2^+ - \frac{A_{22}}{A_{11}}x_2^+ = \left(1 - \frac{A_{22}}{A_{11}}\right)x_2^+$$

And after the second step we get knock x_2 .

$$x_1^+ = x_1 - \frac{1}{A_{22}} \cdot 0 = 0$$

$$x_2^+ = x_2 - \frac{A_{22}}{A_{22}}x_2 = 0$$

⁴Not included in original lectures.

This logic extends to d -dimensional problems and non-diagonal problems where we arrive at the optimum in d -steps. An important takeaway here is that **the order of step-sizes do not matter**. Because, the eigen-spectrum of a quadratic function is identical at all points, the step with step-size reciprocal to the i th eigenvalue will annihilate the i th entry *no matter where it currently is in the domain* i.e. no matter where the previous descent steps went. This is a nice property of quadratic functions but it does not generalize to quadratics, we will see this a common theme that shows up in later algorithms.

2.3.2 Optimal 2-Step Schedules

With the above insight, we can return to the case where we only assume bounds $0 < m \leq \lambda_i \leq M$ on the eigenvalues. 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 .

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

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

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

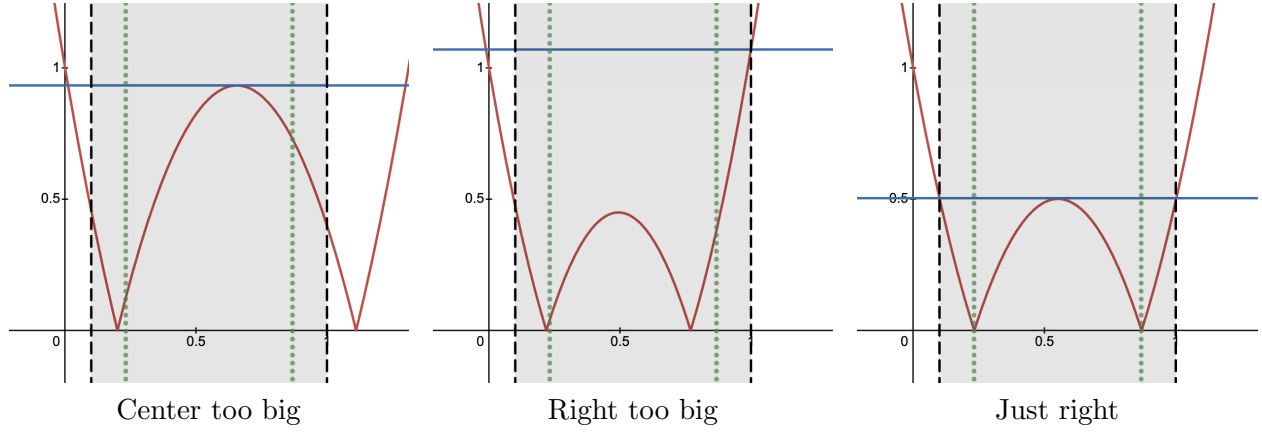
Because \mathbf{H} is symmetric, we may diagonalize as $\mathbf{H} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ and extract the orthogonal matrices out of the polynomial to get $p(\mathbf{H}) = \mathbf{Q}p(\mathbf{\Lambda})\mathbf{Q}^\top$. Additionally, matrix-norms are invariant under orthogonal transforms so $\|\mathbf{Q}p(\mathbf{\Lambda})\mathbf{Q}^\top\| = \|p(\mathbf{\Lambda})\|$ 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](#).



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

2.3.3 The General Case

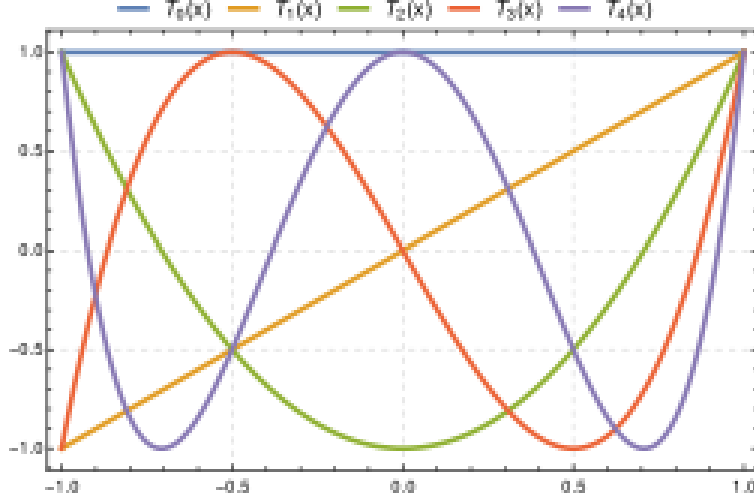
The general problem $n \geq 2$ steps with step sizes given by $\{\alpha_i\}_{i \in [n]}$ can be expressed as

$$R_n^* = \min_{\alpha, \beta} \max_{f, x_0 \neq x^*} \frac{\left\| \left(\prod_{i \in [n]} (I - \alpha_i H) \right) (x_0 - x^*) \right\|}{\|x_0 - x^*\|}$$

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

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

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 are, but for now we will just provide a graph and some rough intuition for why they are the solution.



Chebyshev Polynomials $T_n(x)$ for degrees $n = 1$ through $n = 5$ on $[-1, 1]$.

You can see in the above figure that the extrema of the Chebyshev polynomials lie in $\{-1, 1\}$. It should be intuitive how this corresponds to the absolute value of the polynomial all sharing a common maxima.

Take for granted that the solution to the minimax is

$$p_n^* = \frac{T_n(L(\lambda))}{T_n(L(0))} \quad (1)$$

where L affinely maps the eigenvalue range $[m, M] \rightarrow [-1, 1]$ which is the natural domain of the Chebyshev polynomial.⁵

The maximum value that T_n takes on is 1 so maximizing over λ we get

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

We can understand the limiting optimal rate by taking the large n limit:

We get that

$$\begin{aligned} \lim_{n \rightarrow \infty} R_n^* &= \lim_{n \rightarrow \infty} \frac{1}{T_n(L(0))} \\ &= \lim_{n \rightarrow \infty} \frac{2}{\left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^n + \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^{-n}} \\ &= 2 \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^n \\ &= 2 \left(1 - \frac{2}{\sqrt{\kappa}+1}\right)^n \end{aligned}$$

⁵This works out to be $L(\lambda) = \frac{\lambda-m}{M-m} \cdot 2 - 1$. For $\lambda = 0$ we get $L(0) = \frac{-2m}{M-m} - 1 = \frac{-M-m}{M-m} = \frac{m+M}{m-M} = \frac{1+\kappa}{1-\kappa}$

The second line follows from the first definition of the Chebyshev polynomials which shows up in the following section. For now it suffices to take this as true.

Taking the geometric mean over n steps and dropping constants we see that the final optimal rate is

$$\mathcal{O}\left(\frac{1}{\sqrt{\kappa}}\right)$$

which is the famous “accelerated” rate for gradient descent.

Similarly we can get an iteration complexity $R_n \leq \epsilon$ is $n = \Theta\left(\sqrt{k} \log(1/\epsilon)\right)$

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.

2.4 The Chebyshev Polynomials and Accelerated Methods

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

Various Definitions of Chebyshev Polys:

- **Explicit:**

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

- **Trigonometric:**

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

- **Roots:**

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

– $\left\{ \frac{2t+1}{2n}\pi \right\}_{t=0 \dots n-1}$ are 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.

– Figure 1 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**

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

- **Extremal:**

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

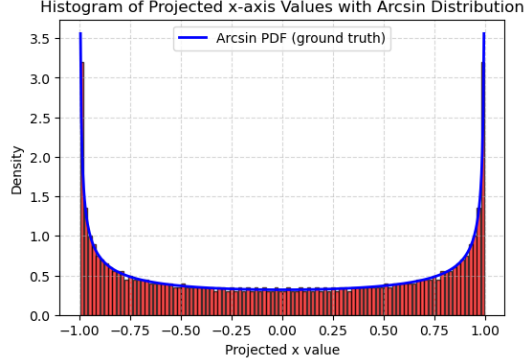


Figure 1: Distribution of roots of Chebyshev polynomials projected onto the x-axis (arcsine distribution).

- 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 define the Chebyshev polynomials. This duality follow directly from the relationship between inverse roots and step-sizes.

2.4.1 Polyak’s Heavy Ball

To begin we can look at the Polyak Heavy Ball Method which introduced the notion of momentum to first-order optimization. Polyak’s Heavy Ball can be dervied via the 3-term recurrence.

Recall that the 3-term recurrence was

$$T_{n+1}(z) = 2zT_n(z) - T_{n-1}(z)$$

We need to shift / rescale in order to apply this recurrence to $p_{n+1}(\mathbf{H})$, our optimal contraction polynomials. We defined $L(\lambda) = a\lambda + b$ to handle scaling the input ⁶. The RHS of the recurrence can then be expressed in terms of lambda $2(a\lambda + b)T_n(L(\lambda)) - T_{n-1}(L(\lambda))$, but we also need to rescale the out put as in (1) so taking this into accoun the optimal polynomial recurrence is

$$p_{n+1}(\lambda) = ((c_1 + c_2\lambda)p_n(\lambda) - c_3p_{n-1}(\lambda))$$

for some values of c_1, c_2, c_3

Next, we want to go from this recurrence into something that we can implement. Let’s start with some wishful thinking: for $n + 1$ steps assume we contract according to $p_{n+1}(\mathbf{H})$ which is optimal, then write the contraction via the recurrence.

$$\begin{aligned} \mathbf{x}_{n+1} - \mathbf{x}^* &= p_{n+1}(\mathbf{H})(\mathbf{x}_0 - \mathbf{x}^*) \\ &= ((c_1 + c_2\mathbf{H})p_n(\mathbf{H}) - c_3p_{n-1}(\mathbf{H}))(\mathbf{x}_0 - \mathbf{x}^*) \\ &= (c_1 + c_2\mathbf{H})\underbrace{p_n(\mathbf{H})(\mathbf{x}_0 - \mathbf{x}^*)}_{=(\mathbf{x}_n - \mathbf{x}^*)} - c_3\underbrace{p_{n-1}(\mathbf{H})(\mathbf{x}_0 - \mathbf{x}^*)}_{=(\mathbf{x}_{n-1} - \mathbf{x}^*)} \end{aligned}$$

⁶for some values a, b (we will only be concerned with the functional form for now)

Where the third line points out that the recurrence allows us to express the optimal $n + 1$ step iterate in terms of our n and $n - 1$ step iterates. Further rearranging gives.

$$\begin{aligned}\mathbf{x}_{n+1} - \mathbf{x}^* &= c_1 (\mathbf{x}_n - \mathbf{x}^*) + c_2 \underbrace{\mathbf{H}(\mathbf{x}_n - \mathbf{x}^*)}_{=\nabla f(\mathbf{x}_n)} - c_3 (\mathbf{x}_{n-1} - \mathbf{x}^*) \\ &= c_1 \mathbf{x}_n + c'_3 \mathbf{x}_{n-1} + c_2 \nabla f(\mathbf{x}_n) - (c'_3 + c_1) \mathbf{x}^*\end{aligned}$$

The first line identifies the second term as the gradient for the quadratic function. In rearranging we substitute $c'_3 = -c_3$ for clarity. If we constraint $c_1 + c'_3 = 1$ we can add \mathbf{x}^* to both sides to get

$$\begin{aligned}\mathbf{x}_{n+1} &= c_1 \mathbf{x}_n + c'_3 \mathbf{x}_{n-1} + c_2 \nabla f(\mathbf{x}_n) \\ &= \mathbf{x}_n - c'_3 \underbrace{(\mathbf{x}_n - \mathbf{x}_{n-1})}_{\text{momentum}} + c_2 \nabla f(\mathbf{x}_n)\end{aligned}$$

The final step rearranges the update rule in terms of a “momentum” term. Inductively, if we assumed that $\mathbf{x}_n, \mathbf{x}_{n-1}$ were the optimal iterates, then by construction it follows that \mathbf{x}_{n+1} is optimal up to $n + 1$ step step-size schedule on vanilla gradient descent, thus as we take many updates $n \rightarrow \infty$ the algorithm achieves the accelerated rate.

If we were to have carried values for the constants throughout, we would have arrived at

$$c'_3 = \left(\frac{2}{\sqrt{M} + \sqrt{m}} \right)^2, \quad c_2 = \left(\frac{\sqrt{M} - \sqrt{m}}{\sqrt{M} + \sqrt{m}} \right)^2$$

The relationship between momentum and step-size parameters c'_3 and c_2 is worth its own section. Fabian Pedregosa’s 2023 ICML Blog *A Hitchhiker’s Guide to Momentum* may be of interest for further reading and includes derivation of the regions for which (c'_3, c_2) achieve the accelerated rate.

2.4.2 Random Stepsizes via Potential Theory

Note that the derivation for optimal step-size schedules in 2.3.3 **never** made an assumption on the ordering of step-sizes. Informally, as number n steps we are optimizing grows large, taking a random ordering of these step-sizes can be approximated by sampling iid from some distribution - that is something like $\mu_n = \frac{1}{n} \sum_t \delta(\alpha_t)$ for large n . The figure 1 would suggest that the limiting distribution $\mu_n \rightarrow \mu$ is the arcsine distribution.

Alternatively, we can give theoretical justification for μ by showing that it is the optimal distribution. To begin, will express our problem using the matrix norm.

$$\begin{aligned}R_n(\omega) &= \max_f \left\| \prod (\mathbf{I} - \alpha_t(\omega) \mathbf{\Lambda}) \right\| \\ &= \max_{\lambda \in [m, M]} \left\| \prod |1 - \alpha_t(\omega) \lambda| \right\|\end{aligned}$$

Where $\alpha_t(\omega)$ are function of the random draw from μ , and the second line diagonalizes. The optimization problem we will be interested in for the random case is

$$\min_{\mu} \lim_{n \rightarrow \infty} \mathbf{E}_{\{\alpha_t\} \sim \mu^n} \left[\sqrt[n]{R_n} \right]$$

We want to study how the geometric mean converges. To do this, we can analyze the log one-step contraction rate. In general, if X_i are distributed i.i.d and $G_n := (\prod X_i)^{1/n}$ is the geometric mean, then $\mathbf{E} \log X_i = m$ implies that $G_n \rightarrow e^m$ almost surely.⁷ Thus it suffices to study the problem minimizing

$$\log R(u) = \max_{\lambda \in [m, M]} \mathbf{E} \left(\log \left| 1 - \frac{\lambda}{\alpha} \right| \right)$$

since it will also minimize the geometric average of the rates. To solve the problem cleanly we can make a connection to electrostatic potential theory. To begin, we write the two problems as follow:

- Step-size Problem: copy paste from aboves but we use $\beta = \alpha^{-1}$ to minimize over the inverse-step-size distribution γ since it allows us to make a cleaner connection to potential theory.

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

- Electrostatics Problem: Minimization of potentials on a line where the $\log 1/|\beta - \alpha|$ kernel defines the energy between particles

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

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

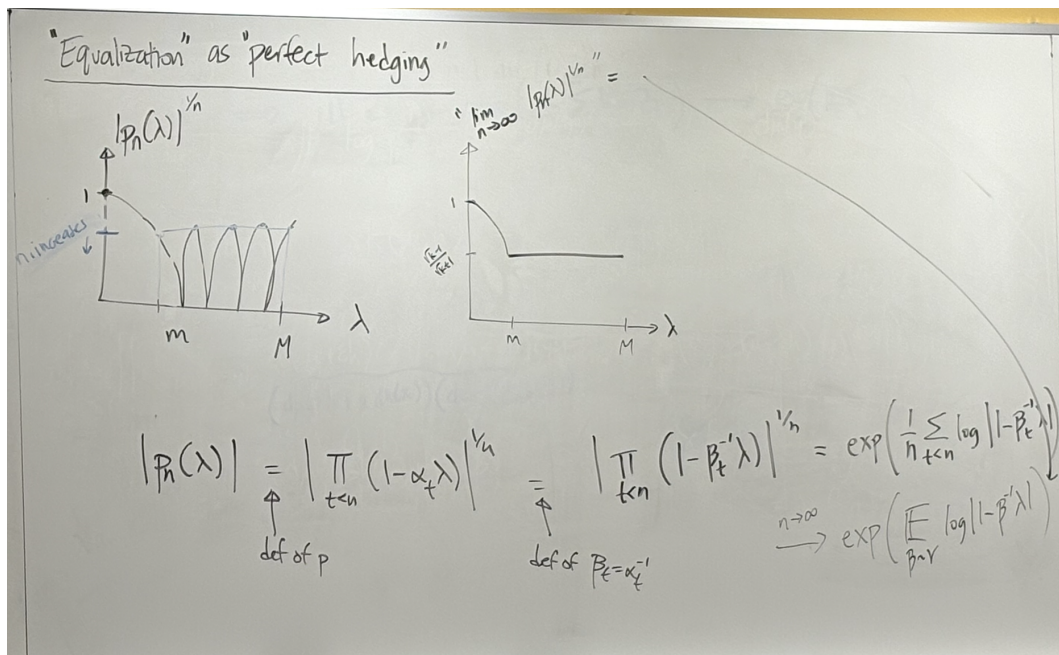
$$\lambda \mapsto \mathbf{E} \log \left(\frac{1}{|\beta - \alpha|} \right)$$

In other words, suppose we sample from the continuum of points leading 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 it is in the distribution.

Jason provided some intuitive argument for the actual conclusion that arcsine is the optimal γ but it was more complex

Connection to perfect hedging:

⁷See further discussion at this Math Stack Exchange post.



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.

2.5 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

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

3 Nesterov’s Accelertion and Curvature

This section is based on work from Su, Boyd and Candes (2015) and Shi et al. (2018).

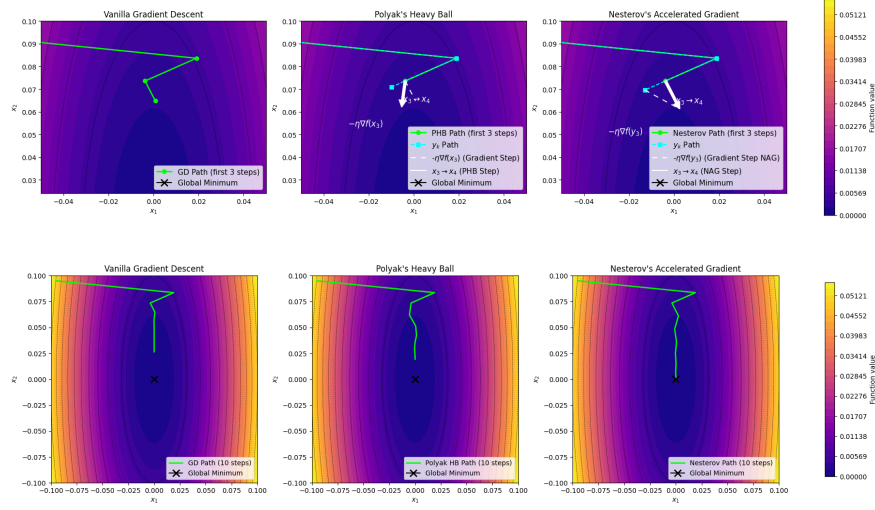


Figure 2: Top row: Illustration of Nesterov's accelerated gradient descent steps on the quadratic function $f(\mathbf{x}) = 10x_1^2 + x_2^2$. Bottom row: Comparison between 10 steps of vanilla gradient descent, Nesterov acceleration, and Polyak heavy ball methods on the same function

Acceleration over the class of L -smooth convex functions has proven to be a more challenging problem. In 1983, Yurii Nesterov proposed a simple modification to the gradient descent rule yielding the $O(1/\sqrt{T})$ rate to match the quadratic setting. Termed, Nesterov acceleration, the update can be written as:

$$\begin{aligned} \mathbf{y}_k &= \mathbf{x}_k + \beta_k(\mathbf{x}_k - \mathbf{x}_{k-1}) \\ \mathbf{x}_{k+1} &= \mathbf{y}_k - \eta \nabla f(\mathbf{y}_k) \end{aligned}$$

where β_k is a momentum parameter (which can be tuned or scheduled) and L is the smoothness parameter of f .

Visually, this update first extrapolates \mathbf{x}_k in the direction of its last iteration ($\mathbf{x}_k - \mathbf{x}_{k-1}$) to get \mathbf{y}_k , and then performs a gradient step from the extrapolated point \mathbf{y}_k . We can visualize the updates on quadratic $f(\mathbf{x}) = 10x_1^2 + x_2^2$ in Figure 2

The proof for NAG is complicated, so we are going to leave it out for now, but there is useful intuition in comparing the algorithm to Polyak's heavy ball. The updates are almost identical except in Nesterov's the gradient is computed at the point \mathbf{y}_k as opposed to \mathbf{x}_k .

$$\begin{aligned} \text{Polyak's Update} \quad \mathbf{x}_{k+1} &= \mathbf{y}_k - \eta \nabla f(\mathbf{x}_k) \\ \text{Nesterov's Update} \quad \mathbf{x}_{k+1} &= \mathbf{y}_k - \eta \nabla f(\mathbf{y}_k) \end{aligned}$$

As in previous algorithms the challenge of adapting quadratic methods to general convex functions is that curvature may vary from one point to another. Intuitively, we might hope that $\nabla f(\mathbf{y}_k)$ provides more of curvature information than $\nabla f(\mathbf{x}_k)$.

We can start to see this by writing $\nabla f(\mathbf{y}_k)$ as a Taylor approximation around $\nabla f(\mathbf{x}_k)$.

$$\begin{aligned}\nabla f(\mathbf{y}_k) &\approx \nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k)(\mathbf{y}_k - \mathbf{x}_k) \\ &= \nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k) \beta_k(\mathbf{x}_k - \mathbf{x}_{k-1})\end{aligned}$$

For the quadratic case, this approximation is exact.