

ONLY TAILS MATTER: AVERAGE-CASE UNIVERSALITY AND ROBUSTNESS IN THE CONVEX REGIME

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent works have studied the average convergence properties of first-order optimization methods on distributions of quadratic problems. The average-case framework allows a more fine-grained and representative analysis of convergence than usual worst-case results, in exchange for a more precise hypothesis over the data generating process, namely assuming knowledge of the expected spectral distribution (e.s.d) of the random matrix associated with the problem. This work shows that a problem’s asymptotic average complexity is determined by the concentration of eigenvalues near the edges of the e.s.d. We argue that having a priori information on this concentration is a more grounded assumption than complete knowledge of the e.s.d. basing our analysis on the approximate concentration is effectively a middle ground between the coarseness of the worst-case scenario convergence and the restrictive previous average-case analysis. We introduce the Generalized Chebyshev method, asymptotically optimal under a hypothesis on this concentration, and globally optimal when the e.s.d. follows a Beta distribution. We compare its performance to classical optimization algorithms, such as Gradient Descent or Nesterov’s scheme, and we show that, asymptotically, Nesterov’s method is universally nearly optimal in the average-case.

1 INTRODUCTION

The analysis of the average complexity of algorithms has a long story in computer science. Average-case complexity, for instance, drives much of the decisions made in cryptography (Bogdanov & Trevisan, 2006).

Despite their relevance, average-case analyses are difficult to extend to other algorithms, partly because of the intrinsic issue of defining a typical distribution over problem instances. Recently though, Pedregosa & Scieur (2020) derived a framework to systemically evaluate the complexity of first-order methods when applied on distributions of quadratic minimization problems. This is done by relating the average-case convergence rate to the *expected spectral distribution* (e.s.d) of the objective function’s Hessian, which is a well-studied object on random matrix theory. Having access to this object in practice is a much stronger hypothesis when compared to the worst-case analysis that relies only on the values of the edges of this distribution.

Paquette et al. (2020) extended the average-case framework by introducing a noisy generative model for the problems. They further derived the average complexity of the Nesterov Accelerated Method (Nesterov, 2003) on a particular distribution. They showed the strong concentration of the metrics around a limiting value as dimensions go to infinity.

Scieur & Pedregosa (2020) showed that for a strongly convex problem with eigenvalues supported on a contiguous interval, the optimal average-case complexity converges asymptotically to the one given by the Polyak Heavy Ball method (Polyak, 1964) in the worst-case.

1.1 CURRENT LIMITATIONS OF THE AVERAGE-CASE ANALYSIS

When analyzing the state of the art of average-case methods on quadratics problems, we observe significant limitations that we address in this paper. First, little is known

about the convergence rate on **convex problems**. Also, optimal average-case algorithms require an **exact estimation of the e.s.d** to guarantee an optimal convergence rate, their convergence rate under inexact e.s.d. is not known. Finally, the **non-smooth** is also discussed in (Pedregosa & Scieur, 2020), but with little details.

Convex problems. The minimization of non-strongly convex problems is drastically slower than their strongly convex counterpart, as Gradient Descent presents worst-case convergence in $\left(\frac{1}{7}\right)$ and Nesterov is $\left(\frac{1}{22}\right)$. In the strongly convex case, both the worst-case and average-case are asymptotically equal. However, little is known on optimal average-case rates for convex problems, as well as the average-case complexity of classical methods such as gradient descent or Nesterov’s method, see (Paquette et al., 2020).

Exact estimation of the e.s.d. In (Pedregosa & Scieur, 2020), the theoretical study of optimal algorithms in the average-case requires an exact estimation of the e.s.d. of the problem class. Such estimation may be hard, nor impossible to obtain in practical scenarios. Despite showing good performance when the e.s.d. is estimated with empirical quantities, there are no theoretical guarantees on the performance of the method when the e.s.d. is poorly estimated. There is therefore a need to analyze the algorithm’s performance under different notions of uncertainty on the spectrum. This allows a practitioner to choose the best algorithm for a practical problem, even with imperfect *a priori* information.

Non-smooth. Pedregosa & Scieur (2020) briefly introduce average-case optimal rates on non-smooth problems, when the e.s.d. is the Laguerre distribution e^{-x} . In this paper, we extend the analysis to the generalized Laguerre distribution $e^{-x} x^\alpha$; $\alpha > 1$.

1.2 CONTRIBUTIONS

Our main contribution is a fine-grained analysis of the average-case complexity on convex quadratic problems: we show that a problem’s complexity depends on the concentration of the eigenvalues of e.s.d. around the edges of their support. From this perspective, we propose a family of **optimal algorithms** in the average-case, analyze their **robustness**, and finally exhibit a **universality** result for Nesterov’s method. More precisely,

- **(Optimal algorithms).** In Section 3, we propose the Generalized Chebyshev Method (GCM, Algorithm 1), a family of algorithms whose parameters depend on the concentration of the e.s.d. around the edges of their support. If the parameters of the GCM method are set properly, the algorithm converges at an optimal average-case rate (Theorem 3 for smooth problems, Theorem 6 for non-smooth problems), a rate that we show is faster than worst-case optimal methods like Nesterov acceleration. We show these rates to be representative of the practical performance of the algorithms in Fig. 6, and retrieve the classical worst-case rates as limits of the average-case (see Table 1).
- **(Robustness).** Developing an optimal algorithm requires the knowledge of the exact e.s.d. However, in practical scenarios, we only have access to an *approximation* of the e.s.d. In Theorem 2 in Section 4 we analyze the rate of GCM in the presence of such a mismatch. We also analyze the optimal average-case rates of distributions representing the smooth convex, non-smooth convex, and strongly convex settings and compare them with the worst-case rates (Table 1).
- **(Universality).** Finally, in Theorem 4, we analyze the asymptotic average-case convergence rate of Nesterov’s method. We show that its convergence rate is nearly optimal (up to a logarithmic factor) under some natural assumptions over the data, namely a concentration of eigenvalues around 0 similar to the Marchenko-Pastur measure. This contributes to the theoretical understanding of the numerical efficiency of Nesterov’s acceleration.

2 AVERAGE-CASE ANALYSIS

In this section, we recall the average-case analysis framework for random quadratic problems. The main result is Theorem 1, which relates the expected error to the *expected spectral*

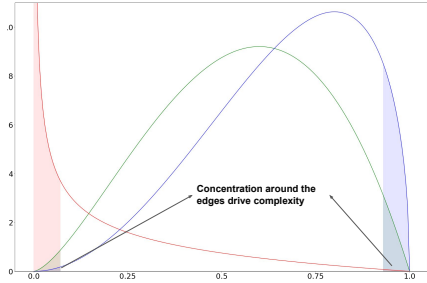


Figure 1: Representation of different spectra with different concentrations of eigenvalues around the edges of the support. The average-case rates for non-strongly problems are determined by these concentrations

Regime	Worst-case	Average-Case
Strongly conv.	$(1 - (1-\sqrt{\kappa}))^t$	$(1 - (1-\sqrt{\kappa}))^t$
Smooth conv.	$1-t^2$	$1-t^2 + \kappa$
Convex	$1-\sqrt{t}$	$1-t^{-1/2}$

Table 1: Comparison between function value worst-case and average-case convergence. κ is the condition number in the smooth strongly convex case. In the smooth convex case $\kappa > 1$ is the concentration of eigenvalues around 0 (see Assumption 1) and in the non-smooth case we consider $d \propto e$

distribution and the residual polynomial. The one-to-one correspondence between the residual polynomials and first-order methods applied to quadratics will allow us to pose the problem of finding an optimal method as the best approximation problem in the space of polynomials.

We define a **random** quadratic problem:

Problem 1. Let $H \in \mathbb{R}^{d \times d}$ be a random symmetric positive-definite matrix independent to $\mathbf{x}^* \in \mathbb{R}^d$, a random vector that is the solution to the problem. We define the random quadratic minimization problem as

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T H (\mathbf{x} - \mathbf{x}^*) \quad (\text{OPT})$$

We are interested in minimizing the expected errors $\mathbb{E} \|f(\mathbf{x}_t) - f(\mathbf{x}^*)\|$, the expected function-value gap, and $\mathbb{E} \|\nabla f(\mathbf{x}_t)\|^2$, the expected gradient norm, where \mathbf{x}_t is the t -th update of a first-order method starting from \mathbf{x}_0 and \mathbb{E} is the expectation over the random variables H ; \mathbf{x}_0 and \mathbf{x}^* .

The expectation we consider is over the problem and not over any randomness of the algorithm.

In this paper, we consider the class of *first-order methods* (F.O.M's) to minimize (OPT). Methods in this class construct the iterates \mathbf{x}_t as

$$\mathbf{x}_t \in \mathbf{x}_0 + \text{span}\{\nabla f(\mathbf{x}_0), \dots, \nabla f(\mathbf{x}_{t-1})\} \quad (1)$$

That is, \mathbf{x}_t belongs to the span of previous gradients. This class of algorithms includes for instance gradient descent and momentum, but not quasi-Newton methods since the preconditioner could allow the iterates to go outside of the span. Furthermore, we will only consider *oblivious* methods, that is, methods in which the coefficients of the update are known in advance and do not depend on previous updates. This leaves out some methods such as conjugate gradient or methods with line-search.

From First-Order Method to Polynomials. There is an intimate link between first-order methods and polynomials that simplifies the analysis of quadratic objectives. The next proposition shows that, with this link, we can assign to each optimization method a polynomial that determines its convergence. Following Fischer (1996), we will say a polynomial P_t is *residual* if $P_t(0) = 1$.

Proposition 1. (Hestenes et al., 1952) Let \mathbf{x}_t be generated by a first-order method. Then there exists a residual polynomial P_t of degree t , that verifies

$$\mathbf{x}_t - \mathbf{x}^* = P_t(H)(\mathbf{x}_0 - \mathbf{x}^*) \quad (2)$$

Remark 1. If the first-order method is further a momentum method, i.e.

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

We can determine the polynomials by the recurrence $P_0 = 1$ and

$$P_{t+1}(\cdot) = P_t(\cdot) + h_t \nabla P_t(\cdot) + m_t(P_t(\cdot) - P_{t-1}(\cdot));$$

We note that while most popular F.O.M's can be posed as a momentum method, the Nesterov method cannot.

A convenient way to collect statistics on the spectrum of a matrix is through its *empirical spectral distribution*.

Definition 1 (Expected spectral distribution (e.s.d)). . Let \mathbf{H} be a random matrix with eigenvalues $\lambda_1, \dots, \lambda_d$. The empirical spectral distribution of \mathbf{H} , called $\mu_{\mathbf{H}}$, is the probability measure

$$\mu_{\mathbf{H}} := \frac{1}{d} \sum_{i=1}^d \delta_{\lambda_i}; \quad (3)$$

where δ_{λ} is the Dirac delta, a distribution equal to zero everywhere except at λ and whose integral over the entire real line is equal to one.

Since \mathbf{H} is random, the empirical spectral distribution $\mu_{\mathbf{H}}$ is a random variable in the space of measures. Its expectation over \mathbf{H} is called the expected spectral distribution and we denote it

$$\mu := \mathbb{E}_{\mathbf{H}}[\mu_{\mathbf{H}}]; \quad (4)$$

We can link the e.s.d. of \mathbf{H} to the convergence of a first-order method on the distribution of \mathbf{H} . In the following we will consider $\mathbf{x}_0 \sim \mathbf{x}^?$ and \mathbf{H} to be independent, with $\mathbf{x}_0 \sim \mathbf{x}^?$ sampled isotropically.

Theorem 1. Let \mathbf{x}_t be generated by a first-order method associated to the polynomial P_t , the measure μ the e.s.d. of \mathbf{H} , and $\mathbb{E}[(\mathbf{x}_0 - \mathbf{x}^?)(\mathbf{x}_0 - \mathbf{x}^?)^T] = R^2 I$ for some constant R . Then we can write the convergence metrics at time step t as

$$\begin{aligned} \mathbb{E}[k\|\mathbf{x}_t - \mathbf{x}^?\|^2] &= R^2 \int P_t^2(\cdot) d(\cdot); & \mathbb{E}[\|\nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}^?)\|^2] &= R^2 \int P_t^2(\cdot) d(\cdot) \\ \text{and} & & \mathbb{E}[\|\nabla f(\mathbf{x}_t)\|_2^2] &= R^2 \int P_t^2(\cdot)^2 d(\cdot); \end{aligned} \quad (5)$$

This shows that polynomials are a powerful abstraction as they allow us to write all of our convergence metrics within the same framework. For simplicity, we set $R^2 = 1$ and we will refer directly to the polynomials associated to a given method. We will refer to objective f as the one associated to the added l term, i.e. the function-value is objective $l = 1$.

This framework is linked to the field of **orthogonal polynomials** by the next proposition. We construct an optimal method w.r.t. a given distribution through a family of orthogonal polynomials associated to it.

Proposition 2 ((Pedregosa & Scieur, 2020)). Let P_t^l be defined as

$$P_t^l := \arg \min_{P_t(0)=1} \int P_t^2(\cdot)^l d(\cdot); \quad (6)$$

Then (P_t^l) is the family of residual orthogonal polynomials w.r.t. to $l+1$ d .

This theorem further implies that the optimal first-order method is a momentum method as Favard's theorem Marcellán & Álvarez-Nodarse (2001) tells us the orthogonal polynomials w.r.t. a given distribution are related through a **three term recurrence**,

$$P_{t+1}(\cdot) = a_t P_t(\cdot) + b_t \nabla P_t(\cdot) + (1 - a_t) P_{t-1}(\cdot); \quad (7)$$

Following Remark 1, the optimal method is derived from this recurrence as

$$\mathbf{x}_{t+1} = \mathbf{x}_t + (a_t - 1)(\mathbf{x}_t - \mathbf{x}_{t-1}) + b_t \nabla f(\mathbf{x}_t); \quad (8)$$

3 METHODS

Being able to write the rates in terms of the *expected spectral distribution* ties the average-case framework to the field of *random matrix theory*. Indeed, because of results from this field, certain e.s.d.'s are considered more natural than others. Indeed, it can be shown that the same distribution arises when we take the gram matrix of random centered i.i.d. features with variance σ^2 : the **Marchenko Pastur** distribution.

Definition 2 (MP distribution). *The Marchenko Pastur distribution associated with the parameter r and with scale σ^2 is given by*

$$d_{MP}(\lambda) = \frac{1}{2\sigma^2} \frac{\rho\left(\frac{\lambda + \sigma^2}{r}\right)\left(\frac{\lambda - \sigma^2}{r}\right)}{r}, \quad (9)$$

with $\rho = \sigma^2(1 + \rho\bar{r})^2$, $\bar{r} = \sigma^2 \max(0; (1 - \rho\bar{r})^2)$.

The Marchenko Pastur distribution d_{MP} can be considered a natural first model for e.s.d.'s as it arises universally from matrices with i.i.d. entries, under mild low moment assumptions, there is no specific distribution of the matrix to be considered. It can be seen as a model for the white-noise in the data. When $r = 1$, i.e. $n = d$, we have $d_{MP} \propto \frac{1}{\sqrt{1-\lambda}}$.

Pedregosa & Scieur (2020) first derived the optimal method w.r.t. d_{MP} , and Paquette et al. (2020) derived Nesterov's rates under the distribution. As we are concerned with being robust, a natural step is to consider the Beta weights.

Definition 3. *The (generalized) Beta weights with parameters α ; β and scale L are given by the (non-normalized) pdf*

$$d(\lambda) = (L - \lambda)^\alpha : \quad (10)$$

This family of distribution generalizes the MP distribution, and both have similar concentrations near 0 when $\alpha = 1/2$.

The optimal method w.r.t. d and objective l is associated to a shifted Jacobi polynomial $P_t^{\alpha, \beta}$ with $\alpha = \alpha + l + 1$; $\beta = \beta$. When $\alpha = \beta = 1/2$, we retrieve the *Chebyshev Method* (Flanders & Shortley, 1950). As such, we name our proposed methods the *Generalized Chebyshev Method* (GCM).

Algorithm 1: GCM(α ; β)

Inputs: Initial vector \mathbf{x}_0 , function f , smoothness parameter estimate L

$\mathbf{x}_1 = \mathbf{0}$; $\alpha = 0$

for $t = 1; \dots; T$ **do**

$$\left[\begin{array}{l} a_t = \frac{2(\alpha^2 + \beta(2t+1)(\alpha + \beta) + 2t^2 + 2)(2t + \alpha + \beta)}{2(t+1)(t + \alpha + \beta)(2t + \alpha)} \\ b_t = \frac{(2t + \alpha + \beta)(2t + \alpha + 2)}{L(t+1)(t + \alpha + \beta)} \\ t = \frac{(t + \alpha)(t + \beta)(2t + \alpha + 2)}{(t+1)(t + \alpha + \beta)(2t + \alpha)} \\ t = \frac{1}{a_t + t t - 1} \\ \mathbf{x}_t = \mathbf{x}_{t-1} + (t a_t - 1)(\mathbf{x}_{t-1} - \mathbf{x}_{t-2}) + t b_t r f(\mathbf{x}_{t-1}) \end{array} \right.$$

We'll consider the Nesterov's method used in Paquette et al. (2020), which is defined by the iterations:

$$\mathbf{x}_{t+1} = \mathbf{y}_t - \frac{1}{L} r f(\mathbf{y}_t) \quad (11)$$

$$\mathbf{y}_{t+1} = \mathbf{x}_{t+1} + \frac{t}{t+3} (\mathbf{x}_{t+1} - \mathbf{x}_t) \quad (12)$$

We also consider the Laguerre method, which is optimal w.r.t. $d(x) = \frac{x e^{-x}}{\Gamma(\alpha+1)}$, taking α as a parameter. This method is proposed to optimize non-smooth functions.

Both these methods are generalizations of ones that have been proposed in Pedregosa &

Scieur (2020). We show that Algorithm 1 corresponded to polynomials \mathcal{P}_t and derive the Laguerre method in appendix B.

Remark 2. *The Generalized Chebyshev takes the largest eigenvalue L as a parameter, but the rates we will show are robust to an overestimation of L .*

4 ROBUST AVERAGE-CASE RATES

We will state our assumption over the spectral distributions. It effectively allows us to parametrize all of our distributions of interest in a way that characterizes the asymptotic convergence, diving them into equivalence classes.

Assumption 1. *We will write \mathcal{D} for a continuous distribution supported in $(0; L]$ s.t. $d_{\mathcal{D}}(x) > 0$ for $x \geq [0; L]$, $d_{\mathcal{D}} = \mathcal{O}(x^{-\alpha})$ near 0 and $d_{\mathcal{D}} = \mathcal{O}((L-x)^{-\beta})$ near L .*

Assumption 1 is quite nonrestrictive, in that, the spectral distribution of the Hessian for any smooth convex problem can be identified with some \mathcal{D} in this class. It is a milder assumption than (1) assuming complete knowledge of the spectrum of the Hessian or (2) the specific distribution on the entries of your data. Moreover the assumption encompasses the frequently used MP (e.g. Martin & Mahoney (2021); Pennington & Bahri (2017)) and Uniform distributions We note there’s no need to consider eigenvalues situated at 0 as they do not contribute to the optimization process.

The β works as a measure of how close we are to the worst-case scenario, as it approaches 1. Samples in finite dimension of distributions with high values of β will work as strongly convex functions in practice.

We show that \mathcal{D} indeed behaves like an equivalence class when considering the asymptotics of the convergence of the methods: only the concentrations near the edge matter. We do this by singling out from each of these classes the beta distributions for which we can compute the rates, then show the rates to be the same inside \mathcal{D} .

Theorem 2 (GCM average-case rates). *A Generalized Chebyshev Method with parameters $(\alpha; \beta)$ applied to a problem with e.s.d. as in Assumption 1 has average-case rates*

$$E[f(\mathbf{x}_t) - f(\mathbf{x}^*)] \leq L C_1 \begin{cases} \mathcal{O}(t^{-1/2}) & \text{if } \alpha < \beta + 1/2 \text{ and } \beta < \beta + 3/2 \\ t^{-2(\alpha+2)} \log t & \text{if } \alpha = \beta + 1/2 \text{ and } \beta = \beta + 3/2 \\ t^{2(\max f)} & \text{if } \alpha > \beta + 1/2 \text{ or } \beta > \beta + 3/2 \end{cases} \quad (13)$$

$$E[\|j_t f(\mathbf{x}_t) - j_2^*\|^2] \leq L^2 C_2 \begin{cases} \mathcal{O}(t^{-1/2}) & \text{if } \alpha < \beta + 1/2 \text{ and } \beta < \beta + 5/2 \\ t^{-2(\alpha+3)} \log t & \text{if } \alpha = \beta + 1/2 \text{ and } \beta = \beta + 5/2 \\ t^{2g} & \text{if } \alpha > \beta + 1/2 \text{ or } \beta > \beta + 5/2 \end{cases} \quad (14)$$

where C_i is a distribution dependent constant.

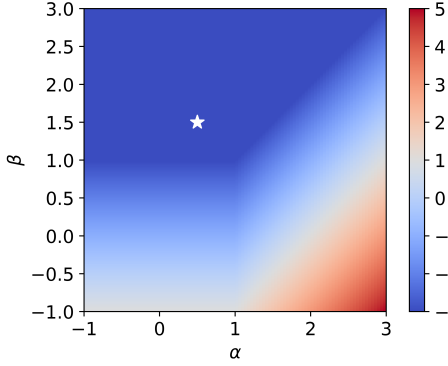
Theorem 2, which is illustrated by fig. 2 shows that overestimating α , and underestimating β will still leave us with the optimal asymptotic rates, so a good rule of thumb for calibrating the algorithm is to use high α and low β .

Theorem 3 shows that a proper choice of $(\alpha; \beta)$ can indeed make the Jacobi polynomial asymptotically optimal w.r.t. to any \mathcal{D} .

Theorem 3 (Optimal Rates). *Let \mathcal{D} follow Assumption 1. The optimal asymptotic average-case rates for $E[f(\mathbf{x}_t) - f(\mathbf{x}^*)]$ and $E[\|j_t f(\mathbf{x}_t) - j_2^*\|^2]$ are attained by the GCM with parameters $(\alpha; \beta + 2)$ and $(\alpha; \beta + 3)$, respectively, and read*

$$E[f(\mathbf{x}_t) - f(\mathbf{x}^*)] = \mathcal{O}(t^{-2(\alpha+2)}); \quad E[\|j_t f(\mathbf{x}_t) - j_2^*\|^2] = \mathcal{O}(t^{-2(\alpha+3)});$$

For the function value ($l = 1$), we find rates that approach t^{-2} as $\beta \rightarrow 1$, showing the worst-case as a limit (over the considered distribution) on the average-case.



Method	Parameters $(\beta; \alpha)$	
	$(\frac{1}{2}; \frac{1}{2})$	$(\frac{1}{2}; \frac{1}{2})$
GCM($\beta = \frac{1}{2}; \alpha = \frac{5}{2}$)	t^{-5}	t^{-3}
GCM($\beta = \frac{1}{2}; \alpha = \frac{3}{2}$)	t^{-4}	t^{-3}
Nesterov	t^{-4}	$t^{-3} \log t$
Gradient Descent	$t^{-\frac{5}{2}}$	$t^{-\frac{3}{2}}$

Figure 3 & Table 1: The figure illustrates the robustness of the Generalized Chebyshev Method with parameters $(\beta; \alpha)$ for a *fixed problem* corresponding to the Marchenko-Pastur distribution $(\beta = \frac{1}{2}; \alpha = \frac{1}{2})$. The color represents the exponent a of the average-case rate $O(t^a)$ of the method for different values of α and β . The white star represents the optimal tuning and the blue area is the set of parameters for which the method converges. Note we have a large region that guarantee the same optimal asymptotic rate. The table compares the asymptotic average-case rates for the function-value for different methods with different $(\beta; \alpha)$ values.

We remark that the above theorems imply that, at least asymptotically, the GCM is robust for a suboptimal choice of parameter α up to $1=2$ below the optimal choice and infinitely above.

For completeness, we also derive worst-case rates for the GCM:

Proposition 3 (GCM worst-case rates). *Let f be a convex, L -smooth quadratic function. Then, for the Generalized Chebyshev Method with parameters $(\beta; \alpha)$, we have worst-case rates*

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq C_1 L \begin{cases} t^{-2} & \text{if } \alpha > 1 \\ t^{-1/2} & \text{if } \alpha = 1 \\ t^{-1} & \text{if } \alpha < 1 \end{cases} \quad (15)$$

$$\|j_{\mathbf{x}} f(\mathbf{x}_t) - j_{\mathbf{x}} f(\mathbf{x}^*)\| \leq C_2 L^2 \begin{cases} t^{-2} & \text{if } \alpha > 2 \\ t^{-3/2} & \text{if } \alpha = 2 \\ t^{-2} & \text{if } \alpha < 2 \end{cases} \quad (16)$$

For a reasonable choice of $\beta; \alpha$, i.e. $\beta = \frac{1}{2}, \alpha = 1$, the function value achieves the theoretical lower bound of t^{-2} .

We now analyze the convergence of the Nesterov method. Nesterov (2003) has shown that it matches up to a constant factor a lower bound on the worst-case complexity of non strongly convex problems. A natural question is if this performance would translate to good average-case rates. To do so, we will extend Paquette et al. (2020) proof for the Nesterov method rates under the MP distribution.

Theorem 4 (Nesterov average-case rates). *Let f as in Assumption 1. Then for the Nesterov method, we have average-case rates*

$$E[f(\mathbf{x}_t) - f(\mathbf{x}^*)] \leq C_1^0 \begin{cases} t^{-2} & \text{if } \alpha < 1 \\ t^{-3} \log t & \text{if } \alpha = 1 \\ t^{-\alpha} & \text{if } \alpha > 1 \end{cases}; \quad E[\|j_{\mathbf{x}} f(\mathbf{x}_t) - j_{\mathbf{x}} f(\mathbf{x}^*)\|^2] \leq C_2^0 t^{-\alpha} \quad (17)$$

The difference between the asymptotic average-case rates of Nesterov and the optimal ones are $t^{-\alpha}$, when $\alpha > 1$, $\log t$ when $\alpha = 1$ and 0 otherwise. This shows that Nesterov is almost optimal when the concentrations near 0 are relatively high, i.e. low β .

Theorem 5 (Gradient Descent average-case rates). *Let f as in Assumption 1. Then for gradient descent*

$$E[f(\mathbf{x}_t) - f(\mathbf{x}^*)] = t^{-\alpha}; \quad E[\|j_{\mathbf{x}} f(\mathbf{x}_t) - j_{\mathbf{x}} f(\mathbf{x}^*)\|^2] = t^{-\alpha} \quad (18)$$

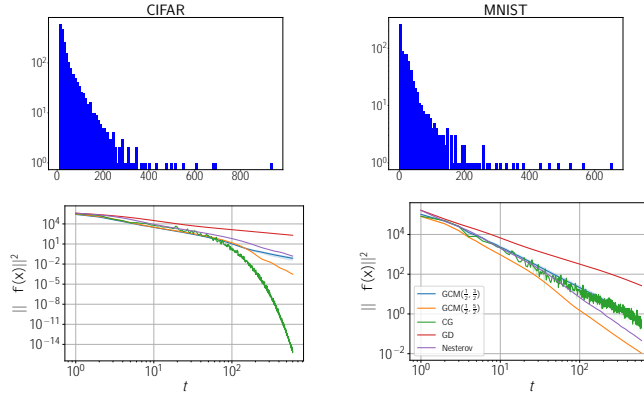


Figure 4: *Above:* Empirical spectrum for the covariance matrix of the features. *Below:* Gradient norms throughout iterations. *Left:* CIFAR-10 Inception features *Right:* MNIST features. Here we choose to compare gradient norms as the minimum function value is not known. The properly tuned GCM achieves remarkable performance under these non-synthetic spectrum’s.

Observe for the function value that the rate for Nesterov is t^{-2} and the rate for Gradient Descent is t^{-1} when $\beta \neq 1$.

Lastly, we consider the optimal rates for a Gamma distribution.

Theorem 6 (Laguerre method rates). *Let $\alpha > 1$ and $\beta > 1$ be a Gamma distribution, i.e. $d(x) = \frac{x^\alpha e^{-x}}{\Gamma(\alpha+1)}$. The optimal rates are given by the Laguerre method of appropriate tuning and*

$$\mathbb{E}[f(\mathbf{x}_t) - f(\mathbf{x}^*)] = O(t^{-(\alpha+2)}); \tag{19}$$

Note that this result does not have the same universality as the others because of the non-compactness of the distribution’s support. These rates are contrasted to the worst-case lower bound on the optimization of non-smooth functions by first-order methods, which gives

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{C}{k}.$$

These rates are not found when $\beta \neq 1$, indicating that the worst-case is especially pessimistic in this scenario.

Remark 3. *All of the expected rates we state are almost deterministic on the high dimensional setting as per the concentration results shown in Paquette et al. (2020)*

5 EXPERIMENTS

We simulate the e.s.d’s in two ways. The Marchenko Pastur distribution, which we sample by taking $\mathbf{H} = \mathbf{X}\mathbf{X}^T$ where \mathbf{X} has i.i.d. gaussian samples. This enables us to simulate $(\beta; \alpha)$ values of $(1=2; \beta=2)$. Other values of $(\beta; \alpha)$ are simulated by sampling $\lambda \in \mathbb{R}^d$ from the corresponding Beta distribution and taking $\mathbf{H} = \mathbf{U} \text{diag}(\lambda) \mathbf{U}^T$, where \mathbf{U} is an independently sampled orthonormal matrix.

We let $\mathbf{x}^* = 0$ and sample \mathbf{x}_0 from a centered Gaussian distribution, the dynamics are the same as in the general case. In all experiments we use the problem’s instance largest eigenvalue to calibrate each method (e.g. Gradient Descent’s stepsize is $1=L$). Our theoretical rates in Theorem 5 and Theorem 4 respectively for the Nesterov method and Gradient Descent are precise under the approximate range $1 < \beta < 2$ as we show in Figure 6. Distributions with higher β need many samples otherwise they behave as strongly convex functions.

The same is not true for the Generalized Chebyshev Method. If $\beta < 2$ or α is low the empirical findings diverge from the theoretical. We believe this is due to numerical instability

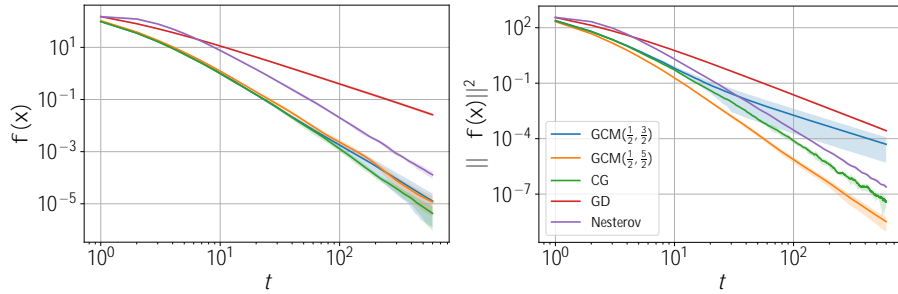


Figure 5: Rates for a synthetic problem, simulating the Marchenko Pastur distribution. Note that both tunings of the GCM achieve performance in function value very close to the one of Conjugate Gradient, which is optimal for every draw of the problem.

under these regimes as the metrics also have much larger variance than in the other regimes. We’ve not been able though to pinpoint the exact source of this supposed instability. This is shown in appendix D.

The GCM with $\alpha > \frac{1}{2}$ performs corresponding to the theory, and it’s non-asymptotically very close to the performance of $\frac{1}{2}$. High values of α also perform very well on non-synthetic data, suggesting in practice we should use these values.

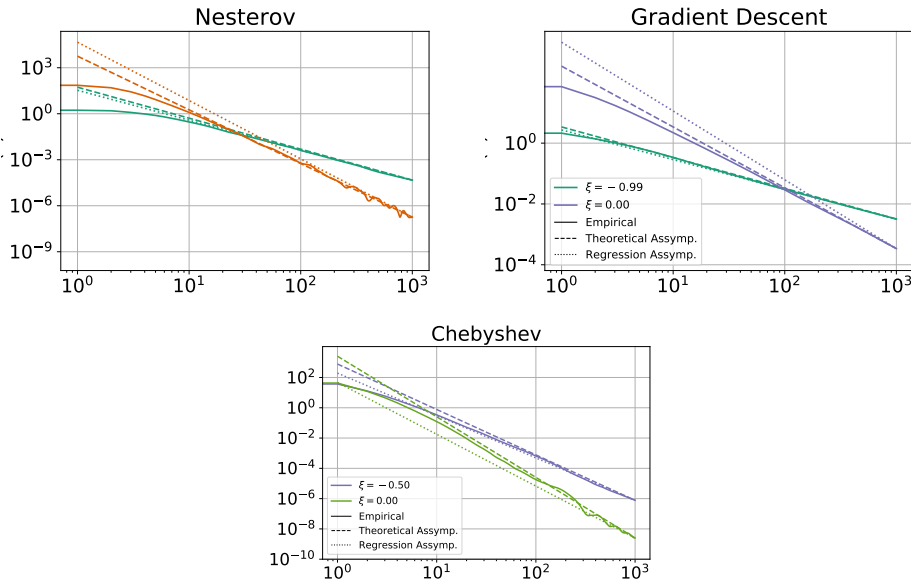


Figure 6: Comparison between experiments run on synthetic Beta distribution and theoretical asymptotic. Y-axis is the function value

6 CONCLUSION AND FURTHER WORK

In this paper, we’ve established that the asymptotic convergence of first order methods on quadratic problems in the convex regime depend on the concentration of the Hessian’s eigenvalue near the edges of the spectrum’s support. We further contributed to the theoretic understanding of the Nesterov’s method performance and established the contrast between the worst-case and average-case in the main regimes considered in Optimization.

BIBLIOGRAPHY

- Andrej Bogdanov and Luca Trevisan. Average-case complexity. *arXiv preprint cs/0606037*, 2006.
- Bernd Fischer. *Polynomial based iteration methods for symmetric linear systems*. SIAM, 1996.
- Donald A Flanders and George Shortley. Numerical determination of fundamental modes. *Journal of Applied Physics*, 21(12):1326–1332, 1950.
- Magnus Rudolph Hestenes, Eduard Stiefel, et al. *Methods of conjugate gradients for solving linear systems*, volume 49. NBS Washington, DC, 1952.
- Francisco Marcellán and Renato Álvarez-Nodarse. On the “favard theorem” and its extensions. *Journal of computational and applied mathematics*, 127(1-2):231–254, 2001.
- Charles H Martin and Michael W Mahoney. Implicit self-regularization in deep neural networks: Evidence from random matrix theory and implications for learning. *Journal of Machine Learning Research*, 22(165):1–73, 2021.
- Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2003.
- Courtney Paquette, Bart van Merriënboer, Elliot Paquette, and Fabian Pedregosa. Halting time is predictable for large models: A universality property and average-case analysis. *arXiv preprint arXiv:2006.04299*, 2020.
- Fabian Pedregosa and Damien Scieur. Acceleration through spectral density estimation. In *International Conference on Machine Learning*, pp. 7553–7562. PMLR, 2020.
- Jeffrey Pennington and Yasaman Bahri. Geometry of neural network loss surfaces via random matrix theory. In *International Conference on Machine Learning*, pp. 2798–2806. PMLR, 2017.
- B.T. Polyak. [Some methods of speeding up the convergence of iteration methods](#). *USSR Computational Mathematics and Mathematical Physics*, 04, 1964.
- Damien Scieur and Fabian Pedregosa. Universal average-case optimality of polyak momentum. *arXiv preprint arXiv:2002.04664*, 2020.
- Gabor Szegő. Orthogonal polynomials, vol. 23. In *American Mathematical Society Colloquium Publications*, 1975.
- Walter Van Assche. Weak convergence of orthogonal polynomials. *Indagationes Mathematicae*, 6(1):7–23, 1995.

A PROOFS OF SECTION 2

Theorem 1. Let \mathbf{x}_t be generated by a first-order method associated to the polynomial P_t , the measure the e.s.d. of H , and $\mathbb{E}[(\mathbf{x}_0 \ \mathbf{x}^\top)(\mathbf{x}_0 \ \mathbf{x}^\top)^\top] = R^2 I$ for some constant R . Then we can write the convergence metrics at time step t as

$$\begin{aligned} \mathbb{E}[k\mathbf{x}_t \ \mathbf{x}^\top k^2] &= R^2 \int P_t^2(\lambda) d(\lambda); & \mathbb{E}[f(\mathbf{x}_t) \ f(\mathbf{x}^\top)] &= R^2 \int P_t^2(\lambda) d(\lambda) \\ \text{and} \quad \mathbb{E}[j\|\nabla f(\mathbf{x}_t)\|_2^2] &= R^2 \int P_t^2(\lambda)^2 d(\lambda); \end{aligned} \quad (5)$$

Proof. We remark that by the definition of the expected spectral distribution of H , we have for continuous g

$$\mathbb{E}_H[g(\text{tr}(\mathbf{H}))] = \int g(\lambda) d(\lambda) \quad (20)$$

We know that $\mathbf{x}_t \ \mathbf{x}^\top = P_t(\mathbf{H})(\mathbf{x}_0 \ \mathbf{x}^\top)$. We can write $j\|\mathbf{x}_t \ \mathbf{x}^\top\|_2^2$ in terms of a trace and use the independence of \mathbf{H} and $\mathbf{x}_0 \ \mathbf{x}^\top$ to connect it to the e.s.d.:

$$\mathbb{E}[j\|\mathbf{x}_t \ \mathbf{x}^\top\|_2^2] = \mathbb{E}[\text{tr}((\mathbf{x}_0 \ \mathbf{x}^\top)^\top P_t(\mathbf{H})^2(\mathbf{x}_0 \ \mathbf{x}^\top))] \quad (21)$$

$$= \mathbb{E}_{H; \mathbf{x}_0 \ \mathbf{x}^\top}[\text{tr}(P_t(\mathbf{H})^2(\mathbf{x}_0 \ \mathbf{x}^\top)(\mathbf{x}_0 \ \mathbf{x}^\top)^\top)] \quad (22)$$

$$= \mathbb{E}_H [P_t(\mathbf{H})^2 \mathbb{E}_{\mathbf{x}_0 \ \mathbf{x}^\top}[(\mathbf{x}_0 \ \mathbf{x}^\top)(\mathbf{x}_0 \ \mathbf{x}^\top)^\top]] \quad (23)$$

$$= R^2 \mathbb{E}_H [P_t(\text{tr}(\mathbf{H}))^2] = R^2 \int P_t(\lambda)^2 d(\lambda) \quad (24)$$

For the gradient and function value the reasoning is the same by noticing that

$$\mathbb{E}[f(\mathbf{x}_t) \ f(\mathbf{x}^\top)] = \mathbb{E}[\text{tr}((\mathbf{x}_0 \ \mathbf{x}^\top)^\top P_t(\mathbf{H}) \mathbf{H} P_t(\mathbf{H})(\mathbf{x}_0 \ \mathbf{x}^\top))] \quad (25)$$

$$= \mathbb{E}_H [(P_t(\text{tr}(\mathbf{H})))^2]; \quad (26)$$

where P_t is also a polynomial. As $\nabla f(\mathbf{x}_t) = \mathbf{H}(\mathbf{x}_t \ \mathbf{x}^\top)$.

$$\mathbb{E}[j\|\nabla f(\mathbf{x}_t)\|_2^2] = \mathbb{E}[\text{tr}((\mathbf{x}_0 \ \mathbf{x}^\top)^\top P_t(\mathbf{H}) \mathbf{H}^2 P_t(\mathbf{H})(\mathbf{x}_0 \ \mathbf{x}^\top))] \quad (27)$$

$$= \mathbb{E}_H [(P_t(\text{tr}(\mathbf{H})))^2] \quad (28)$$

□

Proposition 2 ((Pedregosa & Scieur, 2020)). Let P_t^l be defined as

$$P_t^l := \arg \min_{P_t(0)=1} \int P_t^2(\lambda)^l d(\lambda); \quad (6)$$

Then (P_t^l) is the family of residual orthogonal polynomials w.r.t. to $\lambda^{l+1} d$.

Proof. We differentiate the expression for the metrics w.r.t. to the coefficients of the polynomials:

$$\begin{aligned} \frac{d}{da_k} \int P_t^2(\lambda) d(\lambda) &= \int \frac{d}{da_k} \sum_{k=0}^l a_k^k P_t(\lambda) d(\lambda) = \\ &= 2 \int \lambda^{l+k} P_t(\lambda) d(\lambda) = 0 \end{aligned}$$

This means that $P_t(\lambda)$ is orthogonal to any polynomial of degree $t-1$ w.r.t to the inner product $\int \lambda^{l+1} d$

□

B GCM AND LAGUERRE METHOD DERIVATION

We will first state two lemmas that allow us to construct the optimal polynomials. With them in hand the procedure is trivial.

Lemma 7. *Let (P_t) be a family polynomials following*

$$P_t(x) = (t + x)P_{t-1}(x) + tP_{t-2}(x);$$

with P_0 a constant polynomial and $P_t \notin 0; \forall t$. Then

$$P_t(x) = (a_t + b_t)P_{t-1}(x) + (1 - a_t)P_{t-2}(x) \quad (29)$$

is the recurrence for $P_t(x) = P_t(x) = P_t(0)$. With:

$$a_t = t - t \quad (30)$$

$$b_t = t - t \quad (31)$$

$$t = (t + t - t - 1) \quad (t_0 = 0) \quad (32)$$

The proof of this is presented in [Pedregosa & Scieur \(2020\)](#). Further, we know how to compute the recurrence for the polynomials of a shifted distribution:

Lemma 8. *Let (P_t) be a family polynomials orthogonal w.r.t following*

$$P_t(x) = (t + x)P_{t-1}(x) + tP_{t-2}(x); \quad (33)$$

and define polynomials P_t s.t. :

$$P_t(m(x)) = P_t(x);$$

with $m(x) = a + bx$ a non singular affine transform. Then P_t follows a recurrence like in eq. (33), with:

$$\frac{0}{t} = t + b - t \quad (34)$$

$$\frac{0}{t} = a - t \quad (35)$$

$$\frac{0}{t} = t \quad (36)$$

The lemma is self-evident by considering eq. (33) with argument $m^{-1}(x)$

These results are enough to get the recurrence relation for the residual polynomial w.r.t $x(L - x)$. We begin by the standard Jacobi polynomials, which are orthogonal w.r.t $(1 - x)(1 + x)$ and follow a recurrence according to $t; t; t$ below [Szegő \(1975\)](#):

$$t = \frac{(2n + t + 1)(2n + t - 1)}{2n(n + t)} \quad (37)$$

$$t = \frac{(t^2 - 2t)(2n + t - 1)}{2n(n + t)(2n + t - 2)} \quad (38)$$

$$t = \frac{2(n + 1)(n + 1)(2n + t)}{2n(n + t)(2n + t - 2)} \quad (39)$$

We then shift the distribution according to (x) , and then transform to the residual ones. We slightly simplify these computations and use remark 1 to get Algorithm 1.

We know ([Szegő, 1975](#)) that the Laguerre polynomials L_t , with usual normalization, follow the recurrence

$$L_t(x) = \frac{2t + 1}{t} \frac{1}{t} L_{t-1}(x) + \frac{t + 1}{t} L_{t-2}(x) \quad (40)$$

As we don't have to shift the domains, we have only to apply lemma 7 to get the Laguerre method. Further, we can get an explicit expression for $t = \frac{t}{t+1}$, simplifying the expression.

Algorithm 2: Laguerre()**Inputs:** Initial vector \mathbf{x}_0 , function f $\mathbf{x}_1 = \mathbf{0}$ **for** $t = 1; \dots; T$ **do** $\mathbf{x}_t = \mathbf{x}_{t-1} + \frac{t-1}{t}(\mathbf{x}_{t-1} - \mathbf{x}_{t-2}) - \frac{1}{t}r f(\mathbf{x}_{t-1})$

C PROOFS OF SECTION 3

In the following we will consider shifted versions of the spectral distributions. This shift is written as an affine transform $m(\cdot) : [0; L] \rightarrow [1; 1]$ because most results in the theory of orthogonal polynomials are stated in terms of distributions supported in $[1; 1]$.

This can be seen as an additional layer of abstraction because the quantities evaluated with the shifted distributions and polynomials are proportional, i.e. if $P_t(x) = \mathcal{P}_t(m(x))$ and $\theta(x) = \tilde{\theta}(m(x))$:

$$\int P_t^2(x) \theta(x) dx \propto \int \mathcal{P}_t^2(x) \tilde{\theta}(x) dx \quad (41)$$

So all the asymptotics are the same. The Jacobi polynomials $P_t^{(\alpha, \beta)}$ are those orthogonal w.r.t $d(x) = (1-x)^\alpha (1+x)^\beta$. Most works use the normalization $\mathcal{P}_t^{(\alpha, \beta)}(1) = (t!)^{-\alpha-\beta}$. We will write $\mathcal{P}_t^{(\alpha, \beta)}$ for this normalization and $P_t^{(\alpha, \beta)}$ for the residual polynomials

Theorem 2 (GCM average-case rates). *A Generalized Chebyshev Method with parameters (α, β) applied to a problem with e.s.d. as in Assumption 1 has average-case rates*

$$E[f(\mathbf{x}_t) - f(\mathbf{x}^*)] \leq L C_1^{(\alpha, \beta)} \begin{cases} t^{-1/2} & \text{if } \alpha < -1/2 \text{ and } \beta < -3/2 \\ t^{-2(\alpha+\beta)} \log t & \text{if } \alpha = -1/2 \text{ and } \beta = -3/2 \\ t^{2(\max\{\alpha, \beta\})} & \text{if } \alpha > -1/2 \text{ or } \beta > -3/2 \end{cases} \quad (13)$$

$$E[\|Jr f(\mathbf{x}_t)\|_2^2] \leq L^2 C_2^{(\alpha, \beta)} \begin{cases} t^{-1/2} & \text{if } \alpha < -1/2 \text{ and } \beta < -5/2 \\ t^{-2(\alpha+\beta)} \log t & \text{if } \alpha = -1/2 \text{ and } \beta = -5/2 \\ t^{2(\max\{\alpha, \beta\})} & \text{if } \alpha > -1/2 \text{ or } \beta > -5/2 \end{cases} \quad (14)$$

where $C_i^{(\alpha, \beta)}$ is a distribution dependent constant.

Proof. We will prove that for any $\alpha, \beta > -1$, $l > 0$ and following Assumption 1, we have

$$\int P_t^{(\alpha, \beta)}(x)^2 x^l d_{(\alpha, \beta)}(x) \leq L^l C^{(\alpha, \beta)} \begin{cases} t^{-1/2} & \text{if } \alpha < -1/2 \text{ and } \beta < -1/2 \\ t^{-2(\alpha+\beta)} \log t & \text{if } \alpha = -1/2 \text{ and } \beta = -1/2 \\ t^{2(\max\{\alpha, \beta\})} & \text{if } \alpha > -1/2 \text{ or } \beta > -1/2 \end{cases}$$

We will first show this result for the Beta weights, then show that distributions with the same concentration behave similarly.

The normalization of $\mathcal{P}_t^{(\alpha, \beta)}$ is s.t. [Szegő (1975) (4.3.3)]:

$$\int_{-1}^1 \mathcal{P}_t^{(\alpha, \beta)}(x) (1-x)^\alpha (1+x)^\beta dx = \frac{2^{\alpha+\beta+1}}{2n+\alpha+\beta+1} \frac{(n+\alpha+1)(n+\beta+1)}{(n+1)(n+\alpha+\beta+1)} = (t^{-1}) \quad (42)$$

Further, the residual polynomials are s.t. $J P_t^{(\alpha, \beta)} J = (t!) J \mathcal{P}_t^{(\alpha, \beta)} J$, from the definition of the classical normalization.

We state the result (Exercise 91, Generalisation of 7.34.1) from Szegő (1975):

Lemma 9. *We have*

$$\int_0^1 (1-x) P_t^{(j)}(x)^2 dx \quad (h) \quad (43)$$

$$h := \begin{cases} < t^{-1} \log t & \text{if } j > +1=2 \\ t^{-1} & \text{if } j = +1=2 \\ t^{-1} & \text{if } j < +1=2 \end{cases} \quad (44)$$

Noting that $P_t^{(j)}(x) = (1-x)^t P_t^{(j)}(1-x)$, we can write:

$$\int_0^1 P_t(x)^2 (1-x)(1+x) dx = \int_0^1 (1-x) j P_t^{(j)}(x)^2 dx + \int_0^1 (1-x) j P_t^{(j)}(x)^2 dx \quad (45)$$

We can then show our result for $d_{t,j}(x) = x^{-j}(1-x)$ by carefully considering each of the cases on Lemma 9 and the maximum of each term in eq. 45, and an added t^{-2} from the different normalization. With this, we have the wanted result for the Beta weights

It remains to show:

$$\int_0^1 P_t^{(j)}(x)^2 d_{t,j}(x) = \int_0^1 (1-x) P_t^{(j)}(x)^2 dx \quad (46)$$

And the rest follows from the same arguments. We do this with the help of this lemma shown in Van Assche (1995) relating to the weak convergence of the orthogonal polynomials:

Lemma 10. *Let μ be a measure and (p_t) it's family of orthonormal polynomials s.t. p_t follow the recurrence:*

$$x p_t(x) = a_t p_{t+1}(x) + b_t p_t(x) + a_{t-1} p_{t-1}(x)$$

and a_t, b_t converge respectively to a, b . Then for any f continuous and bounded:

$$\int f(x) p_t^2(x) d\mu(x) \rightarrow \int \frac{f(x)}{1-x^2} dx \quad (47)$$

Let s.t.

$$x(1-x)^j d_{t,j} = A(1-x)^j + B(1-x) \quad (48)$$

We observe that for $0 < x < 1$, $f(x) = \frac{d_{t,j}}{(1-x)(1+x)}$ is bounded.

We get from an application of 10, and the observation that $P_t^{(j)} = N_t p_t^{(j)}$, with $N_t = (t^{-1=2})$:

$$\int_0^1 (1-x) P_t^{(j)}(x)^2 dx = \int_0^1 (1-x) P_t^{(j)}(x)^2 dx + \int_0^1 (1-x) P_t^{(j)}(x)^2 dx \quad (49)$$

$$\int_0^1 (1-x) P_t^{(j)}(x)^2 dx = (h) \quad (50)$$

$$\int_0^1 P_t^{(j)}(x)^2 d_{t,j}(x) = \int_0^1 P_t^{(j)}(x)^2 f(x) (1-x)(1+x) dx + \int_0^1 (1-x) P_t^{(j)}(x)^2 dx \quad (51)$$

□

Theorem 3 (Optimal Rates). *Let f follow Assumption 1. The optimal asymptotic average-case rates for $\mathbb{E}[f(\mathbf{x}_t) - f(\mathbf{x}^*)]$ and $\mathbb{E}[|j_t f(\mathbf{x}_t)|^2]$ are attained by the GCM with parameters $(\gamma; \nu + 2)$ and $(\gamma; \nu + 3)$, respectively, and read*

$$\mathbb{E}[f(\mathbf{x}_t) - f(\mathbf{x}^*)] = O(t^{-2(\nu+2)}); \quad \mathbb{E}[|j_t f(\mathbf{x}_t)|^2] = O(t^{-2(\nu+3)}).$$

Proof. We will prove that for $\gamma > 1$ If $\nu = \gamma - 1$ and $\nu = \gamma + l + 1$ (i.e., γ are optimal), the rate of convergence reads

$$\min_{P_t(0)=1} \int_{-1}^1 P_t^2(x) d(x) = \int_{-1}^1 P_t^2(x) d(x) = O(t^{-2(\nu+l+1)}) \quad (52)$$

Showing the second equality is easy by considering theorem 2, and that is further the minimum asymptotic rate for the Beta distribution $\beta(\gamma, \gamma)$.

By setting ρ_t and $P_t = \frac{P_t}{P_t(0)}$ the optimal orthonormal and residual and polynomials w.r.t. we show that P_t must have the same rate on $\beta(\gamma, \gamma)$ as it does on $\beta(\gamma, \gamma)$, thus the optimal rate of cannot be lower than the optimal rate of $\beta(\gamma, \gamma)$. Indeed, setting $\gamma_1; \gamma_2$ as in eq. 48:

$$\int_{-1}^1 P_t(x)^2 d(x) = \int_{-1}^1 P_t(x)^2 d(\gamma_1; \gamma_2)(x) \quad (53)$$

$$\int_{-1}^1 P_t(x)^2 d(x) = \int_{-1}^1 P_t(x)^2 d(\gamma_1; \gamma_2)(x) \quad (54)$$

$$\int_{-1}^1 P_t(x)^2 d(x) = \int_{-1}^1 P_t(x)^2 d(\gamma_1; \gamma_2)(x) = \frac{1}{\rho_t(\gamma_1)^2} \quad (55)$$

$$(56)$$

Where the first two equations come from the fact that $\beta(\gamma_1; \gamma_2)$ near $\gamma_1 = 1$ and $\gamma_2 = 1$ and the third from lemma 10.

This effectively upper bounds the rates on $\beta(\gamma, \gamma)$ because the rates of P_t on $\beta(\gamma_1; \gamma_2)$ can't be lower than $O(t^{-2(\nu+l+1)})$. \square

Proposition 3 (GCM worst-case rates). *Let f be a convex, L -smooth quadratic function. Then, for the Generalized Chebyshev Method with parameters $(\gamma; \nu)$, we have worst-case rates*

$$\mathbb{E}[f(\mathbf{x}_t) - f(\mathbf{x}^*)] \leq C_1 L \begin{cases} t^{-1/2} & \text{if } \nu > 1 \\ t^{-2} & \text{if } \nu = 1 \\ t^{-1/2} & \text{if } \nu < 1 \end{cases} \quad (15)$$

$$\mathbb{E}[|j_t f(\mathbf{x}_t) - j_t f(\mathbf{x}^*)|^2] \leq C_2 L^2 \begin{cases} t^{-1/2} & \text{if } \nu > 2 \\ t^{-4} & \text{if } \nu = 2 \\ t^{-1/2} & \text{if } \nu < 2 \end{cases} \quad (16)$$

Proof. rates] We will prove that: $\sup_{x \in [0; L]} x^l P_t^2(x) = O(L^l t^{\nu(\gamma; \nu)})$. Where:

$$\nu(\gamma; \nu) = \begin{cases} 2(\nu) & \text{if } \nu > l \\ 1 - 2\nu & \text{if } \nu = l \\ 2l & \text{if } \nu < l \end{cases} \quad (57)$$

From Szegő (1975), Theorem 7.32.2, if $\nu < \frac{1}{2}$:

$$P_t^2(\cos \theta) = \begin{cases} O(t^{-1-2\nu}) & \text{if } \nu < \frac{1}{2} \\ O(t) & \text{if } \nu = \frac{1}{2} \\ O(t^{-1-2\nu}) & \text{if } \nu > \frac{1}{2} \end{cases} \quad (58)$$

We observe that, from the symmetry of the Jacobi polynomials:

$$\sup_{x \in [0; L]} x^l P_t^2(x) = \max_{x \in [0; 1]} \sup_{x \in [0; 1]} P_t^2(x); \sup_{x \in [0; 1]} (1-x)^l P_t^2(x) \quad (59)$$

The $(1-x)^l$ term, corresponds to $(2 \sin(\frac{\theta}{2}))^l$ in the variable θ , which is $O(t^{-2l})$. The rest follows from carefully considering the expressions given by eq. 58. \square

Theorem 4 (Nesterov average-case rates). *Let as in Assumption 1. Then for the Nesterov method, we have average-case rates*

$$E[f(\mathbf{x}_t) - f(\mathbf{x}^*)] \leq C_1^0 \begin{cases} t^{-2} & \text{if } \alpha < 1/2 \\ t^{-3} \log t & \text{if } \alpha = 1/2 \\ t^{-(\alpha+2)} & \text{if } \alpha > 1/2 \end{cases}; \quad E[\|J_t f(\mathbf{x}_t)\|_2^2] \leq C_2^0 t^{-(\alpha+2)}. \quad (17)$$

Proof. We will prove:

$$\int_0^1 P_t^{\text{Nes}}(u)^2 du \leq C^0 \begin{cases} t^{-2} & \text{if } 0 < \alpha < 1/2 \\ t^{-3} \log t & \text{if } \alpha = 1/2 \\ t^{-(\alpha+2)} & \text{if } \alpha > 1/2 \end{cases} \quad (60)$$

Paquette et al. (2020) has shown that the nesterov polynomials P_t are asymptotically, in t :

$$P_t(u) \sim \frac{2J_1(t^{\alpha} u)}{t^{\alpha}} e^{-ut} \quad (61)$$

In the sense that:

$$\int_0^1 u^l P_t^2(u) du = \frac{4J_1^2(t^{\alpha} u)}{t^{2\alpha}} e^{-ut} du + o(t^{-(l+25=12)}) \quad (62)$$

The arguments can be easily used to show that such an integral is $O(t^{-(l+31=12)})$ when evaluated wrt a general d s.t. $l = d$ near 0.

We can thus consider our integral of interest substituting P_t^{Nes} by its Bessel asymptotic and dividing it into three regions, i.e. $[0; 1] = [0; \frac{1}{t^{\alpha}}] \cup [\frac{1}{t^{\alpha}}; \frac{1}{t^{\alpha}}] \cup [\frac{1}{t^{\alpha}}; 1]$ corresponding to two different regimes for the Bessel function. The first region will give us the asymptotic and the others we will bound.

We consider first, for some $\epsilon > 0$:

$$\int_{\frac{1}{t^{\alpha}}}^1 u \frac{4J_1^2(t^{\alpha} u)}{t^{2\alpha}} e^{-ut} du \quad (63)$$

We note the asymptotic for J_1^2 :

$$J_1^2\left(\frac{v}{t^{\alpha}}\right) \sim \frac{1}{t^{\alpha}} (1 + \cos(2\frac{v}{t^{\alpha}} + 2)) \quad (64)$$

Doing the change of variable $v = tu$, and identifying the upper limit of the interval, which is t^{1-2} to 1:

$$\int_{\frac{1}{t^{\alpha}}}^1 u \frac{4J_1^2(t^{\alpha} u)}{t^{2\alpha}} e^{-ut} du = t^{-2} \int_{v=1}^1 v^{-1} J_1^2\left(\frac{v}{t^{\alpha}}\right) e^{-v} dv \quad (65)$$

$$= t^{-2} \int_{v=1}^1 v^{-1} \frac{1}{t^{\alpha}} e^{-v} dv + o(t^{-2}) \quad (66)$$

$$= \frac{1}{t^{\alpha+2}} \int_{v=1}^1 v^{-\frac{3}{2}} \frac{1}{t^{\alpha}} e^{-v} dv + o(t^{-\frac{5}{2}}) \quad (67)$$

Where the cosine term goes to 0 from the Riemann-Lebesgue lemma and Γ is the incomplete Gamma function.

The term corresponding to the interval $[\frac{1}{t^{\alpha}}; 1]$ is exponentially small. Indeed, because of the exponential e^{-ut} it is $O(e^{-\frac{1}{t^{\alpha}}})$. This shows that the integral concentrates in a region that is closer and closer to 0 and that only the behaviour of the distribution near 0 matters. We have for the $[0; \frac{1}{t^{\alpha}}]$ region, doing the change of variables $v = t^{\alpha} u$:

$$\int_0^{\frac{1}{t^{\alpha}}} u \frac{4J_1^2(t^{\alpha} u)}{t^{2\alpha}} e^{-ut} du = t^{-2(\alpha+1)} \int_0^1 v \frac{J_1^2(v)}{v} e^{-\frac{v}{t^{\alpha}}} dv \quad (68)$$

And the $e^{-\frac{v}{t}}$ is (1). We have the following Bessel asymptotics:

$$\frac{J_1^2(\rho\sqrt{v})}{v} \sim \frac{1}{4}; \quad v \neq 0 \quad (69)$$

$$\frac{J_1^2(\rho\sqrt{v})}{v} = O(v^{-3/2}); \quad v \neq 1 \quad (70)$$

So we divide this integral aswell:

$$t^{-2(\alpha+1)} \int_0^t \frac{J_1^2(\rho\sqrt{v})}{v} e^{-\frac{v}{t}} dv = t^{-2(\alpha+1)} \int_0^t v^{-\frac{3}{2}} dv = I(t) t^{-\frac{5}{2}} \quad (71)$$

$$t^{-2(\alpha+1)} \int_0^1 \frac{J_1^2(\rho\sqrt{v})}{v} e^{-\frac{v}{t}} dv = t^{-2(\alpha+1)} \int_0^1 v^{-1} dv = t^{-2(\alpha+1)} \quad (72)$$

Where $I(t) = \log t$ if $\alpha = \frac{1}{2}$ and 1 otherwise.

The nesterov rate is then $I(t) t^{-\frac{5}{2}}$ if $\alpha = \frac{1}{2}$ and $t^{-2(\alpha+1)}$ if $0 < \alpha < \frac{1}{2}$ \square

Theorem 5 (Gradient Descent average-case rates). *Let ρ as in Assumption 1. Then for gradient descent*

$$\mathbb{E}[f(\mathbf{x}_t) - f(\mathbf{x}^*)] = O(t^{-(\alpha+2)}); \quad \mathbb{E}[\|j_t f(\mathbf{x}_t)\|_2^2] = O(t^{-(\alpha+3)}); \quad (18)$$

Proof. Considering that $P_t^{\text{GD}}(\cdot) = (1 - \frac{\cdot}{t})^t$ we will prove :

$$\int_0^1 (1 - \frac{v}{t})^{2t} v^{-l} dv; \quad l = (\alpha + l + 1) \quad (73)$$

We know, for the Beta weights, that:

$$\int_0^1 (1 - \frac{v}{t})^{2t+\alpha} v^{-l} dv = \frac{(l + \alpha + 1)(2t + \alpha + 1)}{(2t + l + \alpha + 2)} = O(t^{-(\alpha+l+1)}) \quad (74)$$

We can identify this asymptotic to the interval \mathbb{R}_0^+ for any α because:

$$\int_0^1 (1 - \frac{v}{t})^{2t+\alpha} v^{-l} dv = O((1 - \frac{v}{t})^{2t}) \quad (75)$$

Then:

$$\int_0^1 (1 - \frac{v}{t})^{2t-l} dv; \quad l = O((1 - \frac{v}{t})^{2t}) \quad (76)$$

$$\int_0^1 (1 - \frac{v}{t})^{2t-l} dv; \quad l = \int_0^1 (1 - \frac{v}{t})^{2t+\alpha} v^{-l} dv = O(t^{-(\alpha+l+1)}) \quad (77)$$

\square

Theorem 6 (Laguerre method rates). *Let $\alpha > 1$ and $\beta > 1$ be a Gamma distribution, i.e. $d(x) = \frac{x^\alpha e^{-x}}{\Gamma(\alpha+1)}$. The optimal rates are given by the Laguerre method of appropriate tuning and*

$$\mathbb{E}[f(\mathbf{x}_t) - f(\mathbf{x}^*)] = O(t^{-(\alpha+2)}); \quad (19)$$

Proof. Let L_t be the Laguerre polynomials with the usual normalization Szegő (1975):

$$\int_0^\infty L_t(x)^2 d(x) = L_t(0) = \binom{n+\alpha}{n} \quad (78)$$

We further now [Szegő (1975) (5.1.13)]:

$$\sum_{k=0}^t L_t(x) = L_{t+1}(x) \quad (79)$$

Thus, letting P_t be the residual laguerre polynomial, we consider:

$$\begin{aligned}
 E[f(\mathbf{x}_t) - f(\mathbf{x}^*)] &= \int_0^1 P_t^{+2}(\xi)^2 d_{+1}(\xi) = \frac{t+2}{t} \int_0^1 L_t^{+2} d_{+1}(\xi) \\
 &= \frac{t+2}{t} \sum_{k=0}^t L_k^{+1}(\xi) d_{+1}(\xi) \\
 &= \frac{t+2}{t} \sum_{k=0}^t \frac{k+1}{k} = \frac{t+2}{t} \sum_{k=0}^t \frac{t+2}{t} \\
 &= \frac{t+2}{t} \sum_{k=0}^t 1 = (t+2)
 \end{aligned} \tag{80}$$

□

D ADDITIONAL EXPERIMENTS

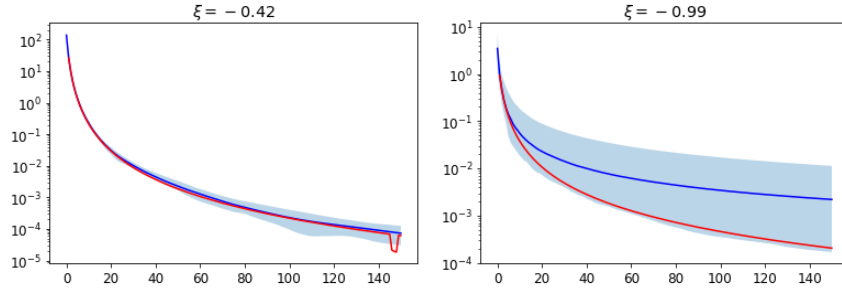


Figure 7: Empirical vs Theoretical function-value performance for $GCM(\xi; ?)$. Red lines are given by numerical integration, shades are minimum and maximum values under 10 runs and the blue line is the mean

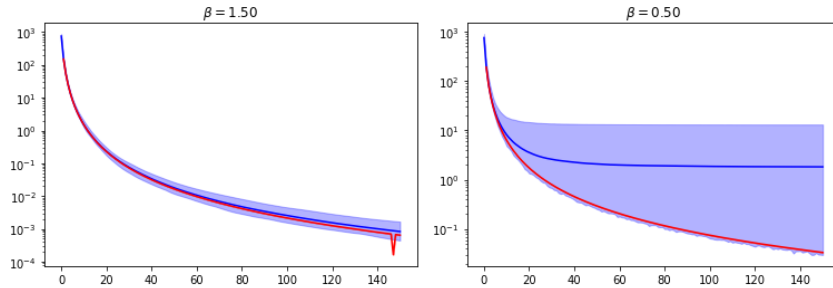


Figure 8: Empirical vs Theoretical function-value performance under Marchenko Pastur distribution. Red lines are given by numerical integration, shades are minimum and maximum values under 10 runs and the blue line is the mean

We note that in the regimes where the empirical average performance doesn't match the theoretical one, we can still find samples of problems who do match. This and the much larger variance on the function-value, this discrepancy is due to numerical instability in these regimes.