

---

# CORNER GRADIENT DESCENT

**Dmitry Yarotsky**

Skoltech

d.yarotsky@skoltech.ru

## ABSTRACT

We consider SGD-type optimization on infinite-dimensional quadratic problems with power law spectral conditions. It is well-known that on such problems deterministic GD has loss convergence rates  $L_t = O(t^{-\zeta})$ , which can be improved to  $L_t = O(t^{-2\zeta})$  by using Heavy Ball with a non-stationary Jacobi-based schedule (and the latter rate is optimal among fixed schedules). However, in the mini-batch Stochastic GD setting, the sampling noise causes the Jacobi HB to diverge; accordingly no  $O(t^{-2\zeta})$  algorithm is known. In this paper we show that rates up to  $O(t^{-2\zeta})$  can be achieved by a generalized stationary SGD with infinite memory. We start by identifying generalized (S)GD algorithms with contours in the complex plane. We then show that contours that have a corner with external angle  $\theta\pi$  accelerate the plain GD rate  $O(t^{-\zeta})$  to  $O(t^{-\theta\zeta})$ . For deterministic GD, increasing  $\theta$  allows to achieve rates arbitrarily close to  $O(t^{-2\zeta})$ . However, in Stochastic GD, increasing  $\theta$  also amplifies the sampling noise, so in general  $\theta$  needs to be optimized by balancing the acceleration and noise effects. We prove that the optimal rate is given by  $\theta_{\max} = \min(2, \nu, \frac{2}{\zeta+1/\nu})$ , where  $\nu, \zeta$  are exponents appearing in the capacity and source spectral conditions. Furthermore, using fast rational approximations of the power functions, we show that ideal corner algorithms can be efficiently approximated by finite-memory algorithms, and demonstrate their practical efficiency on a synthetic problem and MNIST.

## 1 INTRODUCTION

It is well-known that Gradient Descent (GD) on quadratic problems can be accelerated using the additional momentum term (the “Heavy Ball” algorithm, Polyak (1964)). For ill-conditioned problem, Heavy Ball with a suitable non-stationary (“Jacobi”) predefined schedule allows to accelerate a power-law loss converge rate  $O(t^{-\zeta})$  to  $O(t^{-2\zeta})$ , i.e. double the exponent  $\zeta$  (Nemirovskiy & Polyak, 1984; Brakhage, 1987). This acceleration is the best possible for non-adaptive schedules.

On the other hand, for mini-batch *Stochastic* Gradient Descent (SGD) typically used in modern machine learning, the convergence rate picture is much more complicated, and much less is known about possible acceleration. The natural quadratic problem in this case is the fitting of a linear model with a sampled quadratic loss. In the power-law spectral setting, it was found in Berthier et al. (2020) that plain SGD has two distinct convergent phases: either the sampling noise is weak and the SGD rate is the same  $O(t^{-\zeta})$  as for GD, or the convergence is slower due to the prevalence of the sampling noise. We refer to these two scenarios as *signal-* and *noise-dominated*, respectively.

This picture was refined in several other works (Varre et al., 2021; Varre & Flammarion, 2022; Velikanov et al., 2023; Yarotsky & Velikanov, 2024). In particular, Yarotsky & Velikanov (2024) examined generalized SGDs with finite linear memory of any size (generalizing the momentum and similar terms) and proved that with stationary schedules they all have the same phase diagram as plain SGD (Figure 2 left); in particular, they do not accelerate the plain GD/SGD rate  $O(t^{-\zeta})$ .

On the other hand, the non-stationary Jacobi Heavy Ball accelerating deterministic GD from  $O(t^{-\zeta})$  to  $O(t^{-2\zeta})$  fails for mini-batch Stochastic GD: it eventually starts to diverge due to the accumulating sampling noise. Varre & Flammarion (2022) have proposed a non-stationary modification of SGD that achieves a quadratic acceleration, but only on finite-dimensional problems. Yarotsky & Velikanov (2024) have proposed a non-stationary modification of the Heavy Ball/momentum algorithm

that is heuristically expected (but not yet proved) to achieve rates  $O(t^{-\theta\zeta})$  with some  $1 < \theta < 2$  on infinite-dimensional problems.

To sum up, the topic of SGD acceleration in ill-conditioned quadratic problems is far from settled.

In the present paper we propose an entirely new approach to acceleration of (S)GD that both provides a new general geometric viewpoint and proves that, in a certain rigorous sense, SGD in the signal-dominated regime can be accelerated from  $O(t^{-\zeta})$  to  $O(t^{-\theta\zeta})$  with  $\theta$  up to 2.

**Our contributions:**

1. **A view of generalized (S)GD as contours** (Section 3). We show that stationary (S)GD algorithms with an arbitrary-sized linear memory can be identified with contours in the complex plane. This identification leverages the characteristic polynomials  $\chi$  and the loss expansions of memory- $M$  (S)GD from Yarotsky & Velikanov (2024). We show that all the information needed to compute the loss evolution is contained in a map  $\Psi : \{z \in \mathbb{C} : |z| \geq 1\} \rightarrow \mathbb{C}$  associated with  $\chi$ . The map  $\Psi$  gives rise to the contour  $\Psi(\{z \in \mathbb{C} : |z| = 1\})$  and, conversely, can be reconstructed, along with the algorithm, from a given contour.
2. **Corner algorithms** (Section 4). A crucial role is played by contours that have a corner with external angle  $\theta\pi, 1 < \theta < 2$ . We prove that the respective algorithms accelerate the plain GD rate  $O(t^{-\zeta})$  to  $O(t^{-\theta\zeta})$ . However, in Stochastic GD such algorithms have the negative effect of amplifying the sampling noise. By balancing these two effects, we establish the precise phase diagram of feasible accelerations of SGD under power-law spectral assumptions (Figure 1 right). In particular, we identify three natural sub-phases in the signal-dominated phase; in one of them acceleration up to  $O(t^{-2\zeta})$  is theoretically feasible.
3. **Implementation of Corner (S)GD** (Sections 5, 6). Ideal corner algorithms require an infinite memory, but can be fast approximated by finite-memory algorithms using fast rational approximations of the power function  $z^\theta$ . We demonstrate experimentally an acceleration close to theoretically predicted by applying a memory-5 approximation of Corner SGD to a synthetic problem and the MNIST classification.

## 2 BACKGROUND

This section is largely based on the paper Yarotsky & Velikanov (2024) to which we refer for details.

**Gradient descent with memory.** Suppose that we wish to minimize a loss function  $L(\mathbf{w})$  on a linear space  $\mathcal{H}$ . We consider gradient descent with size- $M$  memory that can be written as

$$\begin{pmatrix} \mathbf{w}_{t+1} - \mathbf{w}_t \\ \mathbf{u}_{t+1} \end{pmatrix} = \begin{pmatrix} -\alpha & \mathbf{b}^T \\ \mathbf{c} & D \end{pmatrix} \begin{pmatrix} \nabla L(\mathbf{w}_t) \\ \mathbf{u}_t \end{pmatrix}, \quad t = 0, 1, 2, \dots \quad (1)$$

The vector  $\mathbf{w}_t$  is the current step- $t$  approximation to an optimal vector  $\mathbf{w}_*$ , and  $\mathbf{u}_t$  is an auxiliary vector representing the ‘‘memory’’ of the optimizer. These auxiliary vectors have the form  $\mathbf{u} = (\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(M)})^T$  with  $\mathbf{u}^{(m)} \in \mathcal{H}$  and can be viewed as size- $M$  columns with each component belonging to  $\mathcal{H}$ . We refer to  $M$  as the *memory size*. The parameter  $\alpha$  (learning rate) is scalar, the parameters  $\mathbf{a}, \mathbf{b}, \mathbf{c}$  are  $M$ -dimensional column vectors, and  $D$  is a  $M \times M$  scalar matrix. The algorithm can be viewed as a sequence of transformations of size- $(M + 1)$  column vectors  $\begin{pmatrix} \mathbf{w}_t \\ \mathbf{u}_t \end{pmatrix}$  with  $\mathcal{H}$ -valued components. Throughout the paper, we only consider *stationary* algorithms, in the sense that the parameters  $\alpha, \mathbf{b}, \mathbf{c}, D$  do not depend on  $t$ . The simplest nontrivial special case of GD with memory is Heavy Ball (Polyak, 1964), in which  $M = 1$  and  $\mathbf{u}_t$  is the momentum.

Our theoretical results will rely on the assumption that  $L$  is quadratic:

$$L(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{H} \mathbf{w} - \mathbf{w}^T \mathbf{q}, \quad (2)$$

with a strictly positive definite  $\mathbf{H}$ . Throughout the paper, we will mostly be interested in infinite-dimensional Hilbert spaces  $\mathcal{H}$ , and we slightly abuse notation by interpreting  $\mathbf{w}^T$  as the co-vector (linear functional  $\langle \mathbf{w}, \cdot \rangle$ ) associated with vector  $\mathbf{w}$ . We will assume that  $\mathbf{H}$  has a discrete spectrum with ordered eigenvalues  $\lambda_k \searrow 0$ .

Let  $\mathbf{w}_*$  be the optimal value of  $L$  such that  $\nabla L(\mathbf{w}_*) = \mathbf{H}\mathbf{w}_* - \mathbf{q} = 0$ , and denote  $\Delta\mathbf{w}_t = \mathbf{w}_t - \mathbf{w}_*$ . Then, if  $\Delta\mathbf{w}_t$  and  $\mathbf{u}_t$  are eigenvectors of  $\mathbf{H}$  with eigenvalue  $\lambda$ , then

$$\begin{pmatrix} \Delta\mathbf{w}_{t+1} \\ \mathbf{u}_{t+1} \end{pmatrix} = S_\lambda \begin{pmatrix} \Delta\mathbf{w}_t \\ \mathbf{u}_t \end{pmatrix}, \quad S_\lambda = \begin{pmatrix} 1 & \mathbf{b}^T \\ 0 & D \end{pmatrix} + \lambda \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix} (1, \mathbf{0}^T), \quad (3)$$

and the new vectors  $\Delta\mathbf{w}_{t+1}, \mathbf{u}_{t+1}$  are again eigenvectors of  $\mathbf{H}$  with eigenvalue  $\lambda$ . As a result, performing the spectral decomposition of  $\Delta\mathbf{w}_t, \mathbf{u}_t$  reduces the original dynamics (1) acting in  $\mathcal{H} \otimes \mathbb{R}^{M+1}$  to a  $\lambda$ -indexed collection of independent dynamics each acting in  $\mathbb{R}^{M+1}$ .

For quadratic  $L$ , evolution (1) admits an equivalent representation

$$\mathbf{w}_{t+M+1} = \sum_{m=0}^M p_m \mathbf{w}_{t+m} + \sum_{m=0}^M q_m \nabla L(\mathbf{w}_{t+m}), \quad t = 0, 1, \dots, \quad (4)$$

with constants  $(p_m)_{m=0}^M, (q_m)_{m=0}^M$  such that  $\sum_{m=0}^M p_m = 1$ . These constants are found from the characteristic polynomial

$$\chi(x, \lambda) = \det(x - S_\lambda) = P(x) - \lambda Q(x), \quad P(x) = x^{M+1} - \sum_{m=0}^M p_m x^m, \quad Q(x) = \sum_{m=0}^M q_m x^m. \quad (5)$$

**Batch SGD with memory.** In batch Stochastic Gradient Descent (SGD), it is assumed that the loss has the form  $L(\mathbf{w}) = \mathbb{E}_{\mathbf{x} \sim \rho} \ell(\mathbf{x}, \mathbf{w})$ , where  $\rho$  is some probability distribution of data points  $\mathbf{x}$  and  $\ell(\mathbf{x}, \mathbf{w})$  is the loss at the point  $\mathbf{x}$ . In the algorithm (1), we replace  $\nabla L$  by  $\nabla L_{B_t}$ , where  $B_t$  is a random batch of  $|B|$  points sampled from distribution  $\rho$ , and  $\nabla L_B$  is the empirical approximation to  $L$ , i.e.  $L_B = \frac{1}{|B|} \sum_{\mathbf{x} \in B} \ell(\mathbf{x}, \mathbf{w})$ . The samples  $B_t$  at different steps  $t$  are independent.

We assume  $\ell$  to have the quadratic form  $\ell(\mathbf{x}, \mathbf{w}) = \frac{1}{2}(\mathbf{x}^T \mathbf{w} - y(\mathbf{x}))^2$  for some scalar target function  $y(\mathbf{x})$ . Here, the inner product  $\mathbf{x}^T \mathbf{w}$  can be viewed as a linear model acting on the feature vector  $\mathbf{x}$ . By projecting to the subspace of linear functions, we can assume w.l.o.g. that the target function  $y(\mathbf{x})$  is itself linear in  $\mathbf{x}$ , i.e.  $f(\mathbf{x}) = \mathbf{x}^T \mathbf{w}_*$  with some optimal parameter vector  $\mathbf{w}_*$ . (Later we will slightly weaken this assumption to also cover *unfeasible* solutions  $\mathbf{w}_*$ .) Then the full loss is quadratic as in Eq. (2):  $L(\mathbf{w}) = \mathbb{E}_{\mathbf{x} \sim \rho} \frac{1}{2}(\mathbf{x}^T \Delta\mathbf{w})^2 = \frac{1}{2} \Delta\mathbf{w}^T \mathbf{H} \Delta\mathbf{w}$ , where  $\Delta\mathbf{w} = \mathbf{w} - \mathbf{w}_*$  and the Hessian  $\mathbf{H} = \mathbb{E}_{\mathbf{x} \sim \rho} [\mathbf{x}\mathbf{x}^T]$ .

**Mean loss evolution, SE approximation, and the propagator expansion.** Since the trajectory  $\mathbf{w}_t$  in SGD is random, it is convenient to study the deterministic trajectory of batch-averaged losses  $L_t = \mathbb{E}_{B_1, \dots, B_t} L(\mathbf{w}_t)$ . The sequence  $L_t$  can be described exactly in terms of the second moments of  $\mathbf{w}_t, \mathbf{u}_t$  that admit exact evolution equations. An important aspect of this evolution is that it involves 4<sup>th</sup> order moments of the data distribution  $\rho$  and so cannot in general be solved using only the second-order information available in the Hessian  $\mathbf{H} = \mathbb{E}_{\mathbf{x} \sim \rho} [\mathbf{x}\mathbf{x}^T]$ .

A convenient approach to handle this difficulty is the *Spectrally-Expressible (SE) approximation* proposed in Velikanov et al. (2023). It consists in assuming that there exist constants  $\tau_1, \tau_2$  such that for all positive definite operators  $\mathbf{C}$  in  $\mathcal{H}$

$$\mathbb{E}_{\mathbf{x} \sim \rho} [\mathbf{x}\mathbf{x}^T \mathbf{C} \mathbf{x}\mathbf{x}^T] \approx \tau_1 \text{Tr}[\mathbf{H}\mathbf{C}] \mathbf{H} - (\tau_2 - 1) \mathbf{H}\mathbf{C}\mathbf{H}. \quad (6)$$

In fact, this approximation holds *exactly* for some natural types of distribution  $\rho$  (translation-invariant, gaussian). Otherwise, if the r.h.s. is only an upper or lower bound for the l.h.s., this implies a respective relation between the actual losses and the losses computed under the SE approximation. Theoretical predictions obtained under assumption (6) show good quantitative agreement with experiment on real data. We refer to Velikanov et al. (2023); Yarotsky & Velikanov (2024) for further discussion of the SE approximation.

The main benefit of the SE approximation is that it allows to write a convenient loss expansion

$$L_t = \frac{1}{2} \left( V_{t+1} + \sum_{m=1}^t \sum_{0 < t_1 < \dots < t_m < t+1} U_{t+1-t_m} U_{t_m-t_{m-1}} U_{t_{m-1}-t_{m-2}} \cdots U_{t_2-t_1} V_{t_1} \right) \quad (7)$$

with scalar *noise propagators*  $U_t$  and *signal propagators*  $V_t$ . The signal propagators describe the error reduction during optimization in the absence of sampling noise, while the noise propagators describe the perturbing effect of sampling noise injected at various times.

For our main results in Sections 3, 4, we will assume that  $\tau_2 = 0$ , implying particularly simple formulas for  $U_t, V_t$ :

$$U_t = \frac{\tau_1}{|B|} \sum_{k=1}^{\infty} \lambda_k^2 |(1 \ 0^T) S_{\lambda}^{t-1} \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix}|^2, \quad V_t = \sum_{k=1}^{\infty} \lambda_k (\mathbf{e}_k^T \mathbf{w}_*)^2 |(1 \ 0^T) S_{\lambda}^{t-1} \begin{pmatrix} 1 \\ \mathbf{0} \end{pmatrix}|^2, \quad (8)$$

where  $\mathbf{e}_k$  is a normalized eigenvector for  $\lambda_k$ , and it is also assumed that optimization starts from  $\mathbf{w}_0 = \mathbf{0}$  so that  $\Delta \mathbf{w}_0 = \mathbf{w}_0 - \mathbf{w}_* = -\mathbf{w}_*$ .

Importantly, the batch size  $|B|$  affects  $L_t$  only through the denominator in the coefficient in  $U_t$ . The deterministic GD corresponds to the limit  $|B| \rightarrow \infty$ : in this limit  $U_t \equiv 0$  and  $L_t = \frac{1}{2} V_{t+1}$ .

**Convergence/divergence regimes.** Given expansion (7), we can deduce various convergence properties of the loss from the properties of the propagators  $V_t, U_t$ .

**Theorem 1** (Yarotsky & Velikanov (2024)). *Let numbers  $L_t$  be given by expansion (7) with some  $U_t \geq 0, V_t \geq 0$ . Let  $U_{\Sigma} = \sum_{t=1}^{\infty} U_t$  and  $V_{\Sigma} = \sum_{t=1}^{\infty} V_t$ .*

1. **[Convergence]** *Suppose that  $U_{\Sigma} < 1$ . At  $t \rightarrow \infty$ , if  $V_t = O(1)$  (respectively,  $V_t = o(1)$ ), then also  $L_t = O(1)$  (respectively,  $L_t = o(1)$ ).*
2. **[Divergence]** *If  $U_{\Sigma} > 1$  and  $V_t > 0$  for at least one  $t$ , then  $\sup_{t=1,2,\dots} L_t = \infty$ .*
3. **[Signal-dominated regime]** *Suppose that there exist constants  $\xi_V, C_V > 0$  such that  $V_t = C_V t^{-\xi_V} (1 + o(1))$  as  $t \rightarrow \infty$ . Suppose also that  $U_{\Sigma} < 1$  and  $U_t = O(t^{-\xi_U})$  with some  $\xi_U > \max(\xi_V, 1)$ . Then*

$$L_t = \frac{C_V}{2(1 - U_{\Sigma})} t^{-\xi_V} (1 + o(1)). \quad (9)$$

4. **[Noise-dominated regime]** *Suppose that there exist constants  $\xi_V > \xi_U > 1, C_U > 0$  such that  $U_t = C_U t^{-\xi_U} (1 + o(1))$  and  $V_t = O(t^{-\xi_V})$  as  $t \rightarrow \infty$ . Let also that  $U_{\Sigma} < 1$ . Then*

$$L_t = \frac{V_{\Sigma} C_U}{2(1 - U_{\Sigma})^2} t^{-\xi_U} (1 + o(1)). \quad (10)$$

**Spectral power laws.** The detailed convergence results in items 3, 4 of Theorem 1 require us to know the asymptotics of the propagators  $U_t, V_t$ . To this end we introduce power-law spectral assumptions on the eigenvalues and eigenvectors of  $\mathbf{w}_*$  in our optimization problem:

$$\lambda_k = \Lambda k^{-\nu} (1 + o(1)), \quad (11)$$

$$\sum_{k: \lambda_k < \lambda} \lambda_k (\mathbf{e}_k^T \mathbf{w}_*)^2 = Q \lambda^{\zeta} (1 + o(1)), \quad k \rightarrow \infty, \quad (12)$$

with some constants  $\Lambda, Q > 0$  and exponents  $\nu > 0, \zeta > 0$ . Such power laws are common in kernel methods or overparameterized models, and can be derived theoretically or observed empirically (Atanasov et al., 2021; Bordelon & Pehlevan, 2021; Basri et al., 2020; Velikanov & Yarotsky, 2021). Conditions (11), (12) (or their weaker, inequality forms) are usually referred to as the *capacity* and *source* conditions, respectively (Caponnetto & De Vito, 2007). The exponent  $\zeta$  is akin to an inverse effective condition number: lower  $\zeta$  means that the target and the solution have a heavier spectral tail of eigenvectors with small  $\lambda$ , making the problem harder. The exponent  $\nu$  is akin to an inverse effective dimensionality of the problem: lower  $\nu$  means a larger number of eigenvectors above a given spectral parameter  $\lambda$ . Only the source condition (12) matters for the non-stochastic GD rates, but in SGD the capacity condition (11) also becomes important due to the sampling noise.

If  $0 < \zeta < 1$ , then the source condition (12) is inconsistent with  $\mathbf{w}_*$  having a finite  $\mathcal{H}$ -norm, i.e., strictly speaking,  $\mathbf{w}_*$  is not an element of  $\mathcal{H}$ . Such a solution is called *unfeasible*. In fact, unfeasible scenarios are quite common both theoretically and in practice (see Section 6). The Corner SGD to be proposed in Section 4 will be especially suitable for unfeasible scenarios. Note also that if  $\nu < \frac{1}{2}$ , then  $U_1 = \infty$  and so  $L_t \equiv \infty$ , i.e. the loss immediately diverges.

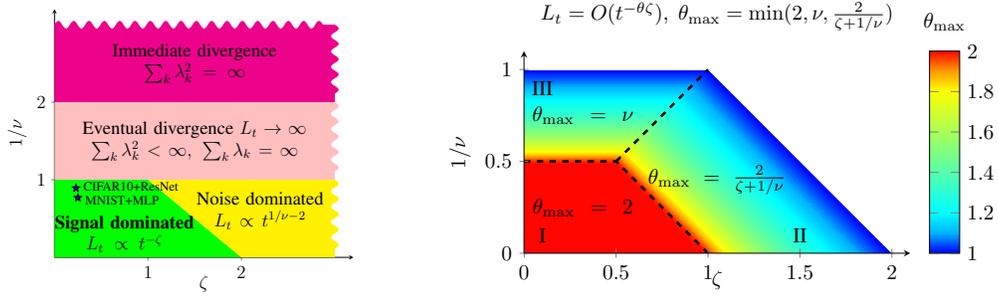


Figure 1: **Left:** The phase diagram of stationary finite-memory SGD from Velikanov et al. (2023); Yarotsky & Velikanov (2024). **Right:** Maximum acceleration factor  $\theta_{\max} = \min(2, \nu, \frac{2}{\zeta + 1/\nu})$  for Corner SGD in the signal-dominated regime (see Theorem 4).

**Stability and asymptotics of the propagators.** Let us say that a square matrix  $A$  is *strictly stable* if all its eigenvalues are less than 1 in absolute value. It is natural to require the matrices  $S_\lambda$  to be strictly stable for all  $\lambda \in \text{spec}(\mathbf{H})$ , since otherwise  $U_t, V_t$ , and hence  $L_t$ , will not generally even converge to 0 as  $t \rightarrow \infty$ . At  $\lambda = 0$  the matrix  $S_{\lambda=0}$  has eigenvalue 1 and additionally eigenvalues of the matrix  $D$ ; accordingly, we will assume that  $D$  is strictly stable.

**Theorem 2.** *Suppose that  $D$  and  $S_\lambda$  are strictly stable for all  $\lambda \in \text{spec}(\mathbf{H})$ . Recalling the characteristic polynomial  $\chi(x, \lambda) = \det(x - S_\lambda) = P(x) - \lambda Q(x)$ , define the effective learning rate*

$$\alpha_{\text{eff}} = -Q(1) \Big/ \frac{dP}{dx}(1), \quad (13)$$

and assume that  $\alpha_{\text{eff}} > 0$ . Then, under spectral assumptions (11), (12) with  $\nu > \frac{1}{2}$ , the propagators  $V_t, U_t$  given by Eq. (8) obey, as  $t \rightarrow \infty$ ,

$$V_t = (1 + o(1)) Q \Gamma(\zeta + 1) (2\alpha_{\text{eff}} t)^{-\zeta}, \quad (14)$$

$$U_t = (1 + o(1)) \frac{(\alpha_{\text{eff}} \Lambda)^{1/\nu} \tau_1 \Gamma(2 - 1/\nu)}{|B|^\nu} (2t)^{1/\nu - 2}. \quad (15)$$

Combined with Theorem 1, this result yields the  $(\zeta, 1/\nu)$ -phase diagram shown in Figure 1 left. In particular, the region  $\nu > 1, 0 < \zeta < 2 - 1/\nu$  represents the signal-dominated phase in which the noise effects are relatively weak and the loss convergence  $L_t \propto t^{-\zeta}$  has the same exponent  $\zeta$  as plain deterministic GD. This holds for all stationary finite- $M$  algorithms and so such algorithms cannot accelerate the exponent. In the present paper we will focus on the signal-dominated phase and propose an “infinite-memory” generalization of SGD that does accelerate the exponent.

### 3 THE CONTOUR VIEW OF GENERALIZED (S)GD

We consider the propagator expansion (7) as a basis for our arguments. Observe that we can write the expression  $(\mathbf{1} \mathbf{0}^T) S_\lambda^t (\frac{-\alpha}{\mathbf{c}})$  appearing in the definition of propagator  $U_t$  in Eq. (8) as

$$(\mathbf{1} \mathbf{0}^T) S_\lambda^t (\frac{-\alpha}{\mathbf{c}}) = \frac{1}{2\pi i} \oint_{|\mu|=r} \mu^t (\mathbf{1} \mathbf{0}^T) (\mu - S_\lambda)^{-1} (\frac{-\alpha}{\mathbf{c}}) d\mu, \quad (16)$$

where  $|\mu| = r$  is a contour in the complex plane encircling all the eigenvalues of  $S_\lambda$ . Next, simple calculation (see Section A) shows that

$$(\mathbf{1} \mathbf{0}^T) (\mu - S_\lambda)^{-1} (\frac{-\alpha}{\mathbf{c}}) = \frac{Q(\mu)}{P(\mu) - \lambda Q(\mu)} = \frac{1}{\frac{P(\mu)}{Q(\mu)} - \lambda} = \frac{1}{\Psi(\mu) - \lambda}, \quad (17)$$

where  $P(x) - \lambda Q(x)$  is the characteristic polynomial of  $S_\lambda$  introduced in Eq. (5), and

$$\Psi(\mu) = \frac{P(\mu)}{Q(\mu)}. \quad (18)$$

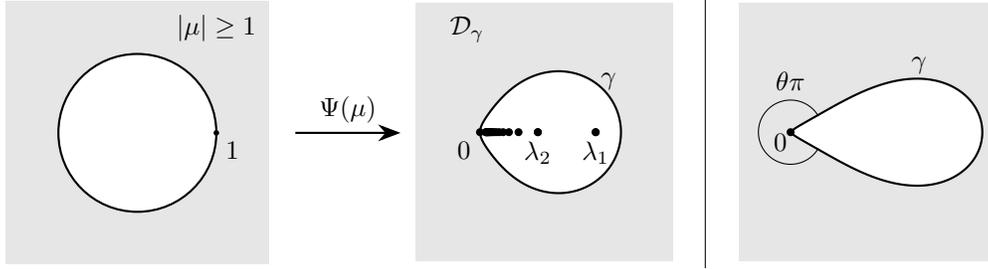


Figure 2: **Left:** The map  $\Psi = \frac{P}{Q}$  for a memory-1 algorithm with  $P(\mu) = (\mu - 1)(\mu - 0.7)$  and  $Q(\mu) = 1.1 - 2\mu$ . The contour  $\gamma = \Psi(\{\mu : |\mu| = 1\})$  encircles  $\text{spec}(\mathbf{H})$ . The map  $\Psi$  bijectively maps  $\{\mu \in \mathbb{C} : |\mu| > 1\}$  to the exterior open domain  $\mathcal{D}_\gamma$  with boundary  $\gamma$ . **Right:** Contour  $\gamma$  corresponding to a corner map  $\Psi$  with angle  $\theta\pi$ .

We see, in particular, that the values  $U_t$  depend on the algorithm parameters only through the function  $\Psi$ . The same observation can also be made regarding the values  $V_t$ . Indeed,  $V_t$ 's are different from  $U_t$ 's in that they involve the expression  $(\mathbf{1} \ \mathbf{0}^T) S_\lambda^t(\frac{\mathbf{1}}{\mathbf{0}})$  instead of  $(\mathbf{1} \ \mathbf{0}^T) S_\lambda^t(-\frac{\mathbf{1}}{\mathbf{c}})$ . The contour representation for  $(\mathbf{1} \ \mathbf{0}^T) S_\lambda^t(\frac{\mathbf{1}}{\mathbf{0}})$  is similar to Eq. (16), and then a simple calculation gives

$$(\mathbf{1} \ \mathbf{0}^T)(\mu - S_\lambda)^{-1}(\frac{\mathbf{1}}{\mathbf{0}}) = \frac{\Psi(\mu)}{(\Psi(\mu) - \lambda)(\mu - 1)}. \quad (19)$$

Recall from Eqs. (4),(5) that  $P$  can be any monic polynomial (i.e., with leading coefficient 1) of degree  $M + 1$  such that  $P(1) = 0$ , while  $Q$  can be any polynomials of degree not greater than  $M$ . Since, by Eq. (7) the loss trajectory  $L_t$  is completely determined by the propagators  $U_t, V_t$ , we see that designing a stationary SGD with memory is essentially equivalent to designing a rational function  $\Psi$  subject to these simple conditions.

Let us consider the map  $\Psi$  from the stability perspective. Recall that we expect  $S_{\lambda_k}$  to be strictly stable for all the eigenvalues  $\lambda_k \in \text{spec}(\mathbf{H})$ . In terms of  $\Psi = \frac{P}{Q}$  this means that  $\Psi(\mu) \neq \lambda_k$  for all  $\mu \in \mathbb{C}$  such that  $|\mu| \geq 1$ . Additionally, if  $D$  is strictly stable, then  $S_0$  has only one simple eigenvalue of unit absolute value,  $\mu = 1$ , and so  $\Psi(\mu) \neq 0$  for  $|\mu| = 1, \mu \neq 1$ . Let us introduce the curve  $\gamma$  as the image of the unit circle under the map  $\Psi$ . Then the last condition means that the curve  $\gamma$  goes through the point 0 only once, at  $\mu = 1$ .

In general, the curve  $\gamma$  can have a complicated shape with self-intersections, and the map  $\Psi$  may not be injective. In particular, the singularity of  $\Psi$  at  $\mu = \infty$  is  $\propto \mu^{M+1-\text{deg}(Q)}$ , so in a vicinity of  $\mu = \infty$  the function  $\Psi$  is injective if and only if  $\text{deg}(Q) = M$ . However, we may expect natural, non-degenerate algorithms to correspond to simple non-intersecting curves  $\gamma$ . This is confirmed by considering simple memory-1 algorithms (Fig. 2 left).

Given a non-intersecting (Jordan) contour  $\gamma$ , denote by  $\mathcal{D}_\gamma$  the respective exterior open domain. Then, by Riemann mapping theorem, there exists a bijective holomorphic map  $\Psi_\gamma : \{\mu \in \mathbb{C} : |\mu| > 1\} \rightarrow \mathcal{D}_\gamma$ . Additionally, by Carathéodory's theorem<sup>1</sup> this map extends continuously to the boundary,  $\Psi_\gamma : \{\mu \in \mathbb{C} : |\mu| = 1\} \rightarrow \gamma$ . Such maps  $\Psi_\gamma$  are non-unique and rather form a three-parameter family  $\Psi_\gamma \circ f$ , where  $f$  is a conformal automorphism of  $\{\mu \in \mathbb{C} : |\mu| > 1\}$ . However, recall that our maps  $\Psi = \frac{P}{Q}$  had the properties  $\Psi(\infty) = \infty$  and  $\Psi(1) = 0$ . The requirements  $\Psi_\gamma(\infty) = \infty$  and  $\Psi_\gamma(1) = 0$  uniquely fix the conformal isomorphism and hence  $\Psi_\gamma$ .

This suggests the following reformulation of the design problem for stationary SGD with memory. Rather than starting with the algorithm in the matrix or sequential forms (1), (4), we can start with a contour  $\gamma$  or the associated Riemann map  $\Psi_\gamma$ , and ensure a fast decay of the respective propagators  $U_t, V_t$  and hence, by Theorem 1, the loss  $L_t$ . Of course, the resulting map  $\Psi_\gamma$  will not be rational in general, but we can hope to subsequently approximate it with a rational function  $\frac{P}{Q}$  and in this way approximately reconstruct the algorithm. We will see that this plan indeed works well.

<sup>1</sup>Carathéodory's theorem considers bounded domains, but our domains  $\{\mu \in \mathbb{C} : |\mu| > 1\}$  and  $\mathcal{D}_\gamma$  are conformally isomorphic to bounded ones by simple transformations  $z = 1/(\mu - \mu_0)$ .

Thus, given a map  $\Psi$  we introduce the values  $U(t, \lambda)$ ,  $V(t, \lambda)$  that generalize the expressions  $(\mathbf{1} \circ^T) S_\lambda^t \binom{-\alpha}{-c}$  and appear in the generalized definition (8) of the propagators  $U_t, V_t$ :

$$U_t = \frac{\tau_1}{|B|} \sum_{k=1}^{\infty} \lambda_k^2 |U(t, \lambda_k)|^2, \quad V_t = \sum_{k=1}^{\infty} (\mathbf{e}_k^T \mathbf{w}_*)^2 |V(t, \lambda_k)|^2 \quad (20)$$

$$U(t, \lambda) = \frac{1}{2\pi i} \oint_{|\mu|=1} \frac{\mu^{t-1} d\mu}{\Psi(\mu) - \lambda} = \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{e^{ixt} dx}{\Psi(e^{ix}) - \lambda} \quad (21)$$

$$V(t, \lambda) = \frac{1}{2\pi i} \oint_{|\mu|=1} \frac{\mu^{t-1} \Psi(\mu) d\mu}{(\Psi(\mu) - \lambda)(\mu - 1)} = \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{e^{ixt} \Psi(e^{ix}) dx}{(\Psi(e^{ix}) - \lambda)(e^{ix} - 1)} \quad (22)$$

The functions  $\frac{1}{\Psi(\mu) - \lambda}$  and  $\frac{\Psi(\mu)}{(\Psi(\mu) - \lambda)(\mu - 1)}$  are holomorphic in  $\{|\mu| > 1\}$  and vanish as  $\mu \rightarrow \infty$ , so the Fourier coefficients  $U(t, \lambda), V(t, \lambda)$  vanish for all negative  $t = -1, -2, \dots$ . In particular, by Parseval's identity, the  $t$  series in the total noise coefficient  $U_\Sigma$  collapses to the squared  $L^2$  norm:

$$U_\Sigma = \frac{\tau_1}{2\pi|B|} \sum_{k=1}^{\infty} \lambda_k^2 \int_{-\pi}^{\pi} \frac{d\phi}{|\Psi(e^{i\phi}) - \lambda_k|^2}. \quad (23)$$

#### 4 CORNER ALGORITHMS

To motivate the algorithms introduced in this section, observe from Theorems 1, 2 that in the signal-dominated regime of stationary memory- $M$  SGD, asymptotically accelerating convergence is essentially equivalent to increasing  $\alpha_{\text{eff}}$  while keeping the total noise coefficient  $U_\Sigma < 1$ . Since  $\Psi(1) = 0$ ,  $\alpha_{\text{eff}}$  can be reformulated in terms of  $\Psi$  as

$$\alpha_{\text{eff}} = -\frac{Q(1)}{\frac{dP}{dx}(1)} = -\left(\frac{d\Psi}{d\mu}(1)\right)^{-1}. \quad (24)$$

Thus, increasing  $\alpha_{\text{eff}}$  means making  $-\frac{d\Psi}{d\mu}(1)$  a possibly smaller positive number. Regarding  $U_\Sigma$ , note that it is inversely proportional to the batch size  $|B|$ , so if the series in Eq. (23) converges, we can always ensure  $U_\Sigma < 1$  by making  $|B|$  large enough.

It is then natural to try  $\Psi = \Psi_\gamma$  with a contour  $\gamma$  having a corner at 0 with a particular angle. Denote the angle by  $\theta\pi$  when measured in the external domain  $\mathcal{D}_\gamma$  (Figure 2 right). Such contours correspond to maps  $\Psi : \{|\mu| > 1\} \rightarrow \mathcal{D}_\gamma$  such that

$$\Psi(\mu) = -c_\Psi(\mu - 1)^\theta(1 + o(1)), \quad \mu \rightarrow 1, \quad (25)$$

with the standard branch of  $(\mu - 1)^\theta$  and some constant  $c_\Psi > 0$ . We will refer to such  $\Psi$  as *corner maps* and to the respective generalized SGD as *corner algorithms*. Formally,

$$-\frac{d\Psi}{d\mu}(\mu = 1) \sim c\theta(\mu - 1)^{\theta-1}|_{\mu=1+} = \begin{cases} +\infty, & \theta < 1 \\ +0, & \theta > 1 \end{cases} \quad (26)$$

so we are interested in  $\theta > 1$ . At the same time, we cannot take  $\theta > 2$ , since this would violate the stability condition  $\Psi\{|\mu| > 1\} \cap \text{spec}(\mathbf{H}) = \emptyset$ . Thus, the relevant range of values for  $\theta$  is  $[1, 2]$ . Within this range, increasing  $\theta$  should have a positive  $\alpha_{\text{eff}}$ -related effect but a negative  $U_\Sigma$ -related effect, since the contour  $\gamma = \Psi(|\mu| = 1)$  is getting closer to the spectral segment  $[0, \lambda_{\text{max}}]$ , thus amplifying the singularity  $|\Psi(e^{i\phi}) - \lambda_k|^{-2}$  in Eq. (23). Our main technical result is

**Theorem 3 (B).** *Let  $\Psi$  be a holomorphic function in  $\{\mu \in \mathbb{C} : |\mu| > 1\}$  commuting with complex conjugation and obeying power law condition (25) with some  $1 < \theta < 2$ . Assume that  $\Psi$  extends continuously to a  $C^1$  function on the closed domain  $|\mu| \geq 1$ ,  $\Psi(\mu) \rightarrow \infty$  as  $\mu \rightarrow \infty$ , and  $\frac{d}{d\mu}\Psi(\mu) = O(|\mu - 1|^{\theta-1})$  as  $\mu \rightarrow 1$ . Assume also that  $\Psi(\{\mu \in \mathbb{C} : |\mu| = 1, \mu \neq 1\}) \cap [0, \lambda_{\text{max}}] = \emptyset$ , where  $\lambda_{\text{max}} = \lambda_1$  is the largest eigenvalue of  $\mathbf{H}$ . Let power-law spectral assumptions (11),(12) hold with some  $\nu > 1, 0 < \zeta < 2$ . Then propagators (20) obey the following  $t \rightarrow \infty$  asymptotics.*

1. **(Noise propagators)**  $U_t = C_U t^{\theta/\nu-2}(1 + o(1))$ , with the coefficient

$$C_U = \frac{\tau_1}{|B|} \Lambda^{1/\nu} \int_{\infty}^0 r^2 F_U^2(r) dr^{-\theta/\nu} < \infty, \quad F_U(r) = \frac{1}{2\pi i} \int_{i\mathbb{R}} \frac{e^{rz} dz}{c_\Psi z^\theta + 1}.$$

2. **(Signal propagators)**  $V_t = C_V t^{-\theta\zeta}(1 + o(1))$ , with the coefficient

$$C_V = Q \int_0^\infty F_V^2(r) dr^{\theta\zeta} < \infty, \quad F_V(r) = \frac{1}{2\pi i} \int_{i\mathbb{R}} \frac{c_\Psi z^{\theta-1} e^{rz} dz}{c_\Psi z^\theta + 1}.$$

We see that the leading  $t \rightarrow \infty$  asymptotics of the propagators are completely determined by the  $\lambda \searrow 0$  spectral asymptotics of the problem and the  $\mu \rightarrow 1$  singularity of the map  $\Psi$ . The functions  $F_U, F_V$  can be written in terms of the Mittag-Leffler functions  $E_{\theta,\theta}, E_\theta$  (see Section B).

Availability of the coefficients  $C_U, C_V$  ensures that the leading asymptotics of  $U_t, V_t$  are strict power laws with specific exponents  $2 - \theta/\nu$  and  $\theta\zeta$ , respectively. Increasing  $\theta$  indeed improves convergence of the signal propagators, but degrades convergence of the noise propagators.

The largest acceleration of the loss exponent  $\zeta$  possibly achievable with corner algorithms is by a factor  $\theta$  arbitrarily close to 2, but in general it will be lower since, by Theorem 1, the exponent of  $L_t$  is the lower of the exponents of  $U_t$  and  $V_t$ ; accordingly, the optimal  $\theta$  is obtained by balancing the two exponents, i.e. setting  $\theta\zeta = 2 - \theta/\nu$ . Also, we need the noise exponent  $2 - \theta/\nu$  to be  $> 1$ , since otherwise the total noise coefficient  $U_\Sigma = \infty$  and  $L_t$  diverges for any batch size  $|B| < \infty$ .

Combining these considerations, we get the phase diagram of feasible accelerations (Figure 1 right).

**Theorem 4.** Consider a problem with power-law spectral conditions (11),(12) in the signal-dominated phase, i.e.  $\nu > 1, 0 < \zeta < 2 - 1/\nu$ . Let  $\theta_{\max}$  denote the supremum of those  $\theta$  for which there exists a corner algorithm and batch size  $B$  such that  $L_t = O(t^{-\theta\zeta})$ . Then

$$\theta_{\max} = \min\left(2, \nu, \frac{2}{\zeta + 1/\nu}\right). \quad (27)$$

The phase diagram thus has three regions:

- I. **Fully accelerated:**  $\theta_{\max} = 2$ , achieved for  $\nu > 2, 0 < \zeta < 1 - 1/\nu$ .
- II. **Signal/noise balanced:**  $\theta_{\max} = \frac{2}{\zeta + 1/\nu} < 2$ ,  $\max(1/\nu, 1 - 1/\nu) < \zeta < 2 - 1/\nu$ . The condition  $1/\nu < \zeta$  ensures that  $U_\Sigma$  is finite and less than 1 for  $|B|$  large enough.
- III. **Limited by  $U_\Sigma$ -finiteness:**  $\theta_{\max} = \nu < 2$ ,  $1 < \nu < 2, 0 < \zeta < 1/\nu$ . The signal exponent  $\theta_{\max}\zeta$  is less than the noise exponent  $2 - \theta_{\max}/\nu$ , but increasing  $\theta$  makes  $U_\Sigma$  diverge.

## 5 FINITE-MEMORY APPROXIMATIONS OF CORNER ALGORITHMS

Though corner maps  $\Psi$  are irrational, they can be efficiently approximated by rational functions. It was originally famously discovered by Newman (1964) that the function  $|x|$  can be approximated by order- $M$  rational functions with error  $O(e^{-c\sqrt{M}})$ . This result was later refined in various ways. In particular, Gopal & Trefethen (2019) establish a rational approximation with a similar error bound for general power functions  $z \mapsto z^\theta$  on complex domains. For  $\theta \in (0, 1)$ , this is done by writing

$$z^\theta = \frac{\sin(\theta\pi)}{\theta\pi} \int_0^\infty \frac{z dt}{t^{1/\theta} + z} = \frac{\sin(\theta\pi)}{\theta\pi} \int_{-\infty}^\infty \frac{z e^{\theta\pi i/2+s} ds}{e^{\pi i/2+s/\theta} + z} \quad (28)$$

and then approximating the last integral by the trapezoidal rule with uniform spacing  $h = \pi\sqrt{2\theta/M}$ .

In our setting, we start by explicitly defining a  $\theta$ -corner map. This can be done in many ways; we find it convenient to set

$$\Psi(\mu) = -A \left( \int_0^1 \frac{d\delta^{2-\theta}}{\mu - 1 + \delta} \right)^{-1} \frac{\mu - 1}{\mu} = A \left( (\theta - 2) \int_0^\infty \frac{e^{-(2-\theta)s} ds}{\mu - 1 + e^{-s}} \right)^{-1} \frac{\mu - 1}{\mu} \quad (29)$$

with a scaling parameter  $A > 0$ .

**Proposition 1 (C).** For any  $1 < \theta < 2$ , Eq. (29) defines a holomorphic map  $\Psi : \mathbb{C} \setminus [0, 1] \rightarrow \mathbb{C}$  such that

$$\Psi(\mu) = \begin{cases} -A\mu(1 + o(1)), & \mu \rightarrow \infty, \\ -\frac{A(2-\theta)\pi}{\sin((2-\theta)\pi)} (\mu - 1)^\theta (1 + o(1)), & \mu \rightarrow 1, \end{cases} \quad (30)$$

where  $z^\theta$  denotes the standard branch in  $\mathbb{C} \setminus (-\infty, 0]$ . Also,  $\Psi(\{|\mu| \geq 1\}) \cap (0, 2A] = \emptyset$ .

Following Gopal & Trefethen (2019), we approximate the last integral in Eq. (29) as

$$\int_0^\infty \phi(s) ds \approx h \sum_{m=1}^M \phi((m - \frac{1}{2})h), \quad h = \frac{l}{\sqrt{M}}, \quad (31)$$

with some fixed constant  $l$ . Note that in contrast to (28), our integral and discretization are “one-sided” ( $s > 0$ ), reflecting the fact that the corner map  $\Psi(\mu)$  is power law only at  $\mu \rightarrow 1$ , which is related to the  $s \rightarrow +\infty$  behavior of the integrand.

Let  $\Psi^{(M)}$  denote the map  $\Psi$  discretized with  $M$  nodes by scheme (31). Observe that  $\Psi^{(M)}$  is a rational function,  $\Psi^{(M)} = \frac{P}{Q}$ , where  $\deg P = M + 1$  and  $\deg Q \leq M$  (in particular,  $P(\mu) = (\mu - 1) \prod_{m=1}^M (\mu - 1 + e^{-(m-1/2)h})$ ). We can then associate to  $\Psi^{(M)}$  a memory- $M$  algorithm (1) with particular  $\alpha, \mathbf{b}, \mathbf{c}, D$ , for example as follows.

**Proposition 2 (D).** *Let  $h = l/\sqrt{M}$  and*

$$D = \text{diag}(1 - e^{-\frac{1}{2}h}, \dots, 1 - e^{-(M-\frac{1}{2})h}), \quad (32)$$

$$\mathbf{b} = (1, \dots, 1)^T, \quad (33)$$

$$\mathbf{c} = (c_1, \dots, c_M)^T, \quad c_m = A^{-1}(2 - \theta)h e^{-(2-\theta)(m-1/2)h} (e^{-(m-1/2)h} - 1), \quad (34)$$

$$\alpha = A^{-1}(2 - \theta)h \frac{1 - e^{-(2-\theta)Mh}}{1 - e^{-(2-\theta)h}} e^{-(2-\theta)h/2}. \quad (35)$$

*Then the respective characteristic polynomial  $\chi(\mu) = P(\mu) - \lambda Q(\mu)$  with  $\frac{P}{Q} = \Psi^{(M)}$ .*

Of course, as any stationary finite-memory algorithm, for very large  $t$  the  $M$ -discretized corner algorithm can only provide a  $O(t^{-\zeta})$  convergence of the loss. But, thanks to the  $O(e^{-c\sqrt{M}})$  rational approximation bound, we expect that even with moderate  $M$ , for practically relevant finite ranges of  $t$  the convergence should be close to  $O(t^{-\theta\zeta})$  of the ideal corner algorithm.

## 6 EXPERIMENTS

The experiments in this section are performed with Corner SGD approximated as in Proposition 2 with memory size  $M = 5$  and spacing parameter  $l = 5$ .

**A synthetic indicator problem.** Suppose that we are fitting the indicator function  $y(x) = \mathbf{1}_{[1/4, 3/4]}(x)$  on the segment  $[0, 1]$  using the shallow ReLU neural network in which only the output layer weights  $w_n$  are trained:

$$\hat{y}(x, \mathbf{w}) = \frac{1}{N} \sum_{n=1}^N w_n (x - \frac{n}{N})_+, \quad (x)_+ \equiv \max(x, 0). \quad (36)$$

This is an exactly linear model that in the limit  $N \rightarrow \infty$  acquires the form

$$\hat{y}(x) = \int_0^1 w(y) (x - y)_+ dy = \mathbf{x}^T \mathbf{w}, \quad (37)$$

where  $\mathbf{x}, \mathbf{w}$  are understood as vectors in  $L^2([0, 1])$ , and  $\mathbf{x} \equiv (x - \cdot)_+$ . We consider the loss  $L(\mathbf{w}) = \mathbb{E}_{x \sim U(0,1)} \frac{1}{2} (\mathbf{x}^T \mathbf{w} - y(x))^2$ , where  $U(0, 1)$  is the uniform distribution on  $[0, 1]$ .

This limiting integral problem obeys asymptotic spectral power laws (11),(12) with precisely computable  $\nu, \zeta$  (see Appendix E):

$$\zeta = \frac{1}{4}, \quad \nu = 4. \quad (38)$$

The problem thus falls into the sub-phase I “full acceleration” of the signal dominated phase, and we expect that it can be accelerated with corner algorithms up to  $\theta_{\max} = 2$ .

In the experiment we set  $N = 5 \cdot 10^4$  and apply corner SGD with  $\theta = 1.8$ , see Figure 3. We observe good agreement of both the theoretical exponent  $\zeta = 0.25$  of plain SGD and the accelerated exponent  $\theta\zeta = 1.8 \cdot 0.25 = 0.45$  of corner SGD with the experimentally estimated values.

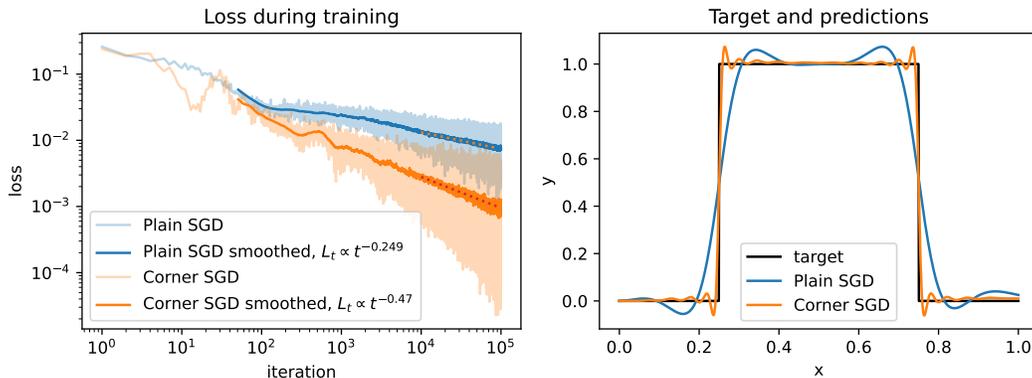


Figure 3: Training loss and final predictions of the linear model (36) trained to fit the target  $y(x) = \mathbf{1}_{[1/4, 3/4]}(x)$  using either plain or corner SGD with batch size  $|B| = 100$ . The loss trajectories oscillate strongly, so their smoothed versions are also shown and used to estimate the exponents  $\zeta$  in power laws  $L_t \propto t^{-\zeta}$ . The corner SGD has  $\theta = 1.8$ , which is close to observed acceleration  $\frac{0.47}{0.249}$ .

**MNIST.** We consider MNIST digit classification performed by a single-hidden-layer ReLU neural network:

$$\hat{y}_r(\mathbf{x}, \mathbf{w}) = \frac{1}{\sqrt{H}} \sum_{n=1}^H w_{rn}^{(2)} \left( \sum_{m=1}^{28 \times 28} w_{nm}^{(1)} x_m \right)_+, \quad r = 0, \dots, 9. \quad (39)$$

Here, the input vector  $\mathbf{x} = (x_m)_{m=1}^{28 \times 28}$  represents a MNIST image, and the outputs  $y_r$  represent the 10 classes. We use the one-hot encoding for the targets  $\mathbf{y}(\mathbf{x})$  and the quadratic pointwise loss  $\ell(\mathbf{x}, \mathbf{w}) = \frac{1}{2} |\hat{\mathbf{y}}(\mathbf{x}, \mathbf{w}) - \mathbf{y}(\mathbf{x})|^2$  for training. The trainable weights include both first- and second-layer weights  $w_{nm}^{(1)}, w_{rn}^{(2)}$ .

Note that the model (39) is nonlinear, but for large width  $H$  and standard independent weight initialization it belongs to the approximately linear NTK regime (Jacot et al., 2018). In Velikanov & Yarotsky (2021) MNIST was found to have an approximate power-law spectrum with

$$\zeta \approx 0.25, \quad \nu \approx 1.3, \quad (40)$$

putting this problem in the sub-phase III “limited by  $U_{\Sigma}$ -finiteness” of the signal-dominated phase (see Figure 1). Theoretically, by Theorem 4, the largest feasible acceleration in this case is  $\theta_{\max} = \nu$ . Note, however, that this theoretical prediction relied on the infinite-dimensionality of the problem and the divergence of the series  $\sum_{t=1}^{\infty} t^{\theta/\nu-2}$ . The actual MNIST problem is finite-dimensional, so its  $U_{\Sigma}$  is always finite (though possibly large) and can be made  $< 1$  if  $|B|$  is large enough. This suggests that corner SGD might practically be used with  $\theta > \nu$  and possibly display acceleration beyond the theoretical bound  $\theta_{\max} = \nu$ . Note also that with exponents (40) the signal/noise balance bound  $\frac{2}{\zeta+1/\nu} \approx 2$ , i.e. it is not an obstacle for increasing the parameter  $\theta$  towards 2.

In Figure 4 we test corner SGD with  $\theta = 1.3$  or  $1.8$  on batch sizes  $|B| = 1000$  and  $100$ . The  $\theta = 1.3$  version shows a stable performance accelerating the plain SGD exponent  $\zeta$  by a factor  $\sim 1.5$ . The  $\theta = 1.8$  version shows lower losses, but does not significantly improve acceleration factor  $1.5$  at  $|B| = 1000$  and is unstable at  $|B| = 100$ . We remark that on the test set the loss and prediction error of Corner SGD also decrease faster compared to plain SGD (see Appendix F).

## 7 ADDITIONAL NOTES AND DISCUSSION

**Extension to SE approximation with  $\tau_2 \neq 0$ .** The key assumption in our derivation and analysis of the contour representation and corner algorithms was the Spectrally Expressible approximation with  $\tau_2 = 0$  for the SGD moment evolution (see Eq. (6)). While the SE approximation in general was justified from several points of view in Velikanov et al. (2023); Yarotsky & Velikanov (2024), a natural question is how important is the condition  $\tau_2 = 0$ . This condition substantially simplifies the representation of propagators  $U_t, V_t$  in Eqs. (8), but does not seem to correspond to any specific

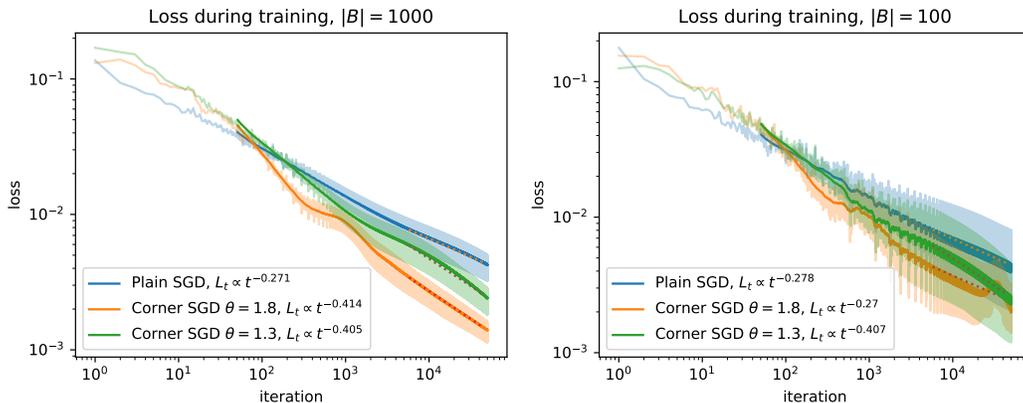


Figure 4: Training loss of neural network (39) on MNIST classification with  $H = 1000$ , with batch size  $|B| = 1000$  (left) or  $100$  (right). The full color curves show the smoothed losses.

natural data distribution  $\rho$ . (In contrast, the cases  $\tau_1 = \tau_2 = 1$  and  $\tau_1 = 1, \tau_2 = -1$  exactly describe translation-invariant and Gaussian distributions; see Velikanov et al. (2023).)

In fact, our analysis of the corner propagators  $U_t, V_t$  can be extended from  $\tau_2 = 0$  to general  $\tau_2$  by a perturbation theory around  $\tau_2 = 0$ . In Appendix G we sketch an argument suggesting that, at least for sufficiently large batch sizes  $|B|$ , Theorem 3 remains valid for general  $\tau_2$ , even with the same coefficients  $C_U, C_V$  (i.e., the contribution from  $\tau_2 \neq 0$  produces only subleading terms in  $U_t, V_t$ ). This implies, in particular, that the acceleration phase diagram in Theorem 4 and Figure 1 (right) is not only  $\tau_1$ -, but also  $\tau_2$ -independent.

**Computational complexity.** The main overhead of finitely-approximated corner algorithms compared to plain SGD lies in the memory requirements: if the model has  $W$  weights (i.e.,  $\dim \mathbf{w}_t = W$  in Eq. (1)), then a memory- $M$  algorithm needs to additionally store  $MW$  scalars in the auxiliary vectors  $\mathbf{u}_t$ . On the other hand, the number of elementary operations (arithmetic operations and evaluations of standard elementary functions) in a single iteration of a finitely-approximated corner algorithm need not be much larger than for plain SGD.

Indeed, an iteration (1) of a memory- $M$  algorithm consists in computing the gradient  $\nabla L(\mathbf{w}_t)$  and performing a linear transformation. In SGD with batch size  $|B|$ , the estimated gradient  $\nabla L_{B_t}(\mathbf{w}_t)$  is computed by backpropagation using  $\propto |B|W$  operations. If Corner SGD is finitely-approximated using a diagonal matrix  $D$  as in Proposition 2, then the number of operations in the linear transformation is  $O(MW)$ . Accordingly, if  $|B| \gg M$  (which should typically be the case in practice), then the computational cost of the linear transformation is negligible compared to the batch gradient estimation, and so the computational overhead of Corner SGD is negligible compared to plain SGD.

**Practical and theoretical acceleration.** Our MNIST experiment in Section 6 shows that finitely-approximated Corner SGD developed in Section 5 can practically accelerate learning even on realistic problems that are not exactly linear. We note, however, that, in contrast to the ideal infinite-memory Corner SGD of Section 4, this finitely-approximated Corner SGD does not theoretically accelerate the convergence exponent  $\zeta$  as  $t \rightarrow \infty$ . (As shown in Yarotsky & Velikanov (2024), this is generally impossible for stationary algorithms with finite linear memory.) Nevertheless, we expect that such an acceleration can be achieved with a suitable *non-stationary* approximation. In Yarotsky & Velikanov (2024), an acceleration with a factor  $\theta$  up to  $2 - 1/\nu$  was heuristically derived for a suitable non-stationary memory-1 SGD algorithm.

We remark also that if the model includes nonlinearities, then even the plain SGD in the signal-dominated regime may show a complex picture of convergence rates depending on the strength of the feature learning effects. In particular, Bordelon et al. (2024) consider a particular model where the “rich training” regime is argued to accelerate the “lazy training” exponent  $\zeta$  by the factor  $\frac{2}{1+\zeta}$ . This is different from our factor  $\theta_{\max} = \min(2, \nu, \frac{2}{\zeta+1/\nu})$  due to a different acceleration mechanism.

---

## ACKNOWLEDGMENT

The author thanks the anonymous reviewers for useful suggestions that helped to improve the paper.

## 8 BIBLIOGRAPHY

- Alexander Atanasov, Blake Bordelon, and Cengiz Pehlevan. Neural networks as kernel learners: The silent alignment effect. *arXiv preprint arXiv:2111.00034*, 2021.
- Ronen Basri, Meirav Galun, Amnon Geifman, David Jacobs, Yoni Kasten, and Shira Kritchman. Frequency bias in neural networks for input of non-uniform density. In *International Conference on Machine Learning*, pp. 685–694. PMLR, 2020.
- Raphaël Berthier, Francis Bach, and Pierre Gaillard. Tight nonparametric convergence rates for stochastic gradient descent under the noiseless linear model. *arXiv preprint arXiv:2006.08212*, 2020.
- M Š Birman and M Z Solomjak. Asymptotic behavior of the spectrum of weakly polar integral operators. *Mathematics of the USSR-Izvestiya*, 4(5):1151–1168, oct 1970. doi: 10.1070/im1970v004n05abeh000948. URL <https://doi.org/10.1070/im1970v004n05abeh000948>.
- Blake Bordelon and Cengiz Pehlevan. Learning curves for sgd on structured features. *arXiv preprint arXiv:2106.02713*, 2021.
- Blake Bordelon, Alexander Atanasov, and Cengiz Pehlevan. How feature learning can improve neural scaling laws. *arXiv preprint arXiv:2409.17858*, 2024.
- Helmut Brakhage. On ill-posed problems and the method of conjugate gradients. In *Inverse and ill-posed Problems*, pp. 165–175. Elsevier, 1987.
- Andrea Caponnetto and Ernesto De Vito. Optimal rates for the regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7(3):331–368, 2007.
- Abinand Gopal and Lloyd N Trefethen. Representation of conformal maps by rational functions. *Numerische Mathematik*, 142:359–382, 2019.
- Hans J Haubold, Arak M Mathai, and Ram K Saxena. Mittag-leffler functions and their applications. *Journal of applied mathematics*, 2011(1):298628, 2011.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *arXiv preprint arXiv:1806.07572*, 2018.
- Arkadi S Nemirovskiy and Boris T Polyak. Iterative methods for solving linear ill-posed problems under precise information. *Eng. Cyber.*, (4):50–56, 1984.
- Donald J Newman. Rational approximation to  $|x|$ . *Michigan Mathematical Journal*, 11(1):11, 1964.
- Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr computational mathematics and mathematical physics*, 4(5):1–17, 1964.
- Aditya Varre and Nicolas Flammarion. Accelerated sgd for non-strongly-convex least squares. *arXiv preprint arXiv:2203.01744*, 2022.
- Aditya Varre, Loucas Pillaud-Vivien, and Nicolas Flammarion. Last iterate convergence of sgd for least-squares in the interpolation regime. *arXiv preprint arXiv:2102.03183*, 2021.
- Maksim Velikanov and Dmitry Yarotsky. Explicit loss asymptotics in the gradient descent training of neural networks. *Advances in Neural Information Processing Systems*, 34, 2021.
- Maksim Velikanov, Denis Kuznedelev, and Dmitry Yarotsky. A view of mini-batch sgd via generating functions: conditions of convergence, phase transitions, benefit from negative momenta. In *The Eleventh International Conference on Learning Representations (ICLR 2023)*, 2023. URL <https://openreview.net/forum?id=bzaPGEllsJJE>.

---

Dmitry Yarotsky. Collective evolution of weights in wide neural networks. *arXiv preprint arXiv:1810.03974*, 2018.

Dmitry Yarotsky and Maksim Velikanov. Sgd with memory: fundamental properties and stochastic acceleration. *arXiv preprint arXiv:2410.04228*, 2024.

## CONTENTS

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Background</b>	<b>2</b>
<b>3</b>	<b>The contour view of generalized (S)GD</b>	<b>5</b>
<b>4</b>	<b>Corner algorithms</b>	<b>7</b>
<b>5</b>	<b>Finite-memory approximations of corner algorithms</b>	<b>8</b>
<b>6</b>	<b>Experiments</b>	<b>9</b>
<b>7</b>	<b>Additional notes and discussion</b>	<b>10</b>
<b>8</b>	<b>Bibliography</b>	<b>12</b>
<b>A</b>	<b>Derivations of Section 3</b>	<b>14</b>
<b>B</b>	<b>Proof of Theorem 3</b>	<b>14</b>
	B.1 The noise propagators . . . . .	14
	B.2 The signal propagators . . . . .	20
<b>C</b>	<b>Proof of Proposition 1</b>	<b>25</b>
<b>D</b>	<b>Proof of Proposition 2</b>	<b>26</b>
<b>E</b>	<b>The synthetic 1D example</b>	<b>26</b>
<b>F</b>	<b>Generalization performance of Corner SGD</b>	<b>28</b>
<b>G</b>	<b>Extending the proof of Theorem 3 to <math>\tau_2 \neq 0</math></b>	<b>29</b>

## A DERIVATIONS OF SECTION 3

We have

$$P(\mu) = \det(\mu - S_0) \quad (41)$$

$$= \det(\mu - S_\lambda + \lambda \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix} (\mathbf{1} \mathbf{o}^T)) \quad (42)$$

$$= \det(\mu - S_\lambda) \det(\mathbf{1} + \lambda \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix} (\mathbf{1} \mathbf{o}^T) (\mu - S_\lambda)^{-1}) \quad (43)$$

$$= (P(\mu) - \lambda Q(\mu)) (\mathbf{1} + \lambda (\mathbf{1} \mathbf{o}^T) (\mu - S_\lambda)^{-1} \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix}). \quad (44)$$

It follows that

$$(\mathbf{1} \mathbf{o}^T) (\mu - S_\lambda)^{-1} \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix} = \frac{1}{\lambda} \left( \frac{P(\mu)}{P(\mu) - \lambda Q(\mu)} - 1 \right) \quad (45)$$

$$= \frac{Q(\mu)}{P(\mu) - \lambda Q(\mu)}. \quad (46)$$

Next, by Sherman-Morrison formula and the above identity,

$$(\mu - S_\lambda)^{-1} = (\mu - S_0 - \lambda \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix} (\mathbf{1} \mathbf{o}^T))^{-1} \quad (47)$$

$$= (\mu - S_0)^{-1} + \lambda \frac{(\mu - S_0)^{-1} \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix} (\mathbf{1} \mathbf{o}^T) (\mu - S_0)^{-1}}{1 - \lambda (\mathbf{1} \mathbf{o}^T) (\mu - S_0)^{-1} \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix}} \quad (48)$$

$$= (\mu - S_0)^{-1} + \lambda \frac{(\mu - S_0)^{-1} \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix} (\mathbf{1} \mathbf{o}^T) (\mu - S_0)^{-1}}{1 - \lambda \frac{Q(\mu)}{P(\mu)}} \quad (49)$$

Using  $(\mathbf{1} \mathbf{o}^T) (\mu - S_0)^{-1} \begin{pmatrix} 1 \\ \mathbf{0} \end{pmatrix} = \frac{1}{\mu - 1}$ , it follows that

$$(\mathbf{1} \mathbf{o}^T) (\mu - S_\lambda)^{-1} \begin{pmatrix} 1 \\ \mathbf{0} \end{pmatrix} = \frac{1}{\mu - 1} + \lambda \frac{\frac{Q(\mu)}{P(\mu)} \frac{1}{\mu - 1}}{1 - \lambda \frac{Q(\mu)}{P(\mu)}} \quad (50)$$

$$= \frac{\frac{P(\mu)}{Q(\mu)}}{(\mu - 1) \left( \frac{P(\mu)}{Q(\mu)} - \lambda \right)}. \quad (51)$$

## B PROOF OF THEOREM 3

### B.1 THE NOISE PROPAGATORS

**The function  $F_U$ .** Recall that by Eq. (21)

$$U(t, \lambda) = \frac{1}{2\pi i} \oint_{|\mu|=1} \frac{\mu^{t-1} d\mu}{\Psi(\mu) - \lambda} = \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{e^{it\phi} d\phi}{\Psi(e^{i\phi}) - \lambda} \quad (52)$$

With the change of variables  $\phi = s\lambda^{1/\theta}$ ,

$$U(t, \lambda) = \frac{-\lambda^{1/\theta-1}}{2\pi} \int_{-\pi/\lambda^{1/\theta}}^{\pi/\lambda^{1/\theta}} \frac{e^{it\lambda^{1/\theta} s} ds}{-\Psi(e^{is\lambda^{1/\theta}})/\lambda + 1} = -\lambda^{1/\theta-1} F_U(t\lambda^{1/\theta}, \lambda), \quad (53)$$

where we have denoted

$$F_U(r, \lambda) = \frac{1}{2\pi} \int_{-\pi/\lambda^{1/\theta}}^{\pi/\lambda^{1/\theta}} \frac{e^{irs} ds}{-\Psi(e^{is\lambda^{1/\theta}})/\lambda + 1}. \quad (54)$$

Recall that we assume  $\Psi(\mu) = -c_\Psi (\mu - 1)^\theta (1 + o(1))$  as  $\mu \rightarrow 1$ . By formally taking the limit  $\lambda \searrow 0$  in the integral, we then expect  $F_U(r, \lambda)$  to converge to

$$F_U(r, 0) \stackrel{\text{def}}{=} F_U(r) \stackrel{\text{def}}{=} \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{e^{irs} ds}{c_\Psi e^{i(\text{sign } s)\theta\pi/2} |s|^\theta + 1} \quad (55)$$

for any fixed  $r$ . This integral can be equivalently written as

$$F_U(r) = \frac{1}{2\pi i} \int_{i\mathbb{R}} \frac{e^{rz} dz}{c_\Psi z^\theta + 1}, \quad (56)$$

assuming the standard branch of  $z^\theta$  holomorphic in  $\mathbb{C} \setminus (-\infty, 0]$ .

The function  $F_U$  can be viewed (up to a coefficient) as the inverse Fourier transform of the function  $s \mapsto (c_\Psi e^{i(\text{sign } s)\theta\pi/2} |s|^\theta + 1)^{-1}$ . Note that, thanks to the condition  $\theta > 1$ , the latter function is Lebesgue-integrable, so  $F_U(r)$  is well-defined and continuous for all  $r \in \mathbb{R}$ . The function  $F_U$  can also be written in terms of the special Mittag-Leffler function  $E_{\theta,\theta}$  (see its integral representation (6.8) in Haubold et al. (2011)):

$$F_U(r) = \frac{r^{\theta-1}}{c_\Psi} E_{\theta,\theta} \left( -\frac{r^\theta}{c_\Psi} \right), \quad E_{a,b}(z) = \frac{1}{2\pi i} \int_\gamma \frac{t^{a-b} e^t dt}{t^a - z}, \quad (57)$$

where the integration path  $\gamma$  encircles the cut  $(-\infty, 0]$  and the singularities of the denominator.

The following asymptotic properties of  $F_U(r)$  can be derived from the general asymptotic expansions of Mittag-Leffler functions (sections 1 and 6 in Haubold et al. (2011)), but we provide proofs for completeness.

**Lemma 1.**

1.  $F_U(r) = 0$  for  $r \leq 0$ .
2.  $F_U(r) = (1 + o(1)) \frac{1}{c_\Psi \Gamma(\theta)} r^{\theta-1}$  as  $r \searrow 0$ .
3.  $F_U(r) = (1 + o(1)) \frac{-c_\Psi}{\Gamma(-\theta)} r^{-\theta-1}$  as  $r \rightarrow +\infty$ .

*Proof.* 1. Consider the function  $f(z)$  integrated in Eq. (56). For any  $r \in \mathbb{R}$  and  $\theta \in (1, 2)$ , the function  $f$  is holomorphic in any strip  $\mathcal{T}_a = \{0 < \Re z < a\}$ ,  $a > 0$ , and is bounded in  $\mathcal{T}_a$  as  $|f(z)| = O(|z|^{-\theta})$ . It follows that the integration line  $i\mathbb{R}$  can be deformed to  $i\mathbb{R} + a$  without changing the integral. If  $r < 0$ , then by letting  $a \rightarrow +\infty$  we can make the integral arbitrarily small.

2. By the change of variables  $rz = z'$ ,

$$F_U(r) = u(r) r^{\theta-1}, \quad (58)$$

where

$$u(r) = \frac{1}{2\pi i c_\Psi} \int_{i\mathbb{R}} \frac{e^{z'} dz'}{z'^\theta + c_\Psi^{-1} r^\theta}. \quad (59)$$

We can find  $\lim_{r \searrow 0} u(r)$  as follows. Observe that the integration line  $i\mathbb{R}$  can be deformed to the line  $\gamma_a$ ,  $a > 0$ , encircling the negative semi-axis:

$$\gamma_a = \gamma_{a,1} \cup \gamma_{a,2} \cup \gamma_{a,3}, \quad (60)$$

$$\gamma_{a,1} = \{z \in \mathbb{C} : \Im z = -a, \Re z \leq 0\}, \quad (61)$$

$$\gamma_{a,2} = \{z \in \mathbb{C} : |z| = a, -\frac{\pi}{2} < \arg z < \frac{\pi}{2}\}, \quad (62)$$

$$\gamma_{a,3} = \{z \in \mathbb{C} : \Im z = a, \Re z \leq 0\}. \quad (63)$$

Indeed, if  $r$  is sufficiently small, then this deformation occurs within the holomorphy domain of the integrated function. The integral is preserved since  $\theta > 0$  and since we deform in the half-plane where the argument of  $e^{z'}$  has  $\Re z' < 0$ .

Thus, for any fixed  $a > 0$  we have

$$\lim_{r \searrow 0} u(r) = \lim_{r \searrow 0} \frac{1}{2\pi i c_\Psi} \int_{\gamma_a} \frac{e^{z'} dz'}{z'^\theta + c_\Psi^{-1} r^\theta} = \frac{1}{2\pi i c_\Psi} \int_{\gamma_a} \frac{e^{z'} dz'}{z'^\theta} = \frac{1}{2\pi i c_\Psi (\theta - 1)} \int_{\gamma_a} \frac{e^{z'} dz'}{z'^{\theta-1}}, \quad (64)$$

where in the last step we integrated by parts. In the last integral, thanks to the weakness of the singularity  $z'^{1-\theta}$  at  $z' = 0$  (note that  $1 - \theta > -1$ ), we can let  $a \rightarrow 0$ :

$$\int_{\gamma_a} \frac{e^{z'} dz'}{z'^{\theta-1}} = \int_0^{+\infty} e^{-s} s^{1-\theta} (e^{-\pi i(1-\theta)} - e^{\pi i(1-\theta)}) ds \quad (65)$$

$$= 2i \sin(\pi(\theta - 1)) \Gamma(2 - \theta) \quad (66)$$

$$= \frac{2\pi i}{\Gamma(\theta - 1)}, \quad (67)$$

where in the last step we used the identity  $\Gamma(z)\Gamma(1-z) = \frac{\pi}{\sin(\pi z)}$ . This is essentially Hankel's representation of the Gamma function, valid for all  $\theta \in \mathbb{C}$  by analytic continuation. Summarizing,

$$\lim_{r \searrow 0} u(r) = \frac{1}{c_\Psi(\theta - 1)\Gamma(\theta - 1)} = \frac{1}{c_\Psi\Gamma(\theta)}. \quad (68)$$

3. We start by performing integration by parts in  $F_U$  :

$$F_U(r) = \frac{-1}{2\pi i r} \int_{i\mathbb{R}} e^{rz} d \frac{1}{c_\Psi z^\theta + 1} = \frac{c_\Psi \theta}{2\pi i r} \int_{i\mathbb{R}} \frac{e^{rz} z^{\theta-1} dz}{(c_\Psi z^\theta + 1)^2}. \quad (69)$$

Performing again the change of variables  $rz = z'$ , we have

$$F_U(r) = v(r)r^{-\theta-1}, \quad (70)$$

where

$$v(r) = \frac{c_\Psi \theta}{2\pi i} \int_{i\mathbb{R}} \frac{e^{z'} z'^{\theta-1} dz'}{(c_\Psi (z'/r)^\theta + 1)^2}. \quad (71)$$

To compute  $\lim_{r \rightarrow \infty} v(r)$ , we again transform the integration line. Let  $\gamma'$  be a line that lies in the domain  $\mathbb{C} \setminus (-\infty, 0)$  and can be represented as the graph of a function  $\Re z = f(\Im z)$  such that

$$f(y) \geq c_1 |y| - c_0 \quad (72)$$

with some constant  $c_1 > 0$  and  $c_0$ .

Note that the integrated function has two singular points  $z' \in \mathbb{C} \setminus (-\infty, 0]$  where the denominator  $c_\Psi (z'/r)^\theta + 1 = 0$ . These two points depend linearly on  $r$ . Require additionally that  $\gamma'$  lie to the right of these points for all  $r > 0$ , so that  $i\mathbb{R}$  can be deformed to  $\gamma'$  without meeting the singularities. This requirement is feasible with a small enough  $c_1 > 0$  since, by the condition  $\theta < 2$ , the imaginary parts of the singular points are negative.

With these assumptions, integration in Eq. (71) can be changed to integration over  $\gamma'$ . Thanks to condition (72), the integrand converges exponentially fast at  $z' \rightarrow \infty$ , and we can take the limit  $r \rightarrow +\infty$  :

$$\lim_{r \rightarrow +\infty} v(r) = \frac{c_\Psi \theta}{2\pi i} \int_{\gamma'} e^{z'} z'^{\theta-1} dz'. \quad (73)$$

The contour  $\gamma'$  can now be transformed to a contour encircling the negative semi-axis, and applying Eq. (67) we get

$$\lim_{r \rightarrow +\infty} v(r) = \frac{c_\Psi \theta}{\Gamma(1-\theta)} = \frac{-c_\Psi}{\Gamma(-\theta)}. \quad (74)$$

□

**The formal leading term in  $U_t$ .** We have

$$U_t = \frac{\tau_1}{|B|} \sum_{k=1}^{\infty} \lambda_k^2 |U(t, \lambda_k)|^2 = \frac{\tau_1}{|B|} \sum_k \lambda_k^{2/\theta} F_U^2(t\lambda_k^{1/\theta}, \lambda_k). \quad (75)$$

To extract the leading term in this expression, we set the second argument in  $F_U(t\lambda_k^{1/\theta}, \lambda_k)$  to 0:

$$U_t^{(1)} \stackrel{\text{def}}{=} \frac{\tau_1}{|B|} \sum_k \lambda_k^{2/\theta} F_U^2(t\lambda_k^{1/\theta}) = \frac{\tau_1}{|B|} a_t t^{\theta/\nu-2}, \quad (76)$$

where

$$a_t = t^{2-\theta/\nu} \sum_k \lambda_k^{2/\theta} F_U^2(t\lambda_k^{1/\theta}) = t^{-\theta/\nu} \sum_k (t\lambda_k^{1/\theta})^2 F_U^2(t\lambda_k^{1/\theta}). \quad (77)$$

**Lemma 2.**

$$\lim_{t \rightarrow \infty} a_t = \Lambda^{1/\nu} \int_0^{\infty} r^2 F_U^2(r) dr^{-\theta/\nu} = \Lambda^{1/\nu} \frac{\theta}{\nu} \int_0^{\infty} r^{1-\theta/\nu} F_U^2(r) dr < \infty. \quad (78)$$

*Proof.* Note first that the integral on the right is convergent. Indeed, by statement 2 of Lemma 1,  $r^{1-\theta/\nu} F^2(r) \propto r^{1-\theta/\nu+2(\theta-1)} = r^{\theta(2-1/\nu)-1}$  near  $r = 0$ . Since we assume  $\nu > 1$  and  $\theta > 1$ , the function  $r^{1-\theta/\nu} F^2(r)$  is bounded near  $r = 0$ . Also, by statement 3 of Lemma 1,  $r^{1-\theta/\nu} F^2(r) \propto r^{1-\theta/\nu-2(\theta+1)} = O(r^{-3})$  as  $r \rightarrow +\infty$ .

For any interval  $I$  in  $\mathbb{R}_+$ , denote by  $S_{I,t}$  the part of the expansion (77) of  $a_t$  corresponding to the terms with  $t\lambda_k^{1/\theta} \in I$ :

$$S_{I,t} = t^{-\theta/\nu} \sum_{k: t\lambda_k^{1/\theta} \in I} (t\lambda_k^{1/\theta})^2 F_U^2(t\lambda_k^{1/\theta}). \quad (79)$$

Recall that the eigenvalues  $\lambda$  are ordered and  $\lambda_k = \Lambda k^{-\nu}(1 + o(1))$  by capacity condition (11). It follows that for a given fixed number  $r > 0$ , the condition  $t\lambda_k^{1/\theta} > r$  holds whenever  $k < k_r$ , where

$$k_r = (1 + o(1))\Lambda^{1/\nu}(t/r)^{\theta/\nu}, \quad t \rightarrow \infty. \quad (80)$$

Then, for  $I = [u, v]$  with  $0 < u < v < \infty$  we have

$$\liminf_{t \rightarrow \infty} S_{I,t} \geq \Lambda^{1/\nu} \inf_{r \in I} [r^2 F_U^2(r)](u^{-\theta/\nu} - v^{-\theta/\nu}), \quad (81)$$

$$\limsup_{t \rightarrow \infty} S_{I,t} \leq \Lambda^{1/\nu} \sup_{r \in I} [r^2 F_U^2(r)](u^{-\theta/\nu} - v^{-\theta/\nu}). \quad (82)$$

Moreover, for any interval  $I = [u, v]$  with  $0 < u < v < \infty$  we can approximate  $\int_I r^2 F_U^2(r) dr^{-\theta/\nu}$  by integral sums corresponding to sub-divisions  $I = I_1 \cup I_2 \cup \dots \cup I_n$ , apply the above inequalities to each  $I_s$ , and conclude that

$$\lim_{t \rightarrow \infty} S_{I,t} = \Lambda^{1/\nu} \int_I r^2 F_U^2(r) dr^{-\theta/\nu}. \quad (83)$$

It remains to handle the two parts of  $a_t$  corresponding to the remaining intervals  $I = [0, u]$  and  $I = [v, \infty)$ . It suffices to show that the associated contributions  $S_{I,t}$  can be made arbitrarily small uniformly in  $t$  by making  $u$  small and  $v$  large enough.

Consider first the interval  $I = [v, \infty)$ . Note that by Lemma 1 for all  $r > 1$  we can write

$$r^2 F_U^2(r) \leq Cr^{-2\theta} \quad (84)$$

with some constant  $C$ , and we also have for all  $k$

$$\Lambda_- k^{-\nu} \leq \lambda_k \leq \Lambda_+ k^{-\nu} \quad (85)$$

for suitable constants  $\Lambda_-$ ,  $\Lambda_+$ . It follows that

$$S_{I,t} \leq t^{-\theta/\nu} \sum_{k: t(\Lambda_+ k^{-\nu})^{1/\theta} > v} C(t(\Lambda_- k^{-\nu})^{1/\theta})^{-2\theta} \quad (86)$$

$$= t^{-\theta/\nu-2\theta} C \Lambda_-^{-2} \sum_{k=1}^{\Lambda_+^{1/\nu}(t/v)^{\theta/\nu}} k^{2\nu} \quad (87)$$

$$= O(1)t^{-\theta/\nu-2\theta}(t/v)^{(\theta/\nu)(2\nu+1)} \quad (88)$$

$$= O(1)v^{-(\theta/\nu)(2\nu+1)}, \quad (89)$$

with  $O(1)$  denoting an expression bounded by a  $t, v$ -independent constant. This is the desired convergence property of  $S_{I,t}$ .

Similarly, for the other interval  $I = [0, u]$  we use the inequality

$$r^2 F_U^2(r) \leq Cr^{2\theta}, \quad r < 1, \quad (90)$$

also following by Lemma 1. Then

$$S_{I,t} \leq t^{-\theta/\nu} \sum_{k:t(\Lambda_- k^{-\nu})^{1/\theta} < u} C(t(\Lambda_+ k^{-\nu})^{1/\theta})^{2\theta} \quad (91)$$

$$= t^{-\theta/\nu+2\theta} C\Lambda_+^2 \sum_{k=\Lambda_-^{1/\nu}(t/u)^{\theta/\nu}}^{\infty} k^{-2\nu} \quad (92)$$

$$= O(1)t^{-\theta/\nu+2\theta}(t/u)^{(\theta/\nu)(1-2\nu)} \quad (93)$$

$$= O(1)u^{(\theta/\nu)(2\nu-1)}, \quad (94)$$

which is the desired convergence property of  $S_{I,t}$  since  $\nu > 1$ .  $\square$

**Completion of proof.** We have shown that if we replace  $F_U(t\lambda_k^{1/\theta}, \lambda_k)$  by  $F_U(t\lambda_k^{1/\theta})$  in Eq. (75), we get desired asymptotics of  $U_t$  in the limit  $t \rightarrow +\infty$ . We will show now that this replacement introduces a lower-order correction  $o(t^{\theta/\nu-2})$ ; this will complete the proof.

We start with a technical lemma (to be applied with  $f = \Psi$ ) giving a lower bound for deviations of asymptotic power law functions with  $\theta < 2$  from real values.

**Lemma 3.** *Suppose that  $f : \{\mu \in \mathbb{C} : |\mu| = 1\} \rightarrow \mathbb{C}$  is continuous,  $f(\mu) = -c(\mu-1)^\theta(1+o(1))$  as  $\mu \rightarrow 1$  with some  $\theta \in [0, 2)$  and  $c > 0$ . Suppose also that  $f(\{\mu \in \mathbb{C} : |\mu| = 1, \mu \neq 1\}) \cap [0, \lambda_{\max}] = \emptyset$  for some  $\lambda_{\max} > 0$ . Then there exist a constant  $C > 0$  such that*

$$|f(e^{is}) - \lambda| \geq C(|s|^\theta + \lambda), \quad s \in [-\pi, \pi], \lambda \in [0, \lambda_{\max}]. \quad (95)$$

*Proof.* If we fix any small  $\epsilon > 0$ , then, by the condition  $f(\{\mu \in \mathbb{C} : |\mu| = 1, \mu \neq 1\}) \cap [0, \lambda_{\max}] = \emptyset$  and a compactness argument, there exist  $C', C > 0$  such that

$$|f(e^{is}) - \lambda| > C' > C(|s|^\theta + \lambda), \quad s \in [-\pi, -\epsilon] \cup [\epsilon, \pi], \lambda \in [0, \lambda_{\max}]. \quad (96)$$

It remains to establish inequality (95) for  $|s| < \epsilon$ . Since  $f(\mu) = c(\mu-1)^\theta(1+o(1))$  and  $\theta \in [0, 2)$ ,

$$|f(e^{is}) - \lambda| = |e^{i \operatorname{sign}(s)\theta\pi/2} c|s|^\theta(1+o(1)) + \lambda| \quad (97)$$

$$= |e^{i \operatorname{sign}(s)\theta\pi/4} c|s|^\theta(1+o(1)) + \lambda e^{-i \operatorname{sign}(s)\theta\pi/4}| \quad (98)$$

$$\geq \Re[e^{i \operatorname{sign}(s)\theta\pi/4} c|s|^\theta(1+o(1)) + \lambda e^{-i \operatorname{sign}(s)\theta\pi/4}] \quad (99)$$

$$= \cos(\theta\pi/4)(c|s|^\theta(1+o(1)) + \lambda) \quad (100)$$

$$\geq \frac{1}{2} \min(c, 1) \cos(\theta\pi/4)(|s|^\theta + \lambda) \quad (101)$$

for  $|s|$  small enough.  $\square$

**Lemma 4.**

1.  $|F_U(r, \lambda) - F_U(r)| = o(1)$  as  $\lambda \rightarrow 0$ , uniformly in all  $r \in \mathbb{R}$ .
2.  $F_U(r, \lambda) = O(\frac{1}{r})$  for all  $r$  of the form  $r = t\lambda^{1/\theta}$ ,  $t = 1, 2, \dots$ , uniformly in all  $\lambda \in (0, \lambda_{\max}]$ .

*Proof.* 1. It suffices to show that, as  $\lambda \searrow 0$ , the functions

$$f_\lambda(s) = -(2\pi)^{-1}(-\Psi(e^{is\lambda^{1/\theta}})/\lambda + 1)^{-1} \mathbf{1}_{[-\pi/\lambda^{1/\theta}, \pi/\lambda^{1/\theta}]}(s) \quad (102)$$

converge in  $L^1(\mathbb{R})$  to

$$f_0(s) = -(2\pi)^{-1}(c_\Psi e^{i(\operatorname{sign} s)\theta\pi/2}|s|^\theta + 1)^{-1}. \quad (103)$$

Let us divide the interval  $[-\pi/\lambda^{1/\theta}, \pi/\lambda^{1/\theta}]$  into two subsets:

$$I_1(\lambda) = [-\lambda^{-h}, \lambda^{-h}], \quad (104)$$

$$I_2(\lambda) = [-\pi/\lambda^{1/\theta}, \pi/\lambda^{1/\theta}] \setminus I_1(\lambda), \quad (105)$$

where  $h$  is some fixed number such that  $\frac{1}{\theta^2} < h < \frac{1}{\theta}$ .

By Lemma 3,  $|\Psi(e^{is\lambda^{1/\theta}})/\lambda - 1| \geq c|s|^\theta$  uniformly for all  $s \in [-\pi/\lambda^{1/\theta}, \pi/\lambda^{1/\theta}]$  and  $\lambda \in (0, \lambda_{\max}]$ . It follows that

$$\inf_{s \in I_2(\lambda)} |\Psi(e^{is\lambda^{1/\theta}})/\lambda - 1| \geq c\lambda^{-h\theta}, \quad \lambda \in (0, \lambda_{\max}], \quad (106)$$

for some constant  $c > 0$ . Using the condition  $\frac{1}{\theta^2} < h$ , it follows that

$$\int_{I_2(\lambda)} |f_\lambda(s)| ds = O(\lambda^{-1/\theta} \lambda^{h\theta}) = o(1), \quad \lambda \searrow 0. \quad (107)$$

Thus, we can assume without loss that the functions  $f_\lambda$  vanish outside the intervals  $I_1(\lambda)$ . On these intervals, thanks to the condition  $h < \frac{1}{\theta}$ , we have

$$f_\lambda(s) = -(2\pi)^{-1} (c_\Psi e^{i(\text{sign } s)\theta\pi/2} |s|^\theta (1 + o(1)) + 1)^{-1} \quad (108)$$

uniformly in  $s \in I_1(\lambda)$ . We can then apply the dominated convergence theorem to the functions  $|f_\lambda - f_0|$ , with a dominating function  $C(1 + |s|^\theta)^{-1}$ , and conclude that  $f_\lambda \rightarrow f_0$  in  $L^1(\mathbb{R})$ , as desired.

2. We start by performing integration by parts in  $U(t, \lambda)$ :

$$U(t, \lambda) = \frac{1}{2\pi it} \oint_{|\mu|=1} \frac{d\mu^t}{\Psi(\mu) - \lambda} = \frac{1}{2\pi it} \oint_{|\mu|=1} \frac{\Psi'(\mu)\mu^t d\mu}{(\Psi(\mu) - \lambda)^2} \quad (109)$$

implying

$$|U(t, \lambda)| \leq \frac{1}{2\pi t} \int_{-\pi}^{\pi} \frac{|\Psi'(e^{is})| ds}{|\Psi(e^{is}) - \lambda|^2}. \quad (110)$$

We will show that this integral is  $O(\frac{1}{\lambda})$ .

Note first that we can replace the integration on  $[-\pi, \pi]$  by integration on  $[-a, a]$  for any  $0 < a < \pi$ . Indeed, by our assumptions  $\Psi$  is  $C^1$  on the unit circle, and  $\Psi(\mu) = 1$  there only if  $\mu = 1$ . Accordingly, the remaining part of the integral is non-singular as  $\lambda \searrow 0$  and so is uniformly bounded for all  $\lambda \in (0, \lambda_{\max}]$ .

Recall that by our assumption  $\Psi'(\mu) = O(|\mu - 1|^{\theta-1})$  as  $\mu \rightarrow 1$ . Applying again Lemma 3,

$$|U(t, \lambda)| \leq \frac{C'}{t} \int_0^\infty \frac{s^{\theta-1} ds}{(s^\theta + \lambda)^2} = \frac{C'}{t\lambda} \quad (111)$$

with some constant  $C'$  independent of  $t, \lambda$ . It follows that

$$|F_U(t\lambda^{1/\theta}, \lambda)| = |\lambda^{1-1/\theta} U(t, \lambda)| \leq \frac{C'}{t\lambda^{1/\theta}}, \quad (112)$$

as claimed.  $\square$

We return now to proving that replacing  $F_U(t\lambda_k^{1/\theta}, \lambda_k)$  by  $F_U(t\lambda_k^{1/\theta})$  in Eq. (75) amounts to a lower-order correction  $o(t^{\theta/\nu-2})$ . It suffices to prove that  $\Delta a_t \rightarrow 0$ , where

$$\Delta a_t = t^{2-\theta/\nu} \sum_k \lambda_k^{2/\theta} (F_U^2(t\lambda_k^{1/\theta}, \lambda_k) - F_U^2(t\lambda_k^{1/\theta})) \quad (113)$$

$$= t^{-\theta/\nu} \sum_k (t\lambda_k^{1/\theta})^2 (F_U^2(t\lambda_k^{1/\theta}, \lambda_k) - F_U^2(t\lambda_k^{1/\theta})). \quad (114)$$

For any interval  $I \subset \mathbb{R}$ , denote by  $\Delta S_{I,t}$  the part of  $\Delta a_t$  corresponding to the terms in (114) such that  $t\lambda_k^{1/\theta} \in I$ . By statement 1 of Lemma 4, for any  $u > 0$  we have, as  $t \rightarrow \infty$ ,

$$|\Delta S_{(0,u),t}| = o(1)t^{2-\theta/\nu} \sum_{k:t\lambda_k^{1/\theta} < u} \lambda_k^{2/\theta} \quad (115)$$

$$= o(1)t^{2-\theta/\nu} O((t/u)^{(\theta/\nu)(1-2\nu/\theta)}) \quad (116)$$

$$= o(1), \quad (117)$$

where we have used the fact that  $2\nu/\theta > \nu > 1$ .

Now consider the remaining interval  $I = [u, +\infty)$ . It suffices to prove that  $|\Delta S_{[u, +\infty), t}|$  can be made arbitrarily small uniformly in  $t$  by choosing  $u$  large enough. By statement 2 of Lemma 4, we can write

$$|\Delta S_{[u, +\infty), t}| \leq C t^{2-\theta/\nu} \sum_{k: t\lambda_k^{1/\theta} > u} \lambda_k^{2/\theta} (t\lambda_k^{1/\theta})^{-2} \quad (118)$$

$$\leq C t^{-\theta/\nu} \sum_{k=1}^{\Lambda_+^{1/\nu} (t/u)^{\theta/\nu}} 1 \quad (119)$$

$$\leq C' u^{-\theta/\nu} \quad (120)$$

with some  $t, u$ -independent constant  $C'$ . This completes the proof of statement 1 of Theorem 3.

## B.2 THE SIGNAL PROPAGATORS

The proof for the signal propagators follows the same ideas as for the noise propagators, with appropriate adjustments.

**The function  $F_V$ .** Recall

$$V(t, \lambda) = \frac{1}{2\pi i} \oint_{|\mu|=1} \frac{\Psi(\mu)\mu^{t-1}d\mu}{(\Psi(\mu) - \lambda)(\mu - 1)} = \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{\Psi(e^{i\phi})e^{it\phi}d\phi}{(\Psi(e^{i\phi}) - \lambda)(e^{i\phi} - 1)}. \quad (121)$$

With the change of variables  $\phi = s\lambda^{1/\theta}$ ,

$$V(t, \lambda) = \frac{\lambda^{1/\theta}}{2\pi} \int_{-\pi/\lambda^{1/\theta}}^{\pi/\lambda^{1/\theta}} \frac{(-\Psi(e^{is\lambda^{1/\theta}})/\lambda)e^{it\lambda^{1/\theta}s}ds}{(-\Psi(e^{is\lambda^{1/\theta}})/\lambda + 1)(e^{is\lambda^{1/\theta}} - 1)} = F_V(t\lambda^{1/\theta}, \lambda), \quad (122)$$

where

$$F_V(r, \lambda) = \frac{\lambda^{1/\theta}}{2\pi} \int_{-\pi/\lambda^{1/\theta}}^{\pi/\lambda^{1/\theta}} \frac{(-\Psi(e^{is\lambda^{1/\theta}})/\lambda)e^{irs}ds}{(-\Psi(e^{is\lambda^{1/\theta}})/\lambda + 1)(e^{is\lambda^{1/\theta}} - 1)}. \quad (123)$$

We again recall that  $\Psi(\mu) = -c_\Psi(\mu - 1)^\theta(1 + o(1))$  as  $\mu \rightarrow 1$  and formally take the pointwise limit  $\lambda \searrow 0$  in the integrand to obtain the expression

$$F_V(r, 0) \stackrel{\text{def}}{=} F_V(r) \stackrel{\text{def}}{=} \frac{1}{2\pi i} \int_{-\infty}^{\infty} \frac{c_\Psi e^{i(\text{sign } s)\theta\pi/2} |s|^\theta e^{irs} ds}{(c_\Psi e^{i(\text{sign } s)\theta\pi/2} |s|^\theta + 1)s} \quad (124)$$

$$= \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{c_\Psi e^{i(\text{sign } s)(\theta-1)\pi/2} |s|^{\theta-1} e^{irs} ds}{(c_\Psi e^{i(\text{sign } s)\theta\pi/2} |s|^\theta + 1)} \quad (125)$$

for any fixed  $r$ . This integral can be equivalently written as

$$F_V(r) = \frac{1}{2\pi i} \int_{i\mathbb{R}} \frac{c_\Psi z^{\theta-1} e^{rz} dz}{c_\Psi z^\theta + 1}, \quad (126)$$

assuming again the standard branch of  $z^\theta$  holomorphic in  $\mathbb{C} \setminus (-\infty, 0]$ . The function  $F_V$  can be written in terms of the Mittag-Leffler function  $E_\theta \equiv E_{\theta,1}$  (the special case of  $E_{a,b}$  given by Eq. (57)):

$$F_V(r) = E_\theta\left(-\frac{r^\theta}{c_\Psi}\right). \quad (127)$$

Note that, in contrast to  $F_U$ , the integrals (125), (126) are not absolutely summable, due to the  $z^{-1}$  fall off of the integrand at  $z \rightarrow \infty$ . However, the integrand is square-summable and so  $F_V$ , as a Fourier transform of such function, is well-defined almost everywhere as a square-integrable function.

In fact,  $F_V$  can be defined for each particular  $r \neq 0$  by restricting the integration in (125) to segments  $[u, v]$  and letting  $u \rightarrow -\infty$  and  $v \rightarrow \infty$ . Indeed, the resulting Fourier transforms  $F_V^{(u,v)}$  converge

to  $F_V$  in  $L^2(\mathbb{R})$ . However, these transforms are continuous functions of  $r$ , and as  $u \rightarrow \infty, v \rightarrow \infty$  they converge pointwise, and even uniformly on the sets  $\{r : |r| > \epsilon\}$ , for any fixed  $\epsilon > 0$ .

To see this last property of uniform pointwise convergence, note that the integrand in (125) has the form  $(s^{-1} + O(s^{-1-\theta}))e^{irs}$  as  $s \rightarrow \infty$ . The component  $O(s^{-1-\theta})$  is in  $L^1$ , so the respective part of  $F_V^{(u,v)}$  converges as  $u \rightarrow -\infty, v \rightarrow \infty$  uniformly for all  $r \in \mathbb{R}$ . Regarding the  $s^{-1}$  component, integrating by parts gives

$$\int_1^v \frac{e^{irs} ds}{s} = \frac{e^{irs}}{irs} \Big|_{s=1}^v + \frac{1}{ir} \int_1^v \frac{e^{irs} ds}{s^2}. \quad (128)$$

This expression converges as  $v \rightarrow \infty$  uniformly for  $\{r : |r| > \epsilon\}$  with any fixed  $\epsilon > 0$ , as claimed. The same argument applies to  $\int_u^{-1}$ .

The above argument shows, in particular, that  $F_V$  is naturally defined as a function continuous on the intervals  $(0, +\infty)$  and  $(-\infty, 0)$ .

We collect further properties of  $F_V(r)$  in the following lemma that parallels Lemma 1 for  $F_U$ . The proofs are also similar to the proofs in Lemma 1.

**Lemma 5.**

1.  $F_V(r) = 0$  for  $r < 0$ .
2.  $F_V(r) \rightarrow 1$  as  $r \searrow 0$ .
3.  $F_V(r) = (1 + o(1)) \frac{c_\Psi}{\Gamma(1-\theta)} r^{-\theta}$  as  $r \rightarrow +\infty$ .

*Proof.* 1. Like in Lemma 1, this follows by deforming the integration line in Eq. (126) towards  $+\infty$ .

2. By the change of variables  $rz = z'$ ,

$$F_V(r) = \frac{1}{2\pi i} \int_{i\mathbb{R}} \frac{z'^{\theta-1} e^{z'} dz'}{z'^{\theta} + c_\Psi^{-1} r^\theta}. \quad (129)$$

As in Lemma 1, the integration line  $i\mathbb{R}$  can be deformed to the line  $\gamma_a, a > 0$ , encircling the negative semi-axis:

$$\gamma_a = \gamma_{a,1} \cup \gamma_{a,2} \cup \gamma_{a,3}, \quad (130)$$

$$\gamma_{a,1} = \{z \in \mathbb{C} : \Im z = -a, \Re z \leq 0\}, \quad (131)$$

$$\gamma_{2,2} = \{z \in \mathbb{C} : |z| = a, -\frac{\pi}{2} < \arg z < \frac{\pi}{2}\}, \quad (132)$$

$$\gamma_{a,1} = \{z \in \mathbb{C} : \Im z = a, \Re z \leq 0\}. \quad (133)$$

Taking the limit  $r \searrow 0$ , we get

$$\lim_{r \searrow 0} F_V(r) = \lim_{r \searrow 0} \frac{1}{2\pi i} \int_{\gamma_a} \frac{z'^{\theta-1} e^{z'} dz'}{z'^{\theta} + c_\Psi^{-1} r^\theta} = \frac{1}{2\pi i} \int_{\gamma_a} \frac{e^{z'} dz'}{z'} = 1, \quad (134)$$

since the last integral simply amounts to the residue of  $e^{z'}/z'$  at  $z' = 0$ .

3. Using the same contour  $\gamma'$  as in Lemma 1,

$$F_V(r) = v(r)r^{-\theta}, \quad v(r) = \frac{1}{2\pi i} \int_{\gamma'} \frac{c_\Psi z'^{\theta-1} e^{z'} dz'}{c_\Psi (z'/r)^\theta + 1}. \quad (135)$$

Taking the limit  $r \rightarrow +\infty$  and deforming the contour to the negative semi-axis as in Lemma 1,

$$\lim_{r \rightarrow +\infty} v(r) = \frac{c_\Psi}{2\pi i} \int_{\gamma'} z'^{\theta-1} e^{z'} dz' = \frac{c_\Psi}{\Gamma(1-\theta)}. \quad (136)$$

□

**The formal leading term in  $V_t$ .** We have

$$V_t = \sum_{k=1}^{\infty} \lambda_k (\mathbf{e}_k^T \mathbf{w}_*)^2 |V(t, \lambda_k)|^2 = \sum_k \lambda_k (\mathbf{e}_k^T \mathbf{w}_*)^2 F_V^2(t\lambda_k^{1/\theta}, \lambda_k). \quad (137)$$

To extract the leading term in this expression, we set the second argument in  $F_V(t\lambda_k^{1/\theta}, \lambda_k)$  to 0:

$$V_t^{(1)} \stackrel{\text{def}}{=} \sum_k \lambda_k (\mathbf{e}_k^T \mathbf{w}_*)^2 F_V^2(t\lambda_k^{1/\theta}) = b_t t^{-\theta\zeta}, \quad (138)$$

where

$$b_t = t^{\theta\zeta} \sum_k \lambda_k (\mathbf{e}_k^T \mathbf{w}_*)^2 F_V^2(t\lambda_k^{1/\theta}). \quad (139)$$

The analog of Lemma 2 is

**Lemma 6.**

$$\lim_{t \rightarrow \infty} b_t = Q \int_0^{\infty} F_V^2(r) dr^{\theta\zeta} = Q\theta\zeta \int_0^{\infty} r^{\theta\zeta-1} F_V^2(r) dr < \infty. \quad (140)$$

*Proof.* First, observe that, by the source condition (12) and Lemma 5, the integral converges near  $r = 0$  since  $\theta\zeta > 0$ , and near  $r = \infty$  since  $\zeta < 2$ .

We can establish convergence of the sequence  $b_t$  using the same steps as in Lemma 2. We first introduce the sums  $S_{I,t}$  comprising the terms of expansion (139) such that  $t\lambda_k^{1/\theta} \in I$ . For intervals  $I = [u, v]$  with  $0 < u < v < \infty$  we show, using the source condition (12) and approximation by integral sums, that

$$\lim_{t \rightarrow \infty} S_{I,t} = t^{\theta\zeta} \int_I F_V^2(r) dQ((r/t)^\theta)^\zeta = Q \int_I F_V^2(r) dr^{\theta\zeta}. \quad (141)$$

After that we show that the contribution of the remaining intervals  $(v, +\infty)$  and  $(0, u)$  can be made arbitrarily small uniformly in  $t$  by adjusting  $u, v$ .

In particular, consider the interval  $I = (v, +\infty)$ . Let  $R(\lambda) = \sum_{k:\lambda_k \leq \lambda} \lambda_k (\mathbf{e}_k^T \mathbf{w}_*)^2$  denote the cumulative distribution function of the spectral measure. Since the spectral measure is compactly supported, assumption (12) implies that  $R(\lambda) \leq Q'\lambda^\zeta$  for all  $\lambda > 0$  with some  $Q' > 0$ . Using statement 3 of Lemma 5 and integration by parts, we can bound

$$S_{(v,+\infty),t} \leq t^{\theta\zeta} \sum_{k:t\lambda_k^{1/\theta} > v} \lambda_k (\mathbf{e}_k^T \mathbf{w}_*)^2 C(t\lambda_k^{1/\theta})^{-2\theta} \quad (142)$$

$$= Ct^{\theta(\zeta-2)} \int_{(v/t)^\theta}^{\infty} \frac{dR(\lambda)}{\lambda^2} \quad (143)$$

$$= Ct^{\theta(\zeta-2)} \left( \frac{R(\lambda)}{\lambda^2} \Big|_{(v/t)^\theta}^{\infty} + 2 \int_{(v/t)^\theta}^{\infty} \frac{R(\lambda) d\lambda}{\lambda^3} \right) \quad (144)$$

$$\leq 2CQ't^{\theta(\zeta-2)} \int_{(v/t)^\theta}^{\infty} \lambda^{\zeta-3} d\lambda \quad (145)$$

$$\leq C'v^{(\zeta-2)\theta} \quad (146)$$

with some constant  $C'$  independent of  $v, t$ .

For the intervals  $I = (0, u)$  we have

$$S_{(0,u),t} \leq t^{\theta\zeta} \sum_{k:t\lambda_k^{1/\theta} < u} \lambda_k (\mathbf{e}_k^T \mathbf{w}_*)^2 C \quad (147)$$

$$\leq Ct^{\theta\zeta} Q((u/t)^\theta)^\zeta \quad (148)$$

$$= C'u^{\theta\zeta}. \quad (149)$$

□

**Completion of proof.** It remains to show that the correction in  $V_t$  due to the replacement of  $F_V(t\lambda_k^{1/\theta}, \lambda_k)$  by  $F_V(t\lambda_k^{1/\theta})$  in Eq. (137) is  $o(t^{-\theta\zeta})$ . We first establish an analog of Lemma 4:

**Lemma 7.** *Assuming that  $r = t\lambda^{1/\theta}$  with  $t = 1, 2, \dots$ :*

1.  $|F_V(r, \lambda) - F_V(r)| = o(1)$  as  $\lambda \rightarrow 0$ , uniformly for  $r > \epsilon$ , for any  $\epsilon > 0$ .
2.  $|F_V(r, \lambda)| \leq C \min(\frac{1}{r}, 1)$  for all  $t = 1, 2, \dots$  and  $\lambda \in (0, \lambda_{\max}]$ , with some  $r, \lambda$ -independent constant  $C$ .

*Proof.* 1. The proof of this property is more complicated than the earlier proof for  $F_U$  because the integrals defining  $F_V$  are not absolutely convergent. Recall the integration by parts argument (128) used to define  $F_V(r)$  as the pointwise limit of the functions  $F_V^{(u,v)}(r)$ . We extend this approach to the functions  $F_V(r, \lambda)$  with  $\lambda > 0$ . Specifically, let  $F_V^{(u)}(r, \lambda)$  be defined as  $F_V(r, \lambda)$  in Eq. (123), but with integration restricted to the segment  $[-u, u]$ . By analogy with our convention  $F_V(r) \equiv F_V(r, \lambda = 0)$ , denote also  $F_V^{(u)}(r) \equiv F_V^{(u)}(r, \lambda = 0)$ . We will establish the following two properties:

- (a)  $|F_V^{(u)}(r, \lambda) - F_V(r, \lambda)| \leq \frac{C}{ru}$  for all  $0 < \lambda < \lambda_{\max}$  with a  $r, u, \lambda$ -independent constant  $C$ .
- (b) For any  $u$ ,  $|F_V^{(u)}(r, \lambda) - F_V^{(u)}(r)| \rightarrow 0$  as  $\lambda \searrow 0$  uniformly for  $r \in \mathbb{R}$ .

Observe first that these two properties imply the claimed uniform convergence  $|F_V(r, \lambda) - F_V(r)| = o(1)$  as  $\lambda \rightarrow 0$ . Indeed, given any  $\delta > 0$ , first set  $u = \frac{3C}{\epsilon}$  so that by (a) we have

$$|F_V^{(u)}(r, \lambda) - F_V(r, \lambda)| \leq \delta/3 \quad (150)$$

for all  $r > \epsilon$  and  $0 < \lambda < \lambda_{\max}$ . This inequality also holds in the limit  $\lambda \searrow 0$ , i.e.

$$|F_V^{(u)}(r) - F_V(r)| \leq \delta/3. \quad (151)$$

Now (b) implies that for sufficiently small  $\lambda$  we have

$$|F_V^{(u)}(r, \lambda) - F_V^{(u)}(r)| \leq \delta/3 \quad (152)$$

uniformly in  $r \in \mathbb{R}$ . Combining all three above inequalities, we see that for sufficiently small  $\lambda$

$$|F_V(r, \lambda) - F_V(r)| \leq \delta \quad (153)$$

uniformly for  $r > \epsilon$ , as desired.

It remains to prove the statements (a) and (b). Statement (b) immediately follows from the uniform  $\lambda \searrow 0$  convergence of the integrand in expression (123) on the interval  $s \in [-u, u]$ .

To prove statement (a), we perform integration by parts, using the  $\frac{2\pi}{\lambda^{1/\theta}}$ -periodicity of the integrand:

$$|F_V^{(u)}(r, \lambda) - F_V(r, \lambda)| \quad (154)$$

$$= \frac{\lambda^{1/\theta}}{2\pi} \left| \int_{[-\frac{\pi}{\lambda^{1/\theta}}, \frac{\pi}{\lambda^{1/\theta}}] \setminus [-u, u]} \frac{(\Psi(e^{is\lambda^{1/\theta}})/\lambda)e^{irs} ds}{(\Psi(e^{is\lambda^{1/\theta}})/\lambda - 1)(e^{is\lambda^{1/\theta}} - 1)} \right| \quad (155)$$

$$= \frac{\lambda^{1/\theta}}{2\pi r} \left| \frac{(\Psi(e^{is\lambda^{1/\theta}})/\lambda)e^{irs}}{(\Psi(e^{is\lambda^{1/\theta}})/\lambda - 1)(e^{is\lambda^{1/\theta}} - 1)} \Big|_{s=u}^{-u} - \int_{[-\frac{\pi}{\lambda^{1/\theta}}, \frac{\pi}{\lambda^{1/\theta}}] \setminus [-u, u]} \frac{i\lambda^{1/\theta} [(-\Psi'(e^{is\lambda^{1/\theta}})/\lambda)(e^{is\lambda^{1/\theta}} - 1) - (\Psi(e^{is\lambda^{1/\theta}})/\lambda)(\Psi(e^{is\lambda^{1/\theta}})/\lambda - 1)e^{irs}]}{(\Psi(e^{is\lambda^{1/\theta}})/\lambda - 1)^2 (e^{is\lambda^{1/\theta}} - 1)^2} ds \right|. \quad (156)$$

By our assumptions on  $\Psi$ , Lemma 3 and standard inequalities, there exist  $\lambda, s$ -independent constants  $C, c > 0$  such that for all  $\lambda \in (0, \lambda_{\max}]$  and  $s \in [-\frac{\pi}{\lambda^{1/\theta}}, \frac{\pi}{\lambda^{1/\theta}}]$

$$|\Psi(e^{is\lambda^{1/\theta}})| \leq C|s|^\theta \lambda, \quad (157)$$

$$|\Psi'(e^{is\lambda^{1/\theta}})| \leq C\theta|s|^{\theta-1}\lambda^{(\theta-1)/\theta}, \quad (158)$$

$$|\Psi(e^{is\lambda^{1/\theta}})/\lambda - 1| \geq c(1 + |s|^\theta), \quad (159)$$

$$|e^{is\lambda^{1/\theta}} - 1| \geq c|s|\lambda^{1/\theta}. \quad (160)$$

Applying these inequalities to Eq. (156), we find that

$$|F_V^{(u)}(r, \lambda) - F_V(r, \lambda)| \leq \frac{C'}{r} \left( \frac{u^\theta}{(1+u^\theta)u} + \int_{[-\frac{\pi}{\lambda^{1/\theta}}, \frac{\pi}{\lambda^{1/\theta}}] \setminus [-u, u]} \frac{|s|^\theta ds}{(1+|s|^\theta)s^2} \right) \quad (161)$$

$$\leq \frac{C''}{ru}, \quad (162)$$

as desired.

2. Note that

$$|F_V(r, \lambda)| \leq \frac{C}{r}, \quad C < \infty, \quad (163)$$

simply by setting  $u = 0$  in the bound (161), since the first term on the r.h.s. of (161) vanishes and the second converges thanks to  $\theta > 1$ .

It remains to prove that  $F_V(r, \lambda)$  is bounded uniformly in  $r, \lambda$ . It suffices to prove this for  $r < \epsilon$  with some fixed  $\epsilon > 0$ , since for larger  $r$  this follows from bound (163). Since  $r = t\lambda^{1/\theta}$ , this means it is sufficient to consider

$$\lambda \leq (\epsilon/t)^\theta. \quad (164)$$

To this end consider the original definition (121) of  $V(t, \lambda)$  in terms of integration over the contour  $\{|\mu| = 1\}$ . We will deform this contour within the analyticity domain  $\{\mu \in \mathbb{C} : |\mu| \geq 1\}$  to another contour  $\gamma$ , to be specified below, that fully encircles the point  $\mu = 1$ :

$$V(t, \lambda) = \frac{1}{2\pi i} \oint_\gamma \frac{\Psi(\mu)\mu^{t-1}d\mu}{(\Psi(\mu) - \lambda)(\mu - 1)}. \quad (165)$$

It is convenient to subtract the residue of  $\mu^{t-1}/(\mu - 1)$  equal to 1:

$$V(t, \lambda) - 1 = \frac{1}{2\pi i} \oint_\gamma \frac{\Psi(\mu)\mu^{t-1}d\mu}{(\Psi(\mu) - \lambda)(\mu - 1)} - \frac{1}{2\pi i} \oint_\gamma \frac{\mu^{t-1}d\mu}{\mu - 1} = \frac{\lambda}{2\pi i} \oint_\gamma \frac{\mu^{t-1}d\mu}{(\Psi(\mu) - \lambda)(\mu - 1)}. \quad (166)$$

We define now  $\gamma$  as the original contour perturbed to include an arc of radius  $1/t$  centered at 1:

$$\gamma = \gamma_1 \cup \gamma_2, \quad (167)$$

$$\gamma_1 = \{e^{i\phi}\}_{\phi_1 \leq \phi \leq 2\pi - \phi_1}, \quad (168)$$

$$\gamma_2 = \{1 + \frac{e^{i\phi}}{t}\}_{-\phi_2 \leq \phi \leq \phi_2}, \quad (169)$$

where  $\phi_1 \in (0, \frac{\pi}{2}), \phi_2 \in (\frac{\pi}{2}, \pi)$  are such that  $\gamma$  is connected. Note that  $\phi_1 \propto \frac{1}{t}$  as  $t \rightarrow \infty$ .

Now we bound separately the contribution to the integral from  $\gamma_1$  and  $\gamma_2$ . For  $\gamma_1$  and  $-\pi \leq \phi \leq \pi$  we use the inequalities

$$|\Psi(e^{i\phi}) - \lambda| \geq c|\phi|^\theta, \quad (170)$$

$$|e^{i\phi} - 1| \geq c|\phi| \quad (171)$$

with a  $\phi, \lambda$ -independent constant  $c > 0$ . This gives, using Eq. (164),

$$\lambda \left| \int_{\gamma_1} \frac{\mu^{t-1}d\mu}{(\Psi(\mu) - \lambda)(\mu - 1)} \right| \leq \lambda C \left| \int_{-\pi}^{-\phi_1} + \int_{\phi_1}^{\pi} \frac{d\phi}{|\phi|^{\theta+1}} \right| \leq C' \frac{\lambda}{\phi_1^\theta} \leq C'' \lambda t^\theta \leq C''' \epsilon^\theta. \quad (172)$$

For the  $\gamma_2$  component we use the inequalities

$$|1 + \frac{e^{i\phi}}{t}|^{t-1} \leq e, \quad (173)$$

$$|\Psi(1 + \frac{e^{i\phi}}{t}) - \lambda| \geq ct^{-\theta}, \quad -\phi_2 \leq \phi \leq \phi_2. \quad (174)$$

(Inequality (174) relies on the assumption  $\theta < 2$  and can be proved similarly to Lemma 3.) This gives

$$\lambda \left| \int_{\gamma_2} \frac{\mu^{t-1} d\mu}{(\Psi(\mu) - \lambda)(\mu - 1)} \right| \leq \lambda C \left| \int_{-\pi}^{\pi} \frac{t^{-1} d\phi}{t^{-\theta} \cdot t^{-1}} \right| \leq C' \lambda t^\theta \leq C'' \epsilon^\theta. \quad (175)$$

Fixing some  $\epsilon > 0$ , we see from Eqs. (172), (175) that under assumption (164) the expressions  $|V(t, \lambda) - 1|$ , and hence  $|V(t, \lambda)|$ , are uniformly bounded, as desired.

This completes the proof of the lemma.  $\square$

This lemma can now be used to show that replacing  $F_V(t\lambda_k^{1/\theta}, \lambda_k)$  by  $F_V(t\lambda_k^{1/\theta})$  in Eq. (137) amounts to a lower-order correction  $o(t^{-\theta\zeta})$  in the propagator  $V_t$ . The argument is similar to the respective argument for  $F_U$  in the end of Section B.1. Statement 1 of Lemma 7 is used to show this for the contribution of the terms  $k$  with  $u < t\lambda_k^{1/\theta} < v$ , for any  $0 < u < v < +\infty$ . Then, for terms with  $t\lambda_k^{1/\theta} < u$  we use the uniform boundedness of  $F_V(r, \lambda)$ , i.e. the part  $F_V(r, \lambda) \leq C$  of statement 2, and show that their contribution can be made arbitrarily small by decreasing  $u$ . Finally, for terms with  $t\lambda_k^{1/\theta} > v$  we use the part  $F_V(r, \lambda) \leq \frac{C}{r}$  of statement 2, and show that their contribution can be made arbitrarily small by increasing  $v$ .

This completes the proof of Theorem 3.

## C PROOF OF PROPOSITION 1

To simplify notation, set  $A = 1$ ; results for general  $A$ 's are easily obtained by rescaling.

Note first that for any  $\mu \in \mathbb{C} \setminus [0, 1]$  the integral in Eq. (29) converges and is nonzero. To see that it is nonzero, note that if  $\mu$  has a nonzero imaginary part, then the integral has a nonzero imaginary part of the opposite sign, hence is nonzero. On the other hand, if  $\mu > 1$  or  $\mu < 0$ , then the integral is strictly positive or negative, so also nonzero. It follows that the expression in parentheses is invertible and so  $\Psi(\mu)$  is well-defined for all  $\mu \in \mathbb{C} \setminus [0, 1]$ .

The asymptotics  $\Psi(\mu) = -\mu(1 + o(1))$  at  $\mu \rightarrow \infty$  is obvious.

To find the asymptotics at  $\mu \rightarrow 1$ , make the substitution  $z = \delta/(\mu - 1)$  in the integral:

$$\int_0^1 \frac{d\delta^{2-\theta}}{\mu - 1 + \delta} = (\mu - 1)^{\theta-1} \int_0^{1/(\mu-1)} \frac{dz^{2-\theta}}{1+z}. \quad (176)$$

As  $\mu \rightarrow 1$  the last integral converges to a standard integral:

$$\int_0^{1/(\mu-1)} \frac{dz^{2-\theta}}{1+z} \rightarrow \int_0^\infty \frac{dz^{2-\theta}}{1+z} = \frac{(2-\theta)\pi}{\sin((2-\theta)\pi)}. \quad (177)$$

The integration line in the last integral is any line connecting 0 to  $\infty$  in  $\mathbb{C} \setminus (-\infty, 0)$ ; the integral does not depend on the line thanks to the condition  $\theta > 1$ .

We prove now that  $\Psi(\{|\mu| \geq 1\}) \cap (0, 2] = \emptyset$ . Let us first show that if  $|\mu| \geq 1$  and  $\Im\mu \neq 0$ , then  $\Psi(\mu) \notin (0, +\infty)$ . To this end write

$$\Psi(\mu) = ab, \quad (178)$$

$$a = - \left( \int_0^1 \frac{(\mu - 1)d\delta^{2-\theta}}{\mu - 1 + \delta} \right)^{-1} = - \left( \int_0^1 \frac{d\delta^{2-\theta}}{1 + \frac{\delta}{\mu-1}} \right)^{-1}, \quad (179)$$

$$b = \frac{(\mu - 1)^2}{\mu} = J(\mu) - 2, \quad (180)$$

where  $J(\mu) = \mu + \frac{1}{\mu}$  is Zhukovsky's function.

Suppose, for definiteness, that  $\Im\mu > 0$ . Regarding  $a$ , note that if  $\Im\mu > 0$ , then the imaginary part of the integrand in Eq. (179) is also positive, and so  $\Im a > 0$ .

Regarding  $b$ , recall that if  $\Im\mu > 0$  and  $|\mu| > 1$ , then  $\Im J(\mu) > 0$ . On the other hand, if  $|\mu| = 1$ , then  $J(\mu) \in [-2, 2]$ . Combining these observations, we see that if  $\Im\mu > 0$  and  $|\mu| \geq 1$ , then either  $\Im b > 0$ , or  $b \leq 0$ . Since  $\Im a > 0$ , it follows that  $ab \notin (0, +\infty)$ .

We see that  $\Psi(\mu)$  can be real and positive only if  $\mu \in \mathbb{R}$ . Clearly,  $\Psi(\mu) > 0$  if  $\mu \leq -1$ , and  $\Psi(\mu) \leq 0$  if  $\mu \geq 1$ . It is easily checked by differentiation that  $\Psi(\mu)$  is monotone decreasing for  $\mu \in (-\infty, -1]$ , so the smallest positive value attained by  $\Psi$  is

$$\Psi(-1) = 2 \left( \int_0^1 \frac{d\delta^{2-\theta}}{2-\delta} \right)^{-1} > 2. \quad (181)$$

## D PROOF OF PROPOSITION 2

In terms of  $\alpha, \mathbf{b}, \mathbf{c}, D$ , the components  $P, Q$  of the characteristic polynomial  $\det(\mu - S_\lambda) = P(\mu) - \lambda Q(\mu)$  can be written as

$$P(\mu) = (\mu - 1) \det(\mu - D), \quad (182)$$

$$Q(\mu) = -\det \begin{pmatrix} \alpha & \mathbf{b}^T \\ \mathbf{c} & \mu - D \end{pmatrix} = \det(\mu - D) (\mathbf{b}^T (\mu - D)^{-1} \mathbf{c} - \alpha). \quad (183)$$

(see Theorem 1 in Yarotsky & Velikanov (2024)). Accordingly,

$$\frac{(\mu - 1)Q(\mu)}{P(\mu)} = \mathbf{b}^T (\mu - D)^{-1} \mathbf{c} - \alpha. \quad (184)$$

If  $D = \text{diag}(d_1, \dots, d_M)$ , then

$$\frac{(\mu - 1)Q(\mu)}{P(\mu)} = \sum_{m=1}^M \frac{b_m c_m}{\mu - d_m} - \alpha. \quad (185)$$

On the other hand, our definition of  $\Psi^{(M)}$  implies that

$$\frac{(\mu - 1)A}{\Psi^{(M)}(\mu)} = (\theta - 2)h\mu \sum_{m=1}^M \frac{e^{-(2-\theta)(m-1/2)h}}{\mu - 1 + e^{-(m-1/2)h}} \quad (186)$$

$$= (\theta - 2)h \left[ \sum_{m=1}^M \frac{e^{-(2-\theta)(m-1/2)h}(1 - e^{-(m-1/2)h})}{\mu - 1 + e^{-(m-1/2)h}} + \sum_{m=1}^M e^{-(2-\theta)(m-1/2)h} \right] \quad (187)$$

$$= (2 - \theta)h \left[ \sum_{m=1}^M \frac{e^{-(2-\theta)(m-1/2)h}(e^{-(m-1/2)h} - 1)}{\mu - 1 + e^{-(m-1/2)h}} - \frac{1 - e^{-(2-\theta)Mh}}{1 - e^{-(2-\theta)h}} e^{-(2-\theta)h/2} \right]. \quad (188)$$

By comparing this expansion with Eq. (185), we see that the values of  $\alpha, \mathbf{b}, \mathbf{c}, D$  given in Eqs. (32)-(35) ensure that  $P/Q = \Psi^{(M)}$ .

## E THE SYNTHETIC 1D EXAMPLE

Recall that in Section 6 we consider the synthetic 1D example in which we fit the target function  $y(x) = \mathbf{1}_{[1/4, 3/4]}(x)$  on the segment  $[0, 1]$  with a model that in the infinite-size limit has the integral form

$$\hat{y}(x) = \int_0^1 w(y)(x - y)_+ dy = \mathbf{x}^T \mathbf{w}, \quad (189)$$

where  $\mathbf{x}, \mathbf{w}$  are understood as vectors in  $L^2([0, 1])$ , and  $\mathbf{x} \equiv (x - \cdot)_+$ . We consider the loss  $L(\mathbf{w}) = \mathbb{E}_{\frac{1}{2}}(\mathbf{x}^T \mathbf{w} - y(x))^2$ , where  $\rho$  is the uniform distribution on  $[0, 1]$ .

The asymptotic power-law structure of this problem can be derived either from general theory of singular operators and target functions, or from the specific eigendecomposition available in this simple 1D setting.

**The eigenvalues.** First observe that the operator  $\mathbf{H} = \mathbb{E}_{\mathbf{x} \sim \rho}[\mathbf{x}\mathbf{x}^T]$  in our case is the integral operator

$$\mathbf{H}f(x) = \int_0^1 K(x, y)f(y)dy, \quad K(x, y) = \int_0^1 (x-z)_+(y-z)_+dz. \quad (190)$$

The operator has eigenvalues (see, e.g., Section A.6 of Yarotsky (2018))  $\lambda_k = \xi_k^{-4}$ , where

$$\xi_k = \frac{\pi}{2} + \pi k + O(e^{-\pi k}), \quad k = 0, 1, \dots \quad (191)$$

Numerically,  $\xi_0 \approx 1.875$  so the leading eigenvalue  $\lambda_0 \approx 0.0809$ .

In particular, the capacity condition (11) holds with  $\nu = 4$ .

In fact, such a power-law asymptotics is a general property of integral operators with diagonal singularities of a particular order (Birman & Solomjak, 1970). It is easily checked that the diagonal singularity of operator (190) is of order  $\alpha = 3$ . In dimension  $d$  the exponent  $\nu$  has the general form  $\nu = 1 + \frac{\alpha}{d}$ , which evaluates to 4 in our case  $d = 1$ .

**The eigencoefficients.** To establish the source condition (12), we can invoke the general theory that says that for targets that are indicator function of smooth domains we have  $\zeta = \frac{1}{d+\alpha} = \frac{1}{4}$  (Velikanov & Yarotsky, 2021). Alternatively, we can directly find  $\zeta$  thanks to the simple structure of the problem.

A short (though not quite rigorous) argument is to observe that the exact minimizer  $\mathbf{w}_*$  making the loss  $L(\mathbf{w}) = 0$  formally has the distributional form

$$\mathbf{w}_*(x) = \delta'(x - 1/4) - \delta'(x - 3/4) \quad (192)$$

with Dirac delta  $\delta(x)$ . This vector  $\mathbf{w}_*$  has an infinite  $L^2([0, 1])$  norm, in agreement with our expectation that  $\zeta = \frac{1}{4} < 1$ . The eigenfunctions of the problem can be explicitly found (Section A.6 of Yarotsky (2018)):

$$\mathbf{e}_k(x) = \cosh(\xi_k x) + \cos(\xi_k x) - \frac{\cosh(\xi_k) + \cos(\xi_k)}{\sinh(\xi_k) + \sin(\xi_k)}(\sinh(\xi_k x) + \sin(\xi_k x)). \quad (193)$$

Then, formally,

$$\mathbf{e}_k^T \mathbf{w}_* = \left. \frac{d\mathbf{e}_k(x)}{dx} \right|_{x=3/4} - \left. \frac{d\mathbf{e}_k(x)}{dx} \right|_{x=1/4} \propto \xi_k. \quad (194)$$

It follows that at small  $\lambda$ , denoting  $k_*(\lambda) = \min\{k : \lambda_k < \lambda\}$ ,

$$\sum_{k: \lambda_k < \lambda} \lambda_k (\mathbf{e}_k^T \mathbf{w}_*)^2 \propto \sum_{k \leq k_*(\lambda)} \xi_k^{-2} \propto \sum_{k \leq k_*(\lambda)} (1/2 + k)^{-2} \propto k_*^{-1}(\lambda) \propto \lambda^{-1/4}, \quad (195)$$

implying again  $\zeta = \frac{1}{4}$ .

A rigorous proof, avoiding Dirac deltas, can be given along the following lines. First note that in the setting of loss function  $L(\mathbf{w}) = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \rho} (\mathbf{x}^T \mathbf{w} - y(\mathbf{x}))^2$  the vector  $\mathbf{q}$  appearing in quadratic form (2) acquires the form  $\mathbf{q} = \mathbb{E}_{\mathbf{x} \sim \rho}[y(\mathbf{x})\mathbf{x}]$ , which in our example gives

$$\mathbf{q}(x) = \int_{1/4}^{3/4} (y-x)_+ dy. \quad (196)$$

We get from the condition  $\mathbf{H}\mathbf{w}_* = \mathbf{q}$  that

$$\mathbf{e}_k^T \mathbf{w}_* = -\frac{\mathbf{e}_k^T \mathbf{q}}{\lambda_k}. \quad (197)$$

The eigenfunctions can be written as

$$\mathbf{e}_k(x) = \cos(\xi_k x) - \sin(\xi_k x) + e^{-\xi_k x} + (-1)^k e^{-\xi_k(1-x)} + O(e^{-\xi_k}), \quad (198)$$

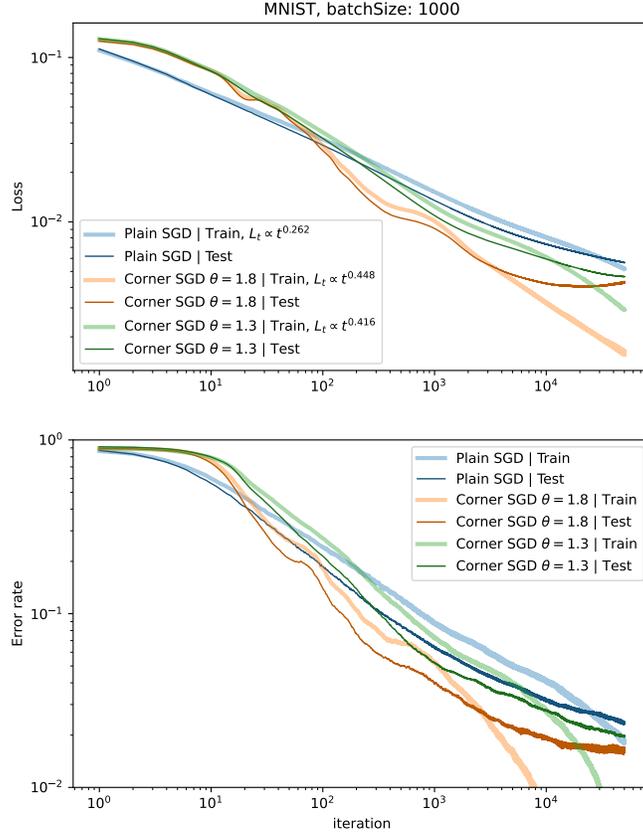


Figure 5: MNIST loss and error rate trajectories on train and test set: batch size 1000.

where the last  $O(e^{-\xi_k})$  is uniform in  $x \in [0, 1]$ . Performing integration by parts twice with vanishing boundary terms, we find that

$$\begin{aligned}
 \mathbf{e}_k^T \mathbf{q} &= \int_0^1 \left( \cos(\xi_k x) - \sin(\xi_k x) + e^{-\xi_k x} + (-1)^k e^{-\xi_k(1-x)} \right) \int_{1/4}^{3/4} (y-x)_+ dy dx + O(e^{-\xi_k}) \\
 &= -\xi_k^{-1} \int_0^1 \left( \sin(\xi_k x) + \cos(\xi_k x) - e^{-\xi_k x} + (-1)^k e^{-\xi_k(1-x)} \right) \int_{1/4}^{3/4} \mathbf{1}_{y>x} dy dx + O(e^{-\xi_k}) \\
 &= \xi_k^{-2} \int_{1/4}^{3/4} (-\cos(\xi_k x) + \sin(\xi_k x)) dx + O(e^{-\xi_k/4}) \tag{199}
 \end{aligned}$$

$$= \xi_k^{-3} \left( -\sin\left(\pi\left(\frac{1}{2} + k\right)x\right) - \cos\left(\pi\left(\frac{1}{2} + k\right)x\right) \right) \Big|_{1/4}^{3/4} + O(e^{-\xi_k/4}) \tag{200}$$

$$\propto \xi_k^{-3}, \tag{201}$$

leading to  $\mathbf{e}_k^T \mathbf{w}_* \propto \xi_k^{-3}/\lambda_k = \xi_k$ , in agreement with Eq. (194).

## F GENERALIZATION PERFORMANCE OF CORNER SGD

In Figures 5 and 6 we show both train and test trajectories of the loss and error rate (fraction of incorrectly classified images) for the MNIST classifier discussed in section 6. The two figures correspond to the batch sizes  $|B| = 1000$  and  $100$ , respectively. The test performance is computed on the standard set of 10000 images, while the training performance is computed on a subset of the full training set that also includes 10000 images.

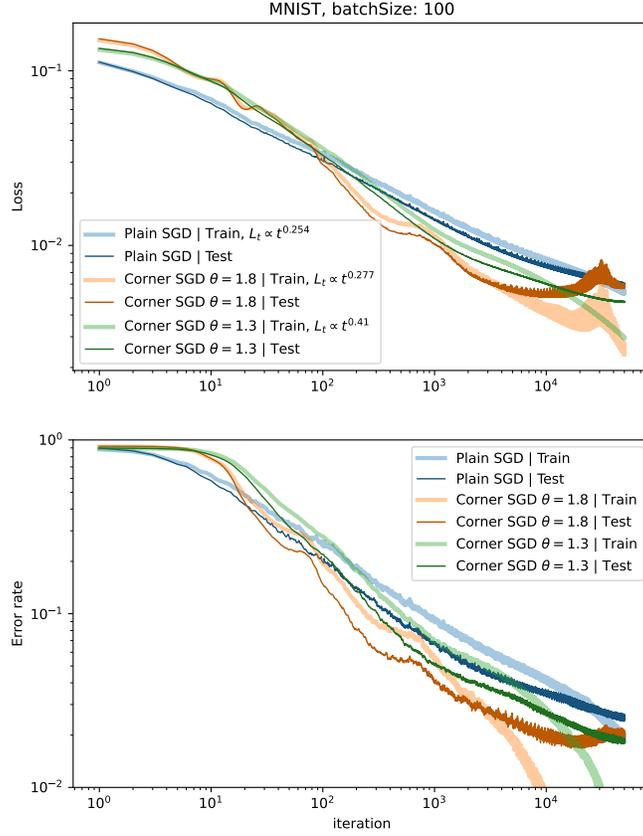


Figure 6: MNIST loss and error rate trajectories on train and test set: batch size 100.

We observe that, similarly to the training set performance, the test performance also improves faster with Corner SGD than with plain SGD. The instability of Corner SGD with  $\theta = 1.8$  and batch size 100 observed in Section 6 on the train set is also visible on the test set.

## G EXTENDING THE PROOF OF THEOREM 3 TO $\tau_2 \neq 0$

In this section we sketch (without much rigor) an argument suggesting that Theorem 3 remains valid under assumption of SE approximation with  $\tau_2 \neq 0$  at least if the batch size  $|B|$  is large enough.

Recall that the assumption  $\tau_2 = 0$  was used to write the propagators  $U_t, V_t$  in the simple form (8). These representations led to the representations (21)-(22) of  $U_t, V_t$  in terms of the contour map  $\Psi$  that were instrumental in proving Theorem 3. While we are not aware of a similar contour representation at  $\tau_2 \neq 0$ , we can expand the general  $\tau_2 \neq 0$  propagators in terms of the spectral components of the  $\tau_2 = 0$  propagators, and in this way reduce the study of the general case to the already analyzed special case.

Specifically, let us introduce the notation

$$G_0(t, \lambda) \equiv U^2(t, \lambda) = |(1 \ 0^T) S_\lambda^{t-1} \begin{pmatrix} -\alpha \\ \mathbf{e} \end{pmatrix}|^2. \quad (202)$$

Then formula (8) for the propagator  $U_t$  can be written as

$$U_t = \frac{\tau_1}{|B|} \sum_{k=1}^{\infty} \lambda_k^2 G_0(t, \lambda_k). \quad (203)$$

In the proof of Theorem 3 it was shown that (see Eqs. (53), (55))

$$G_0(t, \lambda) = U^2(t, \lambda) \approx \lambda^{2/\theta-2} F_U^2(t \lambda^{1/\theta}). \quad (204)$$

Upon substituting  $t\lambda^{1/\theta} = r$  and applying the capacity condition (11), this gave the leading term in  $U_t$ :

$$U_t \approx \frac{\tau_1}{|B|} \sum_{k=1}^{\infty} \lambda_k^{2/\theta} F_U^2(t\lambda_k^{1/\theta}) \quad (205)$$

$$= \left[ \frac{\tau_1}{|B|} \sum_{k=1}^{\infty} (t\lambda_k^{1/\theta})^2 F_U^2(t\lambda_k^{1/\theta}) \right] t^{-2} \quad (206)$$

$$\approx \left[ \frac{\tau_1}{|B|} \int_{\infty}^0 r^2 F_U^2(r) d\Lambda^{1/\nu}(t/r)^{\theta/\nu} \right] t^{-2} \quad (207)$$

$$= \left[ \frac{\tau_1}{|B|} \Lambda^{1/\nu} \int_{\infty}^0 r^2 F_U^2(r) dr^{-\theta/\nu} \right] t^{\theta/\nu-2}. \quad (208)$$

Now, if the SE approximation holds with  $\tau_2 \neq 0$ , then the propagator formulas (8) are no longer valid. Instead (see Yarotsky & Velikanov (2024)), the propagators can be written with the help of the linear transition operators  $A_\lambda$  acting on  $(M+1) \times (M+1)$  matrices  $Z$ :

$$A_\lambda Z = S_\lambda Z S_\lambda^T - \frac{\tau_2}{|B|} \lambda^2 \begin{pmatrix} -\alpha \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{1} \\ \mathbf{0} \end{pmatrix}^T Z \begin{pmatrix} \mathbf{1} \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix}^T. \quad (209)$$

In particular, Eqs. (202), (203) get replaced by

$$U_t = \frac{\tau_1}{|B|} \sum_{k=1}^{\infty} \lambda_k^2 G(t, \lambda_k), \quad (210)$$

$$G(t, \lambda) = \text{Tr} \left[ \begin{pmatrix} \mathbf{1} \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{1} \\ \mathbf{0} \end{pmatrix}^T A_\lambda^{t-1} \left[ \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix} \begin{pmatrix} -\alpha \\ \mathbf{c} \end{pmatrix}^T \right] \right]. \quad (211)$$

Note that Eq. (202) is a special case of Eq. (211) resulting at  $\tau_2 = 0$  thanks to the simple factorized structure of the transformation  $A_\lambda$  with vanishing second term.

Let us now write the binomial expansion of  $G(t, \lambda)$  by choosing one of the two terms on the r.h.s. of Eq. (209) in each of the  $t-1$  iterates of  $A_\lambda$  in Eq. (211). The key observation here is that each term in this binomial expansion can be written as a product of the  $\tau_2 = 0$  factors  $G_0$  with a suitable coefficient:

$$G(t, \lambda) = G_0(t, \lambda) + \sum_{m=1}^{t-1} \left( \frac{-\tau_2 \lambda^2}{|B|} \right)^m \times \quad (212)$$

$$\times \sum_{0 < t_1 < \dots < t_m < t} G_0(t - t_m, \lambda) G_0(t_m - t_{m-1}, \lambda) \cdots G_0(t_2 - t_1, \lambda) G_0(t_1, \lambda). \quad (213)$$

Here,  $0 < t_1 < \dots < t_m < t$  are the iterations at which the second term in Eq. (209) was chosen.

We can now apply again approximation (204) for  $G_0$  in terms of  $F_U$ , and approximate summation by integration:

$$G(t, \lambda) \approx \lambda^{2/\theta-2} \left[ F_U^2(t\lambda^{1/\theta}) + \sum_{m=1}^{\infty} \left( \frac{-\tau_2 \lambda^{1/\theta}}{|B|} \right)^m (F_U^2)^{* (m+1)}(t\lambda^{1/\theta}) \right], \quad (214)$$

where  $(F_U^2)^{* (m+1)}$  is the  $(m+1)$ -fold self-convolution of  $F_U^2$ :

$$(F_U^2)^{* (m+1)}(r) = \int \cdots \int_{0 < r_1 < \dots < r_m < r} F_U^2(r - r_m) F_U^2(r_m - r_{m-1}) \cdots F_U^2(r_1) dr_1 \cdots dr_m. \quad (215)$$

The factor  $\lambda^{1/\theta}$  in (214) results from the respective factor  $\lambda^2$  in Eq. (212), the factor  $\lambda^{2/\theta-2}$  in Eq. (204), and the integration element scaling factor  $\lambda^{-1/\theta}$  due to the substitution  $r_n = t_n \lambda^{1/\theta}$ .

The leading term in expansion (214) corresponds to the case  $\tau_2 = 0$ . Consider the next term,  $m = 1$ . The respective contribution to  $U_t$  is

$$U_t^{(1)} \equiv - \frac{\tau_1 \tau_2}{|B|^2} \sum_{k=1}^{\infty} \lambda_k^{3/\theta} (F_U^2)^{* 2}(t\lambda_k^{1/\theta}). \quad (216)$$

This expression can be analyzed similarly to the leading term in Eq. (205), giving

$$U_t^{(1)} \approx - \left[ \frac{\tau_1 \tau_2}{|B|^2} \Lambda^{1/\nu} \int_{\infty}^0 r^3 (F_U^2)^{*2}(r) dr^{-\theta/\nu} \right] t^{\theta/\nu-3}. \quad (217)$$

Note the faster decay  $t^{\theta/\nu-3}$  compared to  $t^{\theta/\nu-2}$  in the leading term. This difference results from the different exponent  $3/\theta$  on  $\lambda_k$ . It also leads to the factor  $r^3$  rather than  $r^2$  in the integral.

The coefficient in brackets in Eq. (217) is finite unless the integral diverges. To see the convergence, write

$$\int_{\infty}^0 r^3 (F_U^2)^{*2}(r) dr^{-\theta/\nu} = \frac{\theta}{\nu} \int_0^{\infty} r^{2-\theta/\nu} (F_U^2)^{*2}(r) dr \quad (218)$$

and use the inequality  $r^{2-\theta/\nu} \leq (2(r-r_1))^{2-\theta/\nu} + (2r_1)^{2-\theta/\nu}$  valid since  $2-\theta/\nu > 0$ :

$$\int_0^{\infty} r^{2-\theta/\nu} (F_U^2)^{*2}(r) dr \quad (219)$$

$$\leq \int \int_{0 < r_1 < r < \infty} [(2(r-r_1))^{2-\theta/\nu} + (2r_1)^{2-\theta/\nu}] F_U^2(r-r_1) F_U^2(r_1) dr_1 dr \quad (220)$$

$$= 2^{3-\theta/\nu} \left( \int_0^{\infty} r^{2-\theta/\nu} F_U^2(r) dr \right) \left( \int_0^{\infty} F_U^2(r) dr \right) < \infty, \quad (221)$$

since  $F_U(r) \propto r^{-\theta-1}$  as  $r \rightarrow \infty$  by Lemma 1.

Next terms in expansion (214) can be analyzed similarly, but we encounter the difficulty that, due to the associated factor  $\lambda^{m/\theta}$  in Eq. (214), they will contain the integrals  $\int_{\infty}^0 r^{2+m} (F_U^2)^{*(m+1)}(r) dr^{-\theta/\nu}$  that diverge for sufficiently large  $m$ . For this reason, it is convenient to upper bound

$$\lambda^{m/\theta} \leq \lambda_{\max}^{(m-1)/\theta} \lambda^{1/\theta}. \quad (222)$$

Then the contribution  $U_t^{(m)}$  to  $U_t$  from the term  $m$  can be upper bounded by

$$|U_t^{(m)}| \lesssim \left[ \frac{\tau_1 |\tau_2|^m \lambda_{\max}^{(m-1)/\theta}}{|B|^{m+1}} \Lambda^{1/\nu} \int_{\infty}^0 r^3 (F_U^2)^{*(m+1)}(r) dr^{-\theta/\nu} \right] t^{\theta/\nu-3}. \quad (223)$$

Using the inequality  $r^{2-\theta/\nu} \leq ((m+1)(r-r_m))^{2-\theta/\nu} + \dots + ((m+1)r_1)^{2-\theta/\nu}$ , the integral can be bounded as

$$\int_{\infty}^0 r^3 (F_U^2)^{*(m+1)}(r) dr^{-\theta/\nu} \leq \frac{\theta}{\nu} (m+1)^{3-\theta/\nu} \left( \int_0^{\infty} r^{2-\theta/\nu} F_U^2(r) dr \right) \left( \int_0^{\infty} F_U^2(r) dr \right)^m < \infty. \quad (224)$$

Summarizing, the contribution of all the terms in  $U_t$  other than the leading term  $U_t^{(0)}$  can be upper bounded by

$$|U_t - U_t^{(0)}| \lesssim C t^{\theta/\nu-3}, \quad (225)$$

with the constant

$$C = \frac{\tau_1 \theta \Lambda^{1/\nu}}{\nu} \left( \int_0^{\infty} r^{2-\theta/\nu} F_U^2(r) dr \right) \sum_{m=1}^{\infty} \frac{|\tau_2|^m \lambda_{\max}^{(m-1)/\theta}}{|B|^{m+1}} (m+1)^{3-\theta/\nu} \left( \int_0^{\infty} F_U^2(r) dr \right)^m. \quad (226)$$

If

$$|B| > |\tau_2| \lambda_{\max}^{1/\theta} \int_0^{\infty} F_U^2(r) dr, \quad (227)$$

then series (226) converges, and so  $|U_t - U_t^{(0)}| = o(U_t^{(0)})$ , as claimed.

The case of the propagators  $V_t$  can be treated similarly. Starting from  $\tau_2 = 0$ , denote

$$H_0(t, \lambda) = V^2(t, \lambda) = |(1 \ 0^T) S_{\lambda}^{t-1} \begin{pmatrix} 1 \\ 0 \end{pmatrix}|^2, \quad (228)$$

then by Eqs. (122), (124)  $H_0(t, \lambda) \approx F_V^2(t \lambda^{1/\theta})$  and

$$V_t = \sum_{k=1}^{\infty} \lambda_k (\mathbf{e}_k^T \mathbf{w}_*)^2 H_0(t, \lambda) \approx \sum_k \lambda_k (\mathbf{e}_k^T \mathbf{w}_*)^2 F_V^2(t \lambda_k^{1/\theta}). \quad (229)$$

The counterpart of  $H_0$  for general  $\tau_2$  is

$$H(t, \lambda) = \text{Tr}[(\mathbf{1}_0)(\mathbf{1}_0)^T A_\lambda^{t-1} [(\mathbf{1}_0)(\mathbf{1}_0)^T]]. \quad (230)$$

Expansion (212) gets replaced by

$$H(t, \lambda) = H_0(t, \lambda) + \sum_{m=1}^{t-1} \left( \frac{-\tau_2 \lambda^2}{|B|} \right)^m \times \quad (231)$$

$$\times \sum_{0 < t_1 < \dots < t_m < t} G_0(t - t_m, \lambda) G_0(t_m - t_{m-1}, \lambda) \cdots G_0(t_2 - t_1, \lambda) H_0(t_1, \lambda) \quad (232)$$

and expansion (214) gets replaced by

$$H(t, \lambda) \approx F_V^2(t\lambda^{1/\theta}) + \sum_{m=1}^{\infty} \left( \frac{-\tau_2 \lambda^{1/\theta}}{|B|} \right)^m ((F_U^2)^{*m} * F_V^2)(t\lambda^{1/\theta}). \quad (233)$$

The factor  $\lambda^{m/\theta}$  can again be used to extract an extra negative power of  $t$  in the asymptotic bounds. To avoid divergence of the integrals, we can use a bound

$$\lambda^{m/\theta} \leq \lambda_{\max}^{(m-\epsilon)/\theta} \lambda^{\epsilon/\theta} \quad (234)$$

with some sufficiently small  $\epsilon > 0$ . Arguing as before, we then find that for  $|B|$  large enough the contribution of all the terms  $m \geq 1$  is  $O(t^{-\theta\zeta-\epsilon})$ , i.e. asymptotically negligible compared to the leading term  $\propto t^{-\theta\zeta}$ .