
BARRON APPROXIMATION AND LOCALLY OPTIMAL WEIGHT DENSITIES FOR SHALLOW NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Mean field theories provide optimization results for shallow neural networks by analyzing the weight distribution in infinite width limits. Corresponding results for finite sized networks are obtained by particle approximations, for which sharp quantitative bounds are still an open problem. In this paper, we consider a modified mean field loss, which allows a more fine grained control over finite sized networks. We prove convexity and equidistribution properties, which directly lead to Barron type approximation results for networks sampled from local loss minimizers. We demonstrate that particle approximations of the new loss function naturally lead to gradient descent methods with dropout type regularization.

1 INTRODUCTION

Continuum limits have proven to be helpful for both the optimization and approximation theory of shallow neural networks. While there are several possible variants, in this paper we consider the limit

$$\frac{1}{m} \sum_{i=1}^m a_i \sigma(w_i) \rightarrow \int a_\pi(w) \sigma(w) d\pi(w), \quad (1)$$

for probability measure π and some function $a_\pi(w)$. The function $\sigma: \mathcal{W} \rightarrow \mathcal{H}$ maps weights from domain $\mathcal{W} \subset \mathbb{R}^{d_w}$ to a Hilbert space \mathcal{H} , with the main motivation coming from shallow neural networks given by $\mathcal{H} = L_2$ and $\sigma(w)(x) = \text{ReLU}(v^T x - b)$ for $w = (v, b)$.

The limit (1) is sometimes referred to as a two-timescale limit (Marion and Berthier, 2023; Takakura and Suzuki, 2024), if the outer a_i or $a_\pi(w)$ are chosen optimally. This differs slightly from the more common limit $\int a \sigma(w) d\pi(a, w)$ in mean field theory, where a is not a function but included in the measure's domain.

1.1 APPROXIMATION THEORY

We want to approximate a target function f as close as possible by a neural network. For its construction, we assume that f is given as a continuum limit

$$f = \int a_\pi(w) \sigma(w) d\pi(w), \quad (2)$$

and then we find the finite network by sampling $w_i \sim \pi$, effectively reversing the arrow in (1). This leads to error bounds of the form

$$\inf_{a, w} \left\| \frac{1}{m} \sum_{i=1}^m a_i \sigma(w_i) - f \right\| \leq \mathbb{E}_{w \sim \pi^m} \left\| \frac{1}{m} \sum_{i=1}^m a_\pi(w_i) \sigma(w_i) - f \right\| \lesssim m^{-1/2} \|f\|_{B(\pi)}, \quad (3)$$

for some norm $\|\cdot\|_{B(\pi)}$ defined in Section 2.2. This norm matches Barron or variation norms for a careful choice of π . This choice, however, is not available in practice and the error bound (3) works for many other distributions. Indeed, for any probability distribution π' that has a density $d\pi/d\pi'$ with respect to π we can set $a_{\pi'} := a_\pi d\pi/d\pi'$ and obtain $f = \int a_{\pi'}(w) \sigma(w) d\pi'(w)$, as well as the corresponding approximation errors. Despite the similarity in error bounds, the choice of π has

profound practical implications. For unfavourable choices, the norm $\|f\|_{B(\pi)}$ may be excessively large or even infinite.

This problem can already be seen in the simplest examples in one spacial dimension where shallow ReLU networks correspond to splines and the choice of π to a (random) adaptive placement of spline knots. Similar problems also arise in the selection of local refinement in adaptive finite element methods (DeVore, 1998; Cohen and Mirebeau, 2009). In those classical theories, the best choices are achieved by equidistributing local errors or local smoothness. We argue that this principle still applies to shallow neural networks, with equidistribution defined by

$$|a_\pi(w)| \|\sigma(w)\| = \text{constant in } w. \quad (4)$$

As we see in Section C.2, the left hand side can be understood as a local smoothness indicator, relative to the distribution π . The condition also ensures that all samples $w_i \sim \pi$, receive the same outer weight $a_\pi(w_i)$, up to a scaling factor. With equidistribution, the norm $\|f\|_{B(\pi)}$ reduces to the Barron or variation norm found in current approximation results. The main objective of this paper is to show that such theoretical choices are attained by practical optimization methods.

1.2 MEAN FIELD LIMIT

Mean field theory leverages the continuum limit (1) to understand the optimization of neural networks, starting from the objective

$$\min_{\pi} \left\| \int a_\pi(w) \sigma(w) d\pi(w) - f \right\|^2, \quad (5)$$

or more commonly with the alternative continuum limit $\int a\sigma(w) d\pi(a, w)$. This problem is convex in π and optimized by Wasserstein gradient flow or Langevin dynamics. Implications for finite networks are achieved by replacing π with a particle approximation $\pi \approx \frac{1}{m} \sum_{i=1}^m \delta_{w_i}$ for which Wasserstein gradient flow matches the gradient flow of the finite networks.

While the methods converge to global optima, it is not clear that all global optima are good ones: They achieve zero loss and thus match the exact representation (2). But as we have seen, there are many such representations with severely different performance of finite approximations. The very recent paper Takakura and Suzuki (2024), which we became aware of after completion of this paper, controls this choice by an extra penalty term, similar to the $\|f\|_{B(\pi)}$ norm above. We choose an alternative route without penalties but a different loss function.

1.3 NEW CONTRIBUTIONS

Instead of the standard mean field loss (5), we optimize the expected sampling error in the approximation bound (3) directly

$$\min_{\pi} \mathbb{E}_{\mathbf{w} \sim \pi^m} \left\| \frac{1}{m} \sum_{i=1}^m a_\pi(w_i) \sigma(w_i) - f \right\|^2. \quad (6)$$

Unlike mean field theory, the expectation is outside of the norm, which ensures that the loss is non-zero and provides a direct control over finite sized networks. Due to the outer product measure $\pi^m = \pi \otimes \dots \otimes \pi$ and the π -dependent function a_π , this optimization problem may appear overly difficult. Nonetheless, the main results of this paper show that it has several favourable properties:

1. Theorem 3.1: The optimization problem (6) is convex over admissible probability distributions π .
2. Lemmas 3.2, 3.4: Local minimizers π satisfy the *equidistribution* property (4)

$$|a_\pi(w)| \|\sigma(w)\| = \text{constant} + \text{perturbations}, \quad \text{for all } w, \pi - \text{a.s.}$$

3. Theorem 3.3, 3.5: Local minimizers π achieve the approximation errors

$$\mathbb{E}_{\mathbf{w} \sim \pi^m} \left\| \frac{1}{m} \sum_{i=1}^m a_\pi(w_i) \sigma(w_i) - f \right\|^2 \leq \frac{4}{m} \left[|f|_{B(\delta, \epsilon)}^2 + \pi(\delta) \right] + \epsilon^2.$$

for all $m \in \mathbb{N}$, where $|\cdot|_{B(\delta, \epsilon)}$ is a stable variant of classical Barron norms and ϵ, δ perturbation terms dependent on $a_\pi(w)$.

-
- 108 4. Section 4: As for the standard mean field loss, Wasserstein gradient flow matches gradient
109 flow of finite networks. For corresponding gradient descent methods, the outer expectation
110 naturally leads to a dropout regularization.
111

112 We consider two variants for these results, depending on the definition of $a_\pi(w)$. The first is an
113 idealized scenario, with a streamlined analysis, where a_π is given as the Radon-Nikodym derivative
114 of an ideal distribution. The second is a more practical choice, where a_π is defined as the optimal
115 coefficients of finite networks, averaged over off-diagonal weights.
116

117 1.4 LITERATURE REVIEW 118

119 **Approximation** The approximation theory of neural networks seeks to understand how well a
120 target function f can be approximated by a neural network. Universal approximation theorems
121 (Cybenko, 1989; Hornik et al., 1989; Barron, 1993; Zhou, 2020; Lu et al., 2017; Hanin and Sellke,
122 2017) show that this is possible to arbitrary accuracy.

123 Since these early results, lots of effort has been spent on a more fine grained picture that seeks
124 quantitative error bounds given various different smoothness properties of the target function and ar-
125 chitectures of the networks. For Sobolev or Besov smooth targets, classical approximation methods
126 like finite elements or splines are effective and neural networks match their performance (Gribon-
127 val et al., 2022; Gühring et al., 2020; Opschoor et al., 2020; Li et al., 2019; Suzuki, 2019). If one
128 allows discontinuous weight assignments, neural networks can achieve higher approximation rates
129 that classical methods (Yarotsky, 2017; 2018; Yarotsky and Zhevnerchuk, 2020; Daubechies et al.,
130 2022; Shen et al., 2019; Lu et al., 2021). Sobolev and Besov smoothness is not suitable for high
131 dimensional targets and with more tailored smoothness assumptions in Barron and variation norms,
132 neural networks achieve dimension independent approximation rates (Bach, 2017; Klusowski and
133 Barron, 2018; Weinan et al., 2022; Li et al., 2020; Siegel and Xu, 2020; 2022a; Bresler and Nagaraj,
134 2020; Parhi and Nowak, 2021; Unser, 2023; Siegel and Xu, 2024). More information is contained
135 in the reviews (Pinkus, 1999; DeVore et al., 2021; Weinan et al., 2020; Berner et al., 2022).

136 Approximation theorems state the existence of good neural network approximation, but usually do
137 not include a corresponding training mechanism. Early results for gradient descent trained networks
138 (Jentzen and Riekert, 2022; Ibragimov et al., 2022; Drews and Kohler, 2022; Kohler and Krzyzak,
139 2022; Gentile and Welper, 2024; Welper, 2024a; 2025) rely on the convex outer layer, or lazy lin-
140 earized training regimes and therefore cannot match the full approximation power known in theory.
141 Barron smoothness approximation results can be achieved by greedy algorithms (Siegel and Xu,
142 2022b; Siegel et al., 2023) which reduce training to another non-convex sub-problem.

143 **Landscape Analysis** Properties of the training objective are studied in landscape analysis. With
144 strong assumptions or over-parametrization, local minima are global minima (Soudry and Carmon,
145 2016; Kawaguchi, 2016; Nguyen and Hein, 2017; Ge et al., 2018; Du and Lee, 2018; Soltanolkotabi
146 et al., 2019; Venturi et al., 2019; Kawaguchi et al., 2019; Kawaguchi and Huang, 2019). In dif-
147 ferent regimes this is not true (Swirszcz et al., 2017; Safran and Shamir, 2018; He et al., 2020;
148 Ding et al., 2022; Jentzen and Riekert, 2024), or local minima are path connected (He et al., 2020).
149 Approximation properties of critical points, are studied in Welper (2024b). The reference achieves
150 sharp optimization bounds for one dimensional problems by optimizing the network weights. In
151 comparison, we optimize the weight distribution and obtain results in high dimensions.

152 **Lazy Training Regime** For very wide networks, the gradient descent training dynamics are dom-
153 inated by linearization at the initial value. This leads to gradient descent convergence based on the
154 neural tangent kernel, originally introduced in (Jacot et al., 2018; Li and Liang, 2018; Allen-Zhu
155 et al., 2019; Du et al., 2019b;a), and then further developed in e.g. (Zou et al., 2020; Arora et al.,
156 2019a;b; Su and Yang, 2019; Lee et al., 2019; Song and Yang, 2019; Zou and Gu, 2019; Kawaguchi
157 and Huang, 2019; Chizat et al., 2019; Oymak and Soltanolkotabi, 2020; Ji and Telgarsky, 2020;
158 Nguyen and Mondelli, 2020; Bai and Lee, 2020; Cao and Gu, 2020; Chen et al., 2021; Song et al.,
159 2021; Lee et al., 2022; Gentile and Welper, 2024; Welper, 2024a; 2025; Welper and Keene, 2025).
160 Due to the inherent linearization, these approaches do not match the network’s full potential, con-
161 firmed by empirical studies (Vyas et al., 2023) and in (Lee et al., 2020; Seleznova and Kutyniok,
2022).

Beyond the Lazy Regime More recent papers (Damian et al., 2022; Lee et al., 2024; Mousavi-Hosseini et al., 2023) show that neural networks can achieve superior results to linearization or kernel methods. In particular, they can train high dimensional functions of type $g(Ux)$ with some intrinsic lower dimensionality given by a wide matrix $U \in \mathbb{R}^{r \times d}$ with $r \ll d$.

Mean Field Limits Earlier results (Chizat and Bach, 2018; Mei et al., 2018; Rotskoff and Vandeen-Eijnden, 2018; Hu et al., 2020; Sirignano and Spiliopoulos, 2020; Bach and Chizat, 2021) consider the limiting objective (5) and particle discretizations. They show convergence to global minima for gradient flow time $t \rightarrow \infty$ and weight $m \rightarrow \infty$. Newer results contain finer grained quantitative error bounds for gradient flow time (Chizat, 2022; Nitanda et al., 2022) and discretization (Chen et al., 2023; Suzuki et al., 2023; Takakura and Suzuki, 2024).

1.5 NOTATIONS

The notations used throughout the paper are summarized in Appendix A.

2 CONTINUUM LIMIT AND APPROXIMATION

2.1 CONTINUUM LIMIT

Throughout this paper, we assume that the target function f is given by

$$f = \int \sigma(w) d\phi(w), \quad (7)$$

where ϕ is contained in the space \mathcal{M} of finite signed measures on \mathcal{W} . In case ϕ is not unique, we choose a minimizer of the Barron norm in (12), below.

Note that the continuum limit does not include outer weights a . While in mean field theory outer and inner weights $(a, w) \rightarrow v$ are often joined into one single variable, we marginalize instead. E.g. if f is given by a probability density $\rho(a, w)$, we obtain

$$f = \iint \sigma(w) a \rho(a, w) da dw = \int \sigma(w) \left(\int a \rho(a, w) da \right) dw = \int \sigma(w) d\phi(w), \quad (8)$$

with signed measure $\phi(A) = \int_A \left(\int a \rho(a, w) da \right) dw$. Appendix C.2 shows that ϕ is often unique, unlike ρ , which can accommodate different distributions of w by adapting the distribution of a .

2.2 APPROXIMATION: MAUREY SAMPLING

The main interest of continuum limits is to derive properties of their discrete practical counterparts. To study approximation properties, these are constructed by sampling.

In detail, we construct finite networks by replacing the integral representation (7) with a sample mean. However, we cannot sample from ϕ directly because it is signed and not normalized, hence not a probability. Therefore, let ϕ be absolutely continuous with respect to some probability measure π from which we sample instead. Then, we can rewrite the continuum limit with the Radon-Nikodym derivative $a_\pi := d\phi/d\pi$ so that

$$f = \int \sigma(w) \frac{d\phi}{d\pi} d\pi(w) = \mathbb{E}_{w \sim \pi} [a_\pi(w) \sigma(w)] \quad (9)$$

and replace the expectation by mean

$$f \approx \frac{1}{m} \sum_{i=1}^m a_\pi(w_i) \sigma(w_i), \quad (10)$$

with π -i.i.d samples w_i . To state error bounds, we define the *Barron norm relative to π* :

$$\|f\|_{B(\pi)} := \|f\| + |f|_{B(\pi)}, \quad |f|_{B(\pi)}^2 := \mathbb{E}_{w \sim \pi} \left[\|a_\pi(w) \sigma(w)\|^2 \right]. \quad (11)$$

Then, along standard lines, (Weinan et al., 2022; Siegel and Xu, 2024) reviewed in Appendix C.1, we obtain:

Theorem 2.1. *Let f and ϕ be given by (7). Let ϕ be absolutely continuous with respect to some probability measure π , let $a_\pi = d\phi/d\pi$ be the Radon-Nikodym derivative and let $w_i, i \in [m]$ be i.i.d sampled from π . Then*

$$\mathbb{E} \left\| \frac{1}{m} \sum_{i=1}^m a_\pi(w_i) \sigma(w_i) - f \right\|^2 \leq \frac{2}{m} |f|_{B(\pi)}^2.$$

We now turn to the best possible choice of π . Intuitively, we want to place the samples w_i so that all summands in the sample mean (10) have the same weight. Formally, this means that the size of the summands $|a_\pi(w)| \|\sigma(w)\|$ is independent of w , which we refer to as *equidistribution*. The optimality of this choice is confirmed in Lemma 3.2 and leads to standard Barron type approximation error bounds (Weinan et al., 2022; Siegel and Xu, 2024), shown in Appendix C.1:

Corollary 2.2. *Let f and ϕ be given by (7). For measurable sets $W \subset \mathcal{W}$, define the probability measure*

$$\pi(W) = |f|_B^{-1} \int_W \|\sigma(w)\| d|\phi|(w), \quad |f|_B := \int_{\mathcal{W}} \|\sigma(w)\| d|\phi|(w).$$

Let $a_\pi = d\phi/d\pi$ be the Radon-Nikodym derivative and $w_i, i \in [m]$ be i.i.d sampled from π . Then $|a_\pi(w)| \|\sigma(w)\| = |f|_{B(\pi)}$ π -almost surely and

$$\mathbb{E} \left\| \frac{1}{m} \sum_{i=1}^m a_\pi(w_i) \sigma(w_i) - f \right\|^2 \leq \frac{2}{m} |f|_{B(\pi)}^2 = \frac{2}{m} |f|_B^2.$$

Based on the last corollary, we define the *Barron norm*

$$\|f\|_B = \|f\| + |f|_B, \quad |f|_B = \int \|\sigma(w)\| d|\phi|(w). \quad (12)$$

One usually takes the infimum over all eligible ϕ , which we have already done in the definition (7). Our definition is similar to the Barron norms in Ongie et al. (2020), but there are several alternatives in the literature. E.g. the Barron norms in Weinan et al. (2022) use the continuum representation with outer weights as in (8) and define $|f|_B = \inf_\rho \mathbb{E}_{w \sim \pi} |a_\pi(w)| |w|$, which for ReLU activations is slightly larger than ours. An overview over alternative definitions of Barron norms is given in Siegel and Xu (2021).

The remarkable property of Theorem 2.1 and Corollary 2.2 is the dimension d independent error rate $m^{-1/2}$, which can be improved slightly to $m^{-\frac{1}{2} - \frac{2}{2d}}$ for ReLU activations (Siegel and Xu, 2024). However, these results are of theoretical nature because practically we do not know ϕ, π or a_π . Also note that an optimal choice of the sampling distribution π is paramount because the norm $\|f\|_{B(\pi)}$ can be significantly larger than the optimal $\|f\|_B$, even infinite, leading to poor performance of the discrete sample mean.

2.3 EXAMPLES

Appendix C.2 contains a simple model problem, where π can be computed directly and equidistribution is comparable to classical approximation methods like splines or finite elements. The choice of sample distribution is also crucial in Monte Carlo integration $\int f(x) d\pi(x)$, where sampling from π directly is preferred over sampling from generic distributions weighted by π 's density.

3 MAIN RESULTS

As we have seen in the last section, low approximation errors require a careful choice of the sample distribution π . Since all distributions, subject only to absolute continuity, achieve zero mean field loss (5), we use the more fine grained loss

$$\min_{\pi \in \mathcal{M}_{+,1}} \ell(\pi) := \min_{\pi \in \mathcal{M}_{+,1}} \mathbb{E}_{w \sim \pi^m} \left\| \frac{1}{m} \sum_{i=1}^m a_\pi(w_i) \sigma(w_i) - f \right\|^2, \quad (13)$$

which offers direct control over finite sized networks. Despite the product measure and nonlinear map $\pi \rightarrow a_\pi$, we show that local minimizers achieve equidistribution and Barron type approximation errors, first for idealized a_π in Section 3.1 and then for more practical choices in Section 3.2.

3.1 EXACT OUTER WEIGHTS

In this section, for arbitrary probability measure $\phi \ll \pi \in \mathcal{M}_{+,1}$, we define the outer weights by the Radon-Nikodym derivative

$$a_\pi := \frac{d\phi}{d\pi} \quad f = \int \sigma(w) d\phi(w) = \int \sigma(w) \frac{d\phi}{d\pi} d\pi(w) = \mathbb{E}_{w \sim \pi} [a_\pi(w)\sigma(w)], \quad (14)$$

where again we choose ϕ with minimal Barron norm $|f|_B$ in case it is not unique. Despite the product measure $w \sim \pi^m$ and the dependence of a_π on π , the loss $\ell(\pi)$ is convex on the set of permissible probability distributions, proven in Appendix E.1.

Theorem 3.1. *Let ϕ and a_π be given by (14). Then the loss (13) is convex on the set of probability measures $\pi \in \mathcal{M}_{+,1}$ with $\phi \ll \pi$ and $a_\pi \in L^2(\pi)$.*

It follows that all local minima are global. The first order optimality criteria yield equidistribution, independent of m , shown in Appendix E.2.

Lemma 3.2. *Let π be an absolutely continuous $\phi \ll \pi$ local minimizer of (13), with $a_\pi \in L^2(\pi)$ defined in (14) and $\|\sigma(\cdot)\| \in L^\infty(\pi)$. Then*

$$|a_\pi(w)| \|\sigma(w)\| = \lambda, \quad \pi\text{-almost surely,}$$

for some $\lambda \in \mathbb{R}$.

Corollary 2.2 suggests that $\lambda = |f|_B$, which justifies our definition of the Barron norm (12). We directly obtain the following approximation properties for finite networks sampled from optimizers π , proven in Appendix E.3.

Theorem 3.3. *Let $a_\pi \in L^2(\pi)$ given by (14), and π be a local minimizer of the loss (13). Then*

$$\mathbb{E}_{w \sim \pi^m} \left\| \frac{1}{m} \sum_{i=1}^m a_\pi(w_i)\sigma(w_i) - f \right\|^2 \leq \frac{2}{m} |f|_B^2.$$

3.2 AVERAGED OUTER WEIGHTS

The results in the last section are based in the idealized outer weights $a_\pi = d\phi/d\pi$. They provide a streamlined theory, but are not available in practice. Therefore, in this section we consider computable alternatives.

Construction of Outer Weights: Best Approximation To construct practical approximations of a_π , we first consider finite sized neural networks: For width M , outer weights $\mathbf{a} \in \mathbb{R}^M$ and inner weights $w \in \mathcal{W}^M$, define the shallow neural network

$$f_{\mathbf{a},w} := \frac{1}{M} \sum_{i=1}^M a_i \sigma(w_i),$$

with optimal outer weights

$$\mathbf{a}(w) := \arg \min_{\mathbf{a} \in \mathbb{R}^M}^+ \|f_{\mathbf{a},w} - f\|^2, \quad f_w := f_{\mathbf{a}(w),w}, \quad (15)$$

where $\arg \min^+$ picks the candidate with minimal Euclidean norm, in case the minimizers are not unique. From these, we construct $\bar{a}_\pi(w)$, depending only on one single $w = w_1$ and distribution π , by averaging over all other weights

$$\bar{a}_\pi(w) := \mathbb{E}_{w_2, \dots, w_M \sim \pi} [\mathbf{a}_1(w, w_2, \dots, w_M)].$$

We only pick the first component \mathbf{a}_1 , which is multiplied with $\sigma(w) = \sigma(w_1)$ in the network. Unlike a_π in the last section, the new \bar{a}_π can be practically approximated by replacing the expectation by a sample mean.

The idealized coefficients a_π are chosen so that the limiting network $\mathbb{E}_{w \sim \pi} [a_\pi(w)\sigma(w)] = f$ matches the target function f exactly. Lemma F.5 in the appendix shows that the new coefficients \bar{a}_π satisfy this identity approximately

$$\mathbb{E}_{w \sim \pi} [\bar{a}_\pi(w)\sigma(w)] = \mathbb{E}_{w \sim \pi} [f_w],$$

where f_w is the best approximation of f given inner weights w_i . We can ensure that $\mathbb{E}_{w \sim \pi} [f_w]$ is close to f by choosing sufficiently large width $M \gg m$ in the definition of \bar{a}_π , compared to width m in the network loss (13).

Construction of Outer Weights: Abstraction The main results in this section are still true if we replace the least squares minimizer by an arbitrary function $\mathbf{a}: \mathcal{W}^M \rightarrow \mathbb{R}^M$, subject only to the symmetry condition

$$P\mathbf{a}(w) = \mathbf{a}(Pw), \quad \text{for all permutation matrices } P \in \mathbb{R}^{M \times M}. \quad (16)$$

This symmetry is satisfied by our original least squares coefficients (15) as shown in Lemma F.1. As before, we define

$$\bar{a}_\pi(w) := \mathbb{E}_{w_2, \dots, w_M \sim \pi} [\mathbf{a}_1(w, w_2, \dots, w_M)], \quad f_w := f_{\mathbf{a}(w), w} \quad (17)$$

and optimize the loss

$$\min_{\substack{\pi \in \mathcal{M}_+ \\ \pi(1)=1}} \ell(\pi) := \min_{\substack{\pi \in \mathcal{M}_+ \\ \pi(1)=1}} \mathbb{E}_{w \sim \pi^m} \left\| \frac{1}{m} \sum_{i=1}^m \bar{a}_\pi(w_i)\sigma(w_i) - f \right\|^2. \quad (18)$$

Perturbation Terms The approximation \bar{a}_π of a_π results in extra perturbation terms Δ_1 and Δ_2 , which we define next. With $g_\pi(w) = a_\pi(w)\sigma(w) - f$ and notations from Appendix A.2, set

$$\nu(\Delta_1) = 2 \sum_{j=1}^M \iint \langle g_\pi(w_j), \mathbb{E}_{w_{\setminus j}} [\mathbf{a}_j(w)] \sigma(w_j) \rangle d\nu(w_1) d\pi(w_j), \quad (19)$$

$$\nu(\Delta_1) = 2M \langle G_\pi, \nu(\mathbb{E}_{w_{\setminus 1}} [f_w - f_{w_{\setminus 1}}]) \rangle, \quad (20)$$

$$\nu(\Delta_2) = 2(m-1)M \langle \mathbb{E}_w [f_w - f], \nu(\mathbb{E}_{w_{\setminus 1}} [f_w - f_{w_{\setminus 1}}]) \rangle, \quad (21)$$

for all $\nu \in \mathcal{M}$. For Δ_1 , we provide two variants that are equal up to an irrelevant constant by Lemma F.10 in the appendix: The first is elementary but more complicated. Note that for $j = 1$, the w_1 in the expectation is bound by the ν integral and therefore the outer $d\pi(w_1)$ integrates a constant and can be removed. The second variant is a simplification given existence of some $G_\pi \in \mathcal{H}$, independent of w such that

$$\langle \bar{a}_\pi(w)\sigma(w) - f, \sigma(w) \rangle = \langle G_\pi, \sigma(w) \rangle, \quad \pi - \text{a.s. in } w \in \mathcal{W}. \quad (22)$$

The left hand side is a function of w and for common choices of σ , by means of choosing G_π , the right hand side can represent arbitrary functions of w , subject to some regularity like Barron smoothness. E.g. for shallow neural networks without bias on the sphere $\mathcal{H} = L^2(\mathbb{S}^{d-1})$, the expression $\langle G_\pi, \sigma(w) \rangle = \int G_\pi(x) \text{ReLU}(w^T x) dx$ is a function of w in the same way the symmetric variant $f(x) = \int a(w) \text{ReLU}(w^T x) dw$ is a function of x .

We postpone a more careful discussion of the perturbation terms to Appendix B. In short, we expect $|\Delta_1| \lesssim M^{-1/2}$ and $|\Delta_2| \lesssim mM^{-1/2}$ so that they vanish for $m \ll M$. Furthermore, if we choose $m = 1$, the second term is zero.

Equidistribution Our first main result shows equidistribution analogous to Lemma 3.2, up to the two perturbation terms Δ_1 and Δ_2 . The proof is in Appendix F.4.

Lemma 3.4. *Let π be a local minimum of (18) with $\mathbf{a}_1 \in L^2(\pi^m)$, defined in (16), and perturbation terms Δ_1 defined in (19) or (20) and Δ_2 defined in (21). Then*

$$\|\bar{a}_\pi(w)\| \|\sigma(w)\| = \lambda + \Delta_1(w) + \Delta_2(w), \quad \pi\text{-almost surely}$$

for some constant $\lambda \in \mathbb{R}$.

Approximation In Section 2.2, we have seen how equidistribution leads to approximation results in Barron norms. We introduce a stable variant, to account for perturbation errors. For measure $\pi \in \mathcal{M}$, function $\delta: \mathcal{W} \rightarrow \mathbb{R}$, $\epsilon \geq 0$ and $\lambda \in \mathbb{R}$, consider the equidistribution bound

$$\left| \lambda - |a(w)|^2 \|\sigma(w)\|^2 \right| \leq \delta(w) \quad \pi - \text{a.s.} \quad (23)$$

and the approximation bound

$$\|\mathbb{E}_{w \sim \pi} [a(w)\sigma(w)] - f\| \leq \epsilon. \quad (24)$$

Then we define the *stable Barron norm* by

$$|f|_{B(\delta, \epsilon)}^2 := \sup \{ \lambda \in \mathbb{R} \mid \exists \pi \in \mathcal{M}_{+,1}, \exists a \in L^2(\pi) \text{ so that (23) and (24) are satisfied} \}.$$

To compare this with the Barron norm, we choose π as in Corollary 2.2 and $a_\pi = d\phi/d\pi$. Then π is a probability measure and the corollary shows that (23) and (24) are satisfied with perturbations $\delta = \epsilon = 0$ and $\lambda := |f|_B^2$. If the continuum limit is unique so that we can skip the minimum after (7), it follows that

$$|f|_{B(\delta=0, \epsilon=0)} = |f|_B.$$

The second main result shows approximation for stable Barron smoothness, analogous to Theorem 3.3. The proof is in Appendix F.5.

Theorem 3.5. *Let π be a local minimum of (18) with $\mathbf{a}_1 \in L^2(\pi^m)$, defined in (16), and perturbation terms Δ_1 defined in (19) or (20) and Δ_2 defined in (21). Define*

$$\delta(w) := \bar{\lambda} + \Delta_1(w) + \Delta_2(w), \quad \epsilon := \|\mathbb{E}_{w \sim \pi} [f_w - f]\|$$

π -almost surely for arbitrary $\bar{\lambda} \in \mathbb{R}$. Then for all $m \in \mathbb{N}$

$$\mathbb{E}_{w \sim \pi^m} \left\| \frac{1}{m} \sum_{i=1}^m \bar{a}_\pi(w_i) \sigma(w_i) - f \right\|^2 \leq \frac{4}{m} \left[|f|_{B(|\delta|, \epsilon)}^2 + \pi(|\delta|) \right] + 2\epsilon^2.$$

Both Lemmas 3.2 and 3.4 show equidistribution π -almost surely. Hence, it may fail on sets of zero π -measure but non-zero Lebesgue measure. For the choice a_π , the absolute continuity $\phi \ll \pi$ ensures that such sets are irrelevant for the exact representation or approximation. However, the results for \bar{a}_π have no such condition and π may be zero on relevant subsets. However, such sets cannot be too important as long as the approximation error ϵ is small.

4 TOWARDS PRACTICAL ALGORITHMS

The main theorems provide equidistribution and approximation properties of locally optimal distributions π of inner neural network weights. Practically, we need to sample or approximate these distributions. This section provides a preliminary and informal discussion along the lines of mean field theory.

Wasserstein Gradient Flow A natural candidate to optimize the losses (13) or (18) is Wasserstein gradient flow (WGF), given by

$$\dot{\pi} = \text{div} (\pi \nabla_w \nabla_\pi \ell(\pi)),$$

where $\nabla_\pi \ell(\pi)(w)$ denotes the Riesz representation $D[\ell(\pi)]\nu = \int \nabla_\pi \ell(\pi)(w) d\nu(w)$ of the directional derivatives. The WGF is implicitly understood in a distributional sense

$$\iint \{ \dot{\varphi}(t, w) - \nabla_w \varphi(t, w) \cdot \nabla_w \nabla_\pi \ell(\pi)(w) \} d\pi(w) dt = 0, \quad (25)$$

for all smooth and compactly supported functions $\varphi: [0, \infty) \times \mathcal{W} \rightarrow \mathbb{R}$.

To obtain a finite representation of π , we use a particle approximation, determined by weights $w \in \mathbb{R}^m$ for large $m \gg m, M$, trained by regular gradient flow:

$$\pi_w := \frac{1}{m} \sum_{i=1}^m \delta_{w_i}, \quad \dot{w} := -m \nabla_w \ell(\pi_w). \quad (26)$$

On the one hand, discretizing in time and expectation, leads to gradient descent methods with dropout type regularization, as we see in the next paragraph. On the other hand, gradient flow training of \mathbf{w} matches WGF training of the corresponding distribution $\pi_{\mathbf{w}}$. The proof is well known (Chizat and Bach, 2018) and included for completeness in Appendix G.

Lemma 4.1. *Let the discrete measures $\pi_{\mathbf{w}(t)}$ be defined by gradient flow (26). Then they satisfy the Wasserstein gradient flow (25).*

The following result characterizes the stationary points of WGF, see Appendix G.

Lemma 4.2. *Let $\mathbf{w} \in \mathcal{W}^m$ be a stationary point of gradient flow, together with $\pi = \pi_{\mathbf{w}}$ defined in (26). Let \tilde{a}_π be defined by (17), with $\mathbf{a}_1 \in L^2(\pi)$, Δ_1 be defined by (19) or (20) and Δ_2 by (21). Then for all $i \in [m]$ we have*

$$\nabla_{w_i} \left[-m \|\tilde{a}_\pi(w_i)\sigma(w_i)\|^2 + \Delta_1(w_i) + \Delta_2(w_i) \right] = 0.$$

If w was a continuous variable, the zero gradient would imply equidistribution on connected components of π 's support. Confined to discrete points, analogous conclusions require extra regularity, which is left for future work.

Gradient Descent and Dropout To obtain a practical algorithm, we replace expectations with sample means

$$\bar{a}_\pi \approx \tilde{a}_\pi(v) = \frac{1}{N_a} \sum_{j=1}^{N_a} \mathbf{a}_1(v, \mathbf{v}_1^j), \quad \mathbf{v}^j \sim \pi^M, \quad (27)$$

$$\ell(\pi) \approx \tilde{\ell}(\pi) = \frac{1}{N_\ell} \sum_{j=1}^{N_\ell} \left\| \frac{1}{m} \sum_{i=1}^m \tilde{a}_\pi(v_i^j)\sigma(v_i^j) - f \right\|^2, \quad \mathbf{v}^j \sim \pi^m, \quad (28)$$

for some numbers N_a and N_ℓ , and we replace gradient flow by gradient descent with learning rate γ

$$\mathbf{w}^{n+1} = \mathbf{w}^n - \gamma \nabla_{\mathbf{w}} \tilde{\ell}(\pi_{\mathbf{w}^n})$$

Since the samples \mathbf{v}^j are a sub-selection of of the discrete \mathbf{w} , the gradient is well defined. Unravelling the discrete measure $\pi_{\mathbf{w}^n}$, this is a standard gradient descent iteration with two differences: First, instead of optimizing both a_i and w_i simultaneously, we ensure that the outer coefficients are always in an optimized state. Second, we apply two extra sample means. As we see in the following, these are comparable to dropout (Nitish et al., 2014) regularization in standard neural network training.

The connection is easiest for the outer expectation of $\ell(\pi)$. If we add all summands from the mean (28) as separate gradient descent steps and always use the newest possible sample distribution, we obtain the iteration

$$\mathbf{w}^{n+1} = \mathbf{w}^n - \frac{\gamma}{N_\ell} \nabla_{\mathbf{w}} \left\| \frac{1}{m} \sum_{i=1}^m \tilde{a}_\pi(v_i)\sigma(v_i) - f \right\|^2, \quad \mathbf{v} \sim \pi_{\mathbf{w}^n}^m, \quad (29)$$

see Appendix G.3 for details. Since $\pi_{\mathbf{w}}$ is a sum of Dirac deltas, the possible values of v_i are given by the set $\{w_1, \dots, w_m\}$. Hence, the sum over m random elements of this set is identical to dropout of the full sum $\frac{1}{m} \sum_{i=1}^m \tilde{a}_{\pi_{\mathbf{w}}}(w_i)\sigma(w_i)$.

Similarly, intertwining gradient descent optimization and mean for the coefficients \tilde{a}_π yields the iteration

$$\tilde{a}_\pi^n(v) = \tilde{a}_\pi^{n-1}(v) - \frac{\lambda}{N_a} \nabla_a \left\| \frac{1}{M} \sum_{i=1}^M \tilde{a}_\pi^{n-1}(v_i)\sigma(v_i) - f \right\|^2, \quad v_1 = v, \quad \mathbf{v}_{\setminus 1} \sim \pi^{M-1}, \quad (30)$$

for all particles $v \in \{w_1, \dots, w_m\}$, see Appendix G.4 for details. Again the inner sum has random terms dropped out of the full sum over all $v_i \in \{w_1, \dots, w_m\}$.

5 NUMERICAL EXPERIMENTS

In the following preliminary numerical experiments we train the function $f(x, y) = x^2$ on the domain $[-1, 1]^2$. While the domain is two dimensional, the function only depends on one variable to see if gradient descent can find this intrinsic low dimensionality. We train a shallow ReLU network on $n = 100$ uniformly random samples, with 500000 gradient descent steps of learning rate 0.01, widths 16, 32, 64, 128, 256 and dropout 0, 0.5, 0.7. Instead of computing \bar{a}_π , we include the outer weights in the regular gradient descent training.

Convergence rates are shown in Figure 1. Up to some outliers, they stay below the predicted $m^{-1/2}$. As rates in low dimensions can be higher (Siegel and Xu, 2024), this is not unexpected. Without dropout, the overall error is lower, because the outer coefficients are not disrupted by removed neurons.

Figure 2 shows $\|a_i \sigma(w_i)\|$ for all network indices i , sorted by magnitude. This quantity is equidistributed in the continuum limit in Lemmas 3.2 and 3.4 on the support of π . It seems that gradient descent either aligns the inner weights with the support or deactivates neurons by setting the outer coefficients to zero, so that we see two distinct plateaus in the figure. Note also that in Section 4 our loss leads to dropout regularization, without which the equidistribution disappears in the experiments.

The bottom row of Figure 1 shows the inner weights v_i of $\sigma(w_i) = \text{ReLU}(v_i^T x + b_i)$ with $w_i = (v_i, b_i)$. Since the target f is intrinsically one dimensional, all these points should be aligned with the x axis. While this is not the case for unregularized training, dropout does achieve partial alignment. Interestingly, the aligned and nonaligned points belong to different plateaus in the equidistribution plots, as indicated by their color. Since this quantity is easily computable, this may be exploited for new neuron pruning and growing strategies in future work.

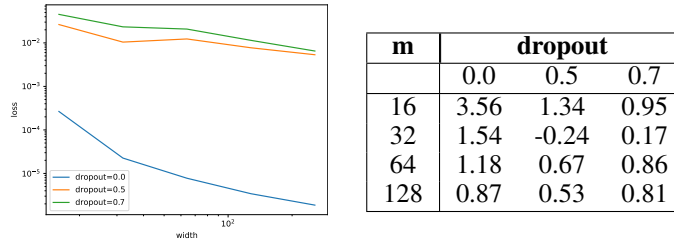


Figure 1: Estimated convergence rates for varying width m and dropout.

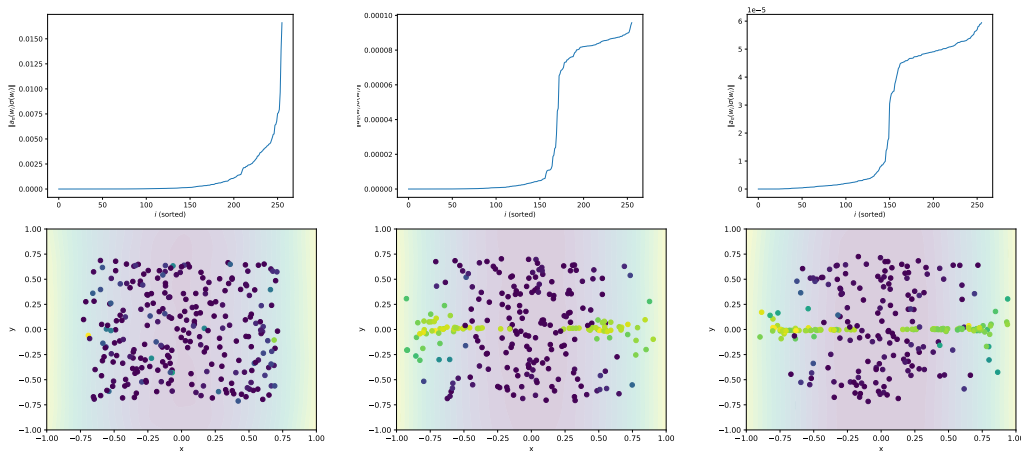


Figure 2: Left to right: Dropout 0.0, 0.5 and 0.7. Top: Equidistribution (4) evaluated at the inner weights $w_i = (v_i, b_i)$. Bottom dots: inner weight v_i of $\text{ReLU}(v_i^T x - b_i)$, colored by equidistribution (4). Bottom background: target function f .

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

REFERENCES

- Z. Allen-Zhu, Y. Li, and Z. Song. A convergence theory for deep learning via over-parameterization. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, page 242–252, Long Beach, California, USA, 09–15 Jun 2019. PMLR. Full version available at <https://arxiv.org/abs/1811.03962>.
- S. Arora, S. Du, W. Hu, Z. Li, and R. Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, page 322–332, Long Beach, California, USA, 09–15 Jun 2019a. PMLR.
- S. Arora, S. S. Du, W. Hu, Z. Li, R. R. Salakhutdinov, and R. Wang. On exact computation with an infinitely wide neural net. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019b.
- F. Bach. Breaking the curse of dimensionality with convex neural networks. *Journal of Machine Learning Research*, 18(19):1–53, 2017.
- F. Bach and L. Chizat. Gradient descent on infinitely wide neural networks: Global convergence and generalization, 2021. <https://arxiv.org/abs/2110.08084>.
- Y. Bai and J. D. Lee. Beyond linearization: On quadratic and higher-order approximation of wide neural networks. In *International Conference on Learning Representations*, 2020.
- A. R. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information Theory*, 39(3):930–945, 1993.
- J. Berner, P. Grohs, G. Kutyniok, and P. Petersen. The Modern Mathematics of Deep Learning. In P. Grohs and G. Kutyniok, editors, *Mathematical Aspects of Deep Learning*, page 1–111. Cambridge University Press, 1 edition, 2022. ISBN 9781009025096 9781316516782.
- G. Bresler and D. Nagaraj. Sharp representation theorems for ReLU networks with precise dependence on depth. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, page 10697–10706. Curran Associates, Inc., 2020.
- Y. Cao and Q. Gu. Generalization Error Bounds of Gradient Descent for Learning Over-Parameterized Deep ReLU Networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04):3349–3356, Apr. 2020.
- F. Chen, Z. Ren, and S. Wang. Uniform-in-time propagation of chaos for mean field langevin dynamics, 2023. <https://arxiv.org/abs/2212.03050>.
- Z. Chen, Y. Cao, D. Zou, and Q. Gu. How much over-parameterization is sufficient to learn deep re{lu} networks? In *International Conference on Learning Representations*, 2021.
- L. Chizat. Mean-field langevin dynamics : Exponential convergence and annealing. *Transactions on Machine Learning Research*, 2022.
- L. Chizat and F. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- L. Chizat, E. Oyallon, and F. Bach. On lazy training in differentiable programming. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- A. Cohen and J.-M. Mirebeau. Adaptive and anisotropic piecewise polynomial approximation. In R. DeVore and A. Kunoth, editors, *Multiscale, Nonlinear and Adaptive Approximation*, page 75–135, Berlin, Heidelberg, 2009. Springer Berlin Heidelberg. ISBN 978-3-642-03413-8.

-
- 594 G. Cybenko. Approximation by superpositions of a sigmoidal function. *Math. Control Signal*
595 *Systems*, 2:303–314, 1989.
- 596
- 597 A. Damian, J. Lee, and M. Soltanolkotabi. Neural networks can learn representations with gradi-
598 ent descent. In P.-L. Loh and M. Raginsky, editors, *Proceedings of Thirty Fifth Conference on*
599 *Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages 5413–5452.
600 PMLR, 02–05 Jul 2022.
- 601 I. Daubechies, R. DeVore, S. Foucart, B. Hanin, and G. Petrova. Nonlinear Approximation and
602 (Deep) ReLU Networks. *Constructive Approximation*, 55(1):127–172, Feb. 2022.
- 603
- 604 R. DeVore, B. Hanin, and G. Petrova. Neural network approximation. *Acta Numerica*, 30:327–444,
605 2021.
- 606 R. A. DeVore. Nonlinear approximation. *Acta Numerica*, 7:51–150, 1998.
- 607
- 608 T. Ding, D. Li, and R. Sun. Suboptimal Local Minima Exist for Wide Neural Networks with Smooth
609 Activations. *Mathematics of Operations Research*, 47(4):2784–2814, Nov. 2022.
- 610 S. Drews and M. Kohler. On the universal consistency of an over-parametrized deep neural network
611 estimate learned by gradient descent, 2022. <https://arxiv.org/abs/2208.14283>.
- 612
- 613 S. Du and J. Lee. On the power of over-parametrization in neural networks with quadratic activation.
614 In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine*
615 *Learning*, volume 80 of *Proceedings of Machine Learning Research*, page 1329–1338. PMLR,
616 10–15 Jul 2018. <https://arxiv.org/abs/1803.01206>.
- 617 S. Du, J. Lee, H. Li, L. Wang, and X. Zhai. Gradient descent finds global minima of deep neural
618 networks. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International*
619 *Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*,
620 page 1675–1685, Long Beach, California, USA, 09–15 Jun 2019a. PMLR.
- 621 S. S. Du, X. Zhai, B. Póczos, and A. Singh. Gradient descent provably optimizes over-parameterized
622 neural networks. In *International Conference on Learning Representations*, 2019b.
- 623
- 624 R. Ge, J. D. Lee, and T. Ma. Learning one-hidden-layer neural networks with landscape design. In
625 *International Conference on Learning Representations*, 2018. [https://arxiv.org/abs/](https://arxiv.org/abs/1711.00501)
626 [1711.00501](https://arxiv.org/abs/1711.00501).
- 627 R. Gentile and G. Welper. Approximation results for Gradient Flow Trained Shallow Neural Net-
628 works in 1d. *Constructive Approximation*, 60(3):547–594, Dec. 2024.
- 629
- 630 R. Gribonval, G. Kutyniok, M. Nielsen, and F. Voigtlaender. Approximation Spaces of Deep Neural
631 Networks. *Constructive Approximation*, 55(1):259–367, Feb. 2022.
- 632
- 633 I. Gühring, G. Kutyniok, and P. Petersen. Error bounds for approximations with deep ReLU neural
634 networks in $w_{s,p}$ norms. *Analysis and Applications*, 18(05):803–859, 2020.
- 635 B. Hanin and M. Sellke. Approximating continuous functions by ReLU nets of minimal width.
636 <https://arxiv.org/abs/1710.11278>, 2017.
- 637
- 638 F. He, B. Wang, and D. Tao. Piecewise linear activations substantially shape the loss surfaces of
639 neural networks. In *International Conference on Learning Representations*, 2020.
- 640 K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approxi-
641 mators. *Neural Networks*, 2(5):359–366, 1989.
- 642
- 643 K. Hu, Z. Ren, D. Siska, and L. Szpruch. Mean-field langevin dynamics and energy landscape of
644 neural networks, 2020. <https://arxiv.org/abs/1905.07769>.
- 645 S. Ibragimov, A. Jentzen, and A. Riekert. Convergence to good non-optimal critical points in the
646 training of neural networks: Gradient descent optimization with one random initialization over-
647 comes all bad non-global local minima with high probability, 2022. [https://arxiv.org/](https://arxiv.org/abs/2212.13111)
[abs/2212.13111](https://arxiv.org/abs/2212.13111).

-
- 648 A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in
649 neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and
650 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran As-
651 sociates, Inc., 2018.
- 652 A. Jentzen and A. Riekert. A proof of convergence for the gradient descent optimization method
653 with random initializations in the training of neural networks with relu activation for piecewise
654 linear target functions. *Journal of Machine Learning Research*, 23(260):1–50, 2022.
- 655 A. Jentzen and A. Riekert. Non-convergence to global minimizers for adam and stochastic gradi-
656 ent descent optimization and constructions of local minimizers in the training of artificial neural
657 networks, 2024. <https://arxiv.org/abs/2402.05155>.
- 658 Z. Ji and M. Telgarsky. Polylogarithmic width suffices for gradient descent to achieve arbitrarily
659 small test error with shallow ReLU networks. In *International Conference on Learning Repre-*
660 *sentations*, 2020.
- 661 K. Kawaguchi. Deep learning without poor local minima. In D. Lee, M. Sugiyama, U. Luxburg,
662 I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, vol-
663 ume 29. Curran Associates, Inc., 2016.
- 664 K. Kawaguchi and J. Huang. Gradient descent finds global minima for generalizable deep neu-
665 ral networks of practical sizes. In *2019 57th Annual Allerton Conference on Communication,*
666 *Control, and Computing (Allerton)*, page 92–99, 2019.
- 667 K. Kawaguchi, J. Huang, and L. P. Kaelbling. Every Local Minimum Value Is the Global Minimum
668 Value of Induced Model in Nonconvex Machine Learning. *Neural Computation*, 31(12):2293–
669 2323, 12 2019.
- 670 J. M. Klusowski and A. R. Barron. Approximation by combinations of ReLU and squared ReLU
671 ridge functions with ℓ^1 and ℓ^0 controls. *IEEE Transactions on Information Theory*, 64(12):
672 7649–7656, 2018.
- 673 M. Kohler and A. Krzyzak. Analysis of the rate of convergence of an over-parametrized deep neural
674 network estimate learned by gradient descent, 2022. [https://arxiv.org/abs/2210.](https://arxiv.org/abs/2210.01443)
675 01443.
- 676 J. Lee, L. Xiao, S. Schoenholz, Y. Bahri, R. Novak, J. Sohl-Dickstein, and J. Pennington. Wide
677 neural networks of any depth evolve as linear models under gradient descent. In H. Wallach,
678 H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in*
679 *Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 680 J. Lee, S. Schoenholz, J. Pennington, B. Adlam, L. Xiao, R. Novak, and J. Sohl-Dickstein. Finite
681 versus infinite neural networks: an empirical study. In H. Larochelle, M. Ranzato, R. Hadsell,
682 M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33,
683 page 15156–15172. Curran Associates, Inc., 2020.
- 684 J. Lee, J. Y. Choi, E. K. Ryu, and A. No. Neural tangent kernel analysis of deep narrow neural
685 networks. In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, editors,
686 *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceed-*
687 *ings of Machine Learning Research*, page 12282–12351. PMLR, 17–23 Jul 2022.
- 688 J. D. Lee, K. Oko, T. Suzuki, and D. Wu. Neural network learns low-dimensional polynomials with
689 sgd near the information-theoretic limit, 2024. <https://arxiv.org/abs/2406.01581>.
- 690 B. Li, S. Tang, and H. Yu. Better approximations of high dimensional smooth functions by deep
691 neural networks with rectified power units. *Communications in Computational Physics*, 27(2):
692 379–411, 2019.
- 693 Y. Li and Y. Liang. Learning overparameterized neural networks via stochastic gradient descent
694 on structured data. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and
695 R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, page 8157–8166.
696 Curran Associates, Inc., 2018.

-
- 702 Z. Li, C. Ma, and L. Wu. Complexity measures for neural networks with general activation functions
703 using path-based norms, 2020. <https://arxiv.org/abs/2009.06132>.
704
- 705 J. Lu, Z. Shen, H. Yang, and S. Zhang. Deep network approximation for smooth functions. *SIAM*
706 *Journal on Mathematical Analysis*, 53(5):5465–5506, 2021.
- 707 Z. Lu, H. Pu, F. Wang, Z. Hu, and L. Wang. The expressive power of neural networks: A view from
708 the width. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and
709 R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, page 6231–6239.
710 Curran Associates, Inc., 2017.
- 711 P. Marion and R. Berthier. Leveraging the two-timescale regime to demonstrate convergence of
712 neural networks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
713
- 714 S. Mei, A. Montanari, and P.-M. Nguyen. A mean field view of the landscape of two-layer neu-
715 ral networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
716 <https://arxiv.org/abs/1804.06561>.
- 717 A. Mousavi-Hosseini, S. Park, M. Girotti, I. Mitliagkas, and M. A. Erdogdu. Neural networks
718 efficiently learn low-dimensional representations with SGD. In *The Eleventh International Con-*
719 *ference on Learning Representations*, 2023.
- 720 Q. Nguyen and M. Hein. The loss surface of deep and wide neural networks. In D. Precup and Y. W.
721 Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70
722 of *Proceedings of Machine Learning Research*, page 2603–2612. PMLR, 06–11 Aug 2017.
723
- 724 Q. N. Nguyen and M. Mondelli. Global convergence of deep networks with one wide layer followed
725 by pyramidal topology. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, edi-
726 tors, *Advances in Neural Information Processing Systems*, volume 33, page 11961–11972. Curran
727 Associates, Inc., 2020.
- 728 A. Nitanda, D. Wu, and T. Suzuki. Convex analysis of the mean field langevin dynamics. In
729 G. Camps-Valls, F. J. R. Ruiz, and I. Valera, editors, *Proceedings of The 25th International Con-*
730 *ference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning*
731 *Research*, pages 9741–9757. PMLR, 2022.
- 732 S. Nitish, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way
733 to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(56):
734 1929–1958, 2014.
- 735 G. Ongie, R. Willett, D. Soudry, and N. Srebro. A function space view of bounded norm infinite
736 width relu nets: The multivariate case. In *International Conference on Learning Representations*,
737 2020.
- 738 J. A. A. Opschoor, P. C. Petersen, and C. Schwab. Deep ReLU networks and high-order finite
739 element methods. *Analysis and Applications*, 18(05):715–770, 2020.
- 740
- 741 S. Oymak and M. Soltanolkotabi. Toward moderate overparameterization: Global convergence
742 guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information*
743 *Theory*, 1(1):84–105, 2020.
- 744
- 745 R. Parhi and R. D. Nowak. Banach space representer theorems for neural networks and ridge splines.
746 *Journal of Machine Learning Research*, 22(43):1–40, 2021.
- 747
- 748 A. Pinkus. Approximation theory of the mlp model in neural networks. *Acta Numerica*, 8:143–195,
749 1999.
- 750 G. M. Rotskoff and E. Vanden-Eijnden. Neural networks as interacting particle systems: Asymp-
751 totic convexity of the loss landscape and universal scaling of the approximation error. *CoRR*,
752 abs/1805.00915, 2018. <https://arxiv.org/abs/1805.00915>.
- 753 I. Safran and O. Shamir. Spurious local minima are common in two-layer ReLU neural networks.
754 In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine*
755 *Learning*, volume 80 of *Proceedings of Machine Learning Research*, page 4433–4441, Stock-
holmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.

-
- 756 M. Seleznova and G. Kutyniok. Analyzing finite neural networks: Can we trust neural tangent
757 kernel theory? In J. Bruna, J. Hesthaven, and L. Zdeborova, editors, *Proceedings of the 2nd*
758 *Mathematical and Scientific Machine Learning Conference*, volume 145 of *Proceedings of Ma-*
759 *chine Learning Research*, page 868–895. PMLR, 16–19 Aug 2022.
- 760 Z. Shen, H. Yang, and S. Zhang. Nonlinear approximation via compositions. *Neural Networks*, 119:
761 74–84, 2019.
- 762 J. W. Siegel and J. Xu. Approximation rates for neural networks with general activation functions.
763 *Neural Networks*, 128:313–321, 2020.
- 764 J. W. Siegel and J. Xu. Characterization of the variation spaces corresponding to shallow neural
765 networks, 2021. <https://arxiv.org/abs/2106.15002>.
- 766 J. W. Siegel and J. Xu. High-order approximation rates for shallow neural networks with cosine and
767 ReLU^k activation functions. *Applied and Computational Harmonic Analysis*, 58:1–26, 2022a.
- 768 J. W. Siegel and J. Xu. Optimal convergence rates for the orthogonal greedy algorithm. *IEEE*
769 *Transactions on Information Theory*, 68(5):3354–3361, 2022b.
- 770 J. W. Siegel and J. Xu. Sharp Bounds on the Approximation Rates, Metric Entropy, and n-Widths
771 of Shallow Neural Networks. *Foundations of Computational Mathematics*, 24(2):481–537, Apr.
772 2024.
- 773 J. W. Siegel, Q. Hong, X. Jin, W. Hao, and J. Xu. Greedy training algorithms for neural networks
774 and applications to PDEs. *Journal of Computational Physics*, 484:112084, July 2023.
- 775 J. Sirignano and K. Spiliopoulos. Mean field analysis of neural networks: A law of large numbers.
776 *SIAM Journal on Applied Mathematics*, 80(2):725–752, 2020.
- 777 M. Soltanolkotabi, A. Javanmard, and J. D. Lee. Theoretical insights into the optimization landscape
778 of over-parameterized shallow neural networks. *IEEE Transactions on Information Theory*, 65
779 (2):742–769, 2019. <https://arxiv.org/abs/1707.04926>.
- 780 C. Song, A. Ramezani-Kebrya, T. Pethick, A. Eftekhari, and V. Cevher. Subquadratic overparame-
781 terization for shallow neural networks. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and
782 J. W. Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, page
783 11247–11259. Curran Associates, Inc., 2021.
- 784 Z. Song and X. Yang. Quadratic suffices for over-parametrization via matrix chernoff bound, 2019.
785 <https://arxiv.org/abs/1906.03593>.
- 786 D. Soudry and Y. Carmon. No bad local minima: Data independent training error guarantees for
787 multilayer neural networks, 2016. <https://arxiv.org/abs/1605.08361>.
- 788 L. Su and P. Yang. On learning over-parameterized neural networks: A functional approximation
789 perspective. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Gar-
790 nett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates,
791 Inc., 2019.
- 792 T. Suzuki. Adaptivity of deep reLU network for learning in besov and mixed smooth besov spaces:
793 optimal rate and curse of dimensionality. In *International Conference on Learning Representa-*
794 *tions*, 2019.
- 795 T. Suzuki, A. Nitanda, and D. Wu. Uniform-in-time propagation of chaos for the mean-field gradient
796 langevin dynamics. In *The Eleventh International Conference on Learning Representations*, 2023.
- 797 G. Swirszcz, W. M. Czarnecki, and R. Pascanu. Local minima in training of neural networks, 2017.
798 <https://arxiv.org/abs/1611.06310v2>.
- 799 S. Takakura and T. Suzuki. Mean-field analysis on two-layer neural networks from a kernel perspec-
800 tive. In *Forty-first International Conference on Machine Learning*, 2024.

810 M. Unser. Ridges, neural networks, and the radon transform. *Journal of Machine Learning Research*,
811 24(37):1–33, 2023.

812 L. Venturi, A. S. Bandeira, and J. Bruna. Spurious valleys in one-hidden-layer neural network
813 optimization landscapes. *Journal of Machine Learning Research*, 20(133):1–34, 2019. <https://arxiv.org/abs/1802.06384>.

814
815

816 N. Vyas, Y. Bansal, and P. Nakkiran. Empirical limitations of the NTK for understanding scaling
817 laws in deep learning. *Transactions on Machine Learning Research*, 2023.

818 E. Weinan, M. Chao, W. Lei, and S. Wojtowytsch. Towards a mathematical understanding of neural
819 network-based machine learning: What we know and what we don’t. *SIAM Transactions on*
820 *Applied Mathematics*, 1(4):561–615, 2020.

821 E. Weinan, C. Ma, and L. Wu. The Barron Space and the Flow-Induced Function Spaces for Neural
822 Network Models. *Constructive Approximation*, 55(1):369–406, Feb. 2022.

823
824

825 G. Welper. Approximation results for gradient flow trained neural networks. *Journal of Machine*
826 *Learning*, 3(2):107–175, 2024a.

827 G. Welper. Nonlinear behaviour of critical points for a simple neural network. *Transactions on*
828 *Machine Learning Research*, 2024b.

829 G. Welper. Approximation and gradient descent training with neural networks. *Sampling Theory,*
830 *Signal Processing, and Data Analysis*, 23(2):20, Dec. 2025.

831 G. Welper and B. Keene. Approximation, estimation and optimization errors for a deep neural
832 network. *Transactions on Machine Learning Research*, 2025.

833
834

835 D. Yarotsky. Error bounds for approximations with deep ReLU networks. *Neural Networks*, 94:
836 103–114, 2017.

837 D. Yarotsky. Optimal approximation of continuous functions by very deep ReLU networks. In
838 S. Bubeck, V. Perchet, and P. Rigollet, editors, *Proceedings of the 31st Conference On Learning*
839 *Theory*, volume 75 of *Proceedings of Machine Learning Research*, page 639–649. PMLR, 06–09
840 Jul 2018.

841 D. Yarotsky and A. Zhevnerchuk. The phase diagram of approximation rates for deep neural net-
842 works. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in*
843 *Neural Information Processing Systems*, volume 33, page 13005–13015. Curran Associates, Inc.,
844 2020.

845 D.-X. Zhou. Universality of deep convolutional neural networks. *Applied and Computational Har-*
846 *monic Analysis*, 48(2):787–794, 2020.

847 D. Zou and Q. Gu. An improved analysis of training over-parameterized deep neural networks.
848 In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors,
849 *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

850
851

852 D. Zou, Y. Cao, D. Zhou, and Q. Gu. Gradient descent optimizes over-parameterized deep ReLU
853 networks. *Machine Learning*, 109(3):467 – 492, 2020.

854
855
856
857
858
859
860
861
862
863

864 A NOTATIONS

865 A.1 STANDARD NOTATIONS

866 For $m, M \in \mathbb{N}$, we define $[m] := \{1, \dots, m\}$. For a vector $\mathbf{w} = [w_1, \dots, w_M]$ and index set
867 $\Lambda \subset [M]$, we write $\mathbf{w}_{\setminus \Lambda} := [w_i]_{i \in [M] \setminus \Lambda}$. For small sets we abbreviate $\mathbf{w}_{\setminus i} = \mathbf{w}_{\setminus \{i\}}$, $\mathbf{w}_{\setminus i, j} =$
868 $\mathbf{w}_{\setminus \{i, j\}}$ etc. For $a_\pi: \mathcal{W} \rightarrow \mathbb{R}$, the expression $a_\pi(\mathbf{w})$, applied to a vector $\mathbf{w} \in \mathcal{W}^M$ is evaluated
869 component-wise.

870 Throughout the paper \mathcal{H} denotes a Hilbert space with norm $\|\cdot\|$ and inner product $\langle \cdot, \cdot \rangle$. Finite
871 signed measures on \mathcal{W} are denoted by \mathcal{M} , restriction to positive measures by \mathcal{M}_+ and probability
872 measures by $\mathcal{M}_{+,1}$. For measures $\mu, \nu \in \mathcal{M}$, we denote the positive and negative parts by μ_+ and
873 μ_- and the variation by $|\mu|$. We abbreviate $\mu(f) := \int f(x) d\mu(x)$. We denote absolute continuity
874 by $\mu \ll \nu$ and singularity by $\mu \perp \nu$. The product measure is $\mu^m := \bigotimes_{i=1}^m \mu$. We use standard
875 $L^p(\pi)$ spaces for functions $\mathcal{W} \rightarrow \mathbb{R}$ and Bochner spaces $L^p(\pi; \mathcal{H})$ for functions $\mathcal{W} \rightarrow \mathcal{H}$. We
876 denote the Gateaux derivative of $\ell(\pi)$ in direction ν by $D[\ell(\pi)]\nu$.

877 A.2 FREE AND BOUND VARIABLES

878 We make judicious use of free and bound variables. For example in the expression

$$879 w \int_0^1 \sin(w) dw$$

880 the outer w is a free variable and can be replaced with any number. The inner w is bound by the
881 integral and can be renamed

$$882 w \int_0^1 \sin(v) dv$$

883 to avoid name collisions. The convention is typically used as follows: For some function $a: \mathbb{R}^M \rightarrow$
884 \mathbb{R} , the expression $\mathbb{E}_{\mathbf{w}_{\setminus 1}} [a(\mathbf{w})]$ has bound variables w_2, \dots, w_M and free variable w_1 . This allows us
885 to make sense of expressions like $\nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [a(\mathbf{w})])$ for an arbitrary measure ν . The outer $\nu(\cdot)$ denotes
886 an integral with respect to ν and takes in a function. The inner expression is a function of the single
887 free variable $w_1 \rightarrow \mathbb{E}_{\mathbf{w}_{\setminus 1}} [a(\mathbf{w})]$. We automatically bound the ν integration to the only free variable
888 so that

$$889 \nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [a(\mathbf{w})]) = \int \mathbb{E}_{w_2, \dots, w_M} [a(w_1, w_2, \dots, w_M)] d\nu(w_1).$$

890 This convention is analogous to free and bound variables in mathematical logic and after getting
891 used to, it provides a compact notation that avoids a serious amount of indices and technicalities
892 throughout this paper.

893 B PERTURBATIONS

894 The main results Lemma 3.4 and Theorem 3.5 demonstrate equidistribution and approximation sub-
895 ject to two perturbation terms Δ_1 and Δ_2 . In this section, we consider some preliminary simplifica-
896 tions and estimates.

897 For both Δ_1 and Δ_2 , we need an estimate for $\mathbb{E}_{\mathbf{w}_{\setminus 1}} [\|f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}\|]$. Standard approximation results
898 from Theorem 2.1 yield

$$899 \mathbb{E}_{\mathbf{w}_{\setminus 1}} [\|f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}\|] \leq \mathbb{E}_{\mathbf{w}_{\setminus 1}} [\|f_{\mathbf{w}} - f\|] + \mathbb{E}_{\mathbf{w}_{\setminus 1}} [\|f - f_{\mathbf{w}_{\setminus 1}}\|] = \mathcal{O}(M^{-1/2}),$$

900 which is insufficient because both Δ_1 and Δ_2 contain a factor of M . However, the estimate is rather
901 crude since usually $f_{\mathbf{w}_{\setminus 1}}$ is much closer to $f_{\mathbf{w}}$ than f , which is exploited in the following lemma
902 and yields an improved error of $\mathcal{O}(M^{-1})$. The proof is in Appendix H.2.

903 **Lemma B.1.** *Let \mathbf{a} and $f_{\mathbf{w}}$ be given by (15). Then, we have*

$$904 \mathbb{E}_{\mathbf{w}_{\setminus 1}} \|f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}\| \leq M^{-1} \mathbb{E}_{\mathbf{w}_{\setminus 1}} [\|\mathbf{a}_1(\mathbf{w})\|] \|\sigma(w_1)\|.$$

The main obstacle to proceed further is the absolute value around \mathbf{a}_1 . Since $f_{\mathbf{w}} \rightarrow f$ for $M \rightarrow \infty$, one may expect that also $\mathbf{a}_1(\mathbf{w}) \rightarrow a_\pi(w_1)$. Then, in the limit coefficients are non-oscillatory in the sense that $\mathbb{E}_{\mathbf{w}_{\setminus 1}} [|\mathbf{a}(\mathbf{w})_1|] \lesssim |\mathbb{E}_{\mathbf{w}_{\setminus 1}} [a(\mathbf{w})_1]| = |a_\pi(w_1)|$ and the right hand side is roughly the Barron norm $M^{-1} \|f\|_B$. This allows us to further bound the perturbation terms as in the following corollary. To state the result, we use the space

$$\mathcal{H}(\mathbf{w}) := \text{span}\{\sigma(w_i) : i \in [M]\} \quad (31)$$

and likewise for shortened vectors like $\mathbf{w}_{\setminus i}$. For the proof, see Appendix H.3.

Corollary B.2. *Let π be a local minimum of (18) with G_π defined in (22), best approximations $\mathbf{a}_1 \in L^2(\pi)$, $f_{\mathbf{w}}$ defined in (15) and $\mathcal{H}(\mathbf{w})$ defined in (31). Assume*

$$\mathbb{E}_{\mathbf{w}_{\setminus 1}} [|\mathbf{a}_1(\mathbf{w})|] \leq c |\mathbb{E}_{\mathbf{w}_{\setminus 1}} [a_1(\mathbf{w})]|$$

for some $c \geq 0$ and define

$$\begin{aligned} \Delta &:= 2c \left[\inf_{h \in \mathcal{H}(\mathbf{w}_{\setminus 1})} \|G_\pi - h\| + (m-1) \|\mathbb{E}_{\mathbf{w}} [f_{\mathbf{w}} - f]\| \right] \|\bar{a}_\pi \sigma\|_{L^\infty(\pi; \mathcal{H})}, \\ \epsilon &:= \|\mathbb{E}_{\mathbf{w} \sim \pi} [f_{\mathbf{w}} - f]\|. \end{aligned}$$

Then for all $m \in \mathbb{N}$

$$\mathbb{E}_{\mathbf{w} \sim \pi^m} \left\| \frac{1}{m} \sum_{i=1}^m \bar{a}_\pi(w_i) \sigma(w_i) - f \right\|^2 \leq \frac{4}{m} \left[\|f\|_{B(\Delta, \epsilon)}^2 + \pi(\Delta) \right] + \epsilon^2.$$

We may expect that $\|\bar{a}_\pi \sigma\|_{L^\infty(\pi; \mathcal{H})}$ is bounded, that $\|\mathbb{E}_{\mathbf{w}} [f_{\mathbf{w}} - f]\| \lesssim M^{-1/2}$ and, similar to approximation by neural networks, that $\inf_{h \in \mathcal{H}(\mathbf{w}_{\setminus 1})} \|G_\pi - h\| \lesssim M^{-1/2}$. Inspecting the proof, this yields the perturbations $\Delta = |\Delta_1 + \Delta_2|$ with $|\Delta_1| \lesssim M^{-1/2}$, $|\Delta_2| \lesssim mM^{-1/2}$. In conclusion, for non-oscillatory coefficients, and M sufficiently larger than m , the perturbation terms are negligible and we obtain Barron type approximation results.

C MAUREY SAMPLING AND BARRON SPACES

C.1 PROOF OF THEOREM 2.1 AND COROLLARY 2.2

In this section, we prove Theorem 2.1 and Corollary 2.2. These are minor adaptations from Theorem 4 in Weinan et al. (2022) and Theorem 1 in Siegel and Xu (2024).

Lemma C.1. *Let $X_i \in \mathcal{H}$ be mean zero independent random variables. Then for independent copies X'_i of X_i*

$$\mathbb{E} \left\| \frac{1}{m} \sum_{i=1}^m X_i \right\|^2 = \frac{1}{2m^2} \sum_{i=1}^m \mathbb{E} \mathbb{E}' \|X_i - X'_i\|^2 = \frac{1}{m^2} \sum_{i=1}^m \mathbb{E} \|X_i\|^2$$

Proof. For the first and last term we have

$$\mathbb{E} \left\| \frac{1}{m} \sum_{i=1}^m X_i \right\|^2 = \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m \mathbb{E} \langle X_i, X_j \rangle = \frac{1}{m^2} \sum_{i=1}^m \mathbb{E} \|X_i\|^2.$$

For the middle term we have

$$\frac{1}{2m^2} \sum_{i=1}^m \mathbb{E} \mathbb{E}' \|X_i - X'_i\|^2 = \frac{1}{2m^2} \sum_{i=1}^m \mathbb{E} \mathbb{E}' \left[\|X_i\|^2 + \|X'_i\|^2 - 2 \langle X_i, X'_i \rangle \right] = \frac{1}{m^2} \sum_{i=1}^m \mathbb{E} \|X_i\|^2,$$

which concludes the proof. \square

972 *Proof of Theorem 2.1.* Define the random variables $X_i := a_\pi(w_i)\sigma(w_i) - f$. Then by construction
 973 (9) the random variables have zero mean $\mathbb{E}[X_i] = 0$, are i.i.d. and by Lemma C.1 we have

$$974 \mathbb{E} \left\| \frac{1}{m} \sum_{i=1}^m X_i \right\|^2 \leq \frac{1}{2m^2} \mathbb{E} \mathbb{E}' \sum_{i=1}^m \|X_i - X'_i\|^2 \leq \frac{2}{m^2} \mathbb{E} \sum_{i=1}^m \|a_\pi(w_i)\sigma(w_i)\|^2 = \frac{2}{m} \|f\|_{B(\pi)}^2,$$

975 where the last equality follows from the definition (11) of the Barron norm.

976 \square

977 *Proof of Corollary 2.2.* First note that the normalization factor $|f|_B^{-1}$ ensures that $\pi(\mathcal{W}) = 1$ so that
 978 π is a probability measure. Since for all continuous functions f by the definition of π we have

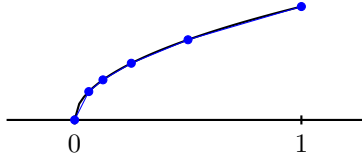
$$979 \int f(w) d|\phi|(w) = \int f(w) \frac{|f|_B}{\|\sigma(w)\|} \frac{\|\sigma(w)\|}{|f|_B} d|\phi|(w) = \int f(w) \frac{|f|_B}{\|\sigma(w)\|} d\pi(w),$$

980 we must have $\frac{d|\phi|}{d\pi} = \frac{|f|_B}{\|\sigma(w)\|}$, almost surely and thus $|a_\pi| = \left| \frac{d\phi}{d\pi} \right| = \frac{d|\phi|}{d\pi} = \frac{|f|_B}{\|\sigma(w)\|}$. It follows that
 981 $|a_\pi(w)| \|\sigma(w)\| = |f|_B$ almost surely, which directly implies $|f|_{B(\pi)} = |f|_B$ and the stated error
 982 bound by Theorem 2.1.

983 \square

984 C.2 EXAMPLE

985 To better understand the role of π , it is instructive to consider the approximation of a target function
 986 $f: [0, 1] \rightarrow \mathbb{R}$ in one dimension with $\sigma(w) := \text{ReLU}(x - w)$. In this setup, shallow networks of
 987 width m are identical to much better understood linear splines with m free knots at locations w_i .
 988 We would like to place most knots, where f is complicated, measured by high derivatives, e.g. for
 989 $f(x) = \sqrt{x}$ as follows:



990 Such an adaptive placement of knots is achieved in the construction of Barron approximation results
 991 in Corollary 2.2, which can be explicitly computed as follows. For ϕ given by density $d\phi(w) =$
 992 $\rho_\phi(w) dw$ so that $f(x) = \int \text{ReLU}(x - w)\rho_\phi(w) dw$, differentiating twice, yields

$$993 f''(x) = \int \delta(x - w)\rho_\phi(w) dw = \rho_\phi(x)$$

994 and thus the density $\rho_\phi = f''$ is the second derivative. Then the sample distribution π is given
 995 by density $\rho_\pi(w) \sim \|\sigma(w)\| |\rho_\phi(w)| \sim \|\sigma(w)\| |f''(w)|$. Note carefully that regions where f is
 996 difficult to approximate, indicated by a large derivative $|f''(w)|$, have equally large density. Hence,
 997 sampling places most knots w_i in difficult regions as desired.

998 For other sample distributions π , the approximation results can be worse or even fail. Indeed, the
 999 Barron norm with respect to arbitrary π is

$$1000 \|f\|_{B(\pi)} = \mathbb{E} \|a_\pi(w)\sigma(w)\|^2, \quad a_\pi(w) = \frac{d\phi}{d\pi}(w) = \frac{f''(w)}{\rho_\pi(w)},$$

1001 which can be infinite, even if the optimal choice $\|f\|_B$ is finite. E.g. for uniform π and $f(x) =$
 1002 $|x - 1/2|^\alpha$ with $\alpha \in (1.0, 1.5)$. In this case, Theorem 2.1 and Corollary 2.2 provide approximation
 1003 rate $m^{-1/2}$ for sampling from the optimal distribution but none for sampling from the uniform
 1004 distribution.

1005 The argument for this simple model problem relies in the observation that the second derivative of
 a ReLU activation is a Dirac delta. This remains true for standard shallow networks in multiple

dimensions if we replace the second derivative with higher derivatives, a ramp filter and a Radon transform, as shown in Parhi and Nowak (2021, Lemma 17). The problem of finding a good sample distribution π is even more important in high dimensions as it must reveal intrinsic low dimensional structures of the target f .

D AUXILIARY RESULTS

D.1 GRADIENTS

This section provides several basic gradients that are used repeatedly in the proofs of the main results. Recall that \mathcal{M} denotes the vector space of finite signed measures on \mathcal{W} and \mathcal{H} denotes a Hilbert space.

Lemma D.1. *For $\pi, \nu \in \mathcal{M}$, let $g \in L^2(|\pi|; \mathcal{H}) \cap L^2(|\nu|; \mathcal{H})$. Then the flowing two Gateaux derivatives are well defined:*

$$\begin{aligned} D_\pi \left[\int \|g(w)\|^2 d\pi(w) \right] \nu &= \int \|g(w)\|^2 d\nu(w), \\ D_\pi \left[\int \langle g(v), g(w) \rangle d\pi^2(v, w) \right] \nu &= 2 \int \langle g(v), g(w) \rangle d(\pi \otimes \nu)(v, w). \end{aligned}$$

Proof. Follows directly from elementary computation of the difference quotients. All necessary integrals are finite by the given assumptions. \square

Lemma D.2. *For $\pi, \nu \in \mathcal{M}$, let $g \in L^2(|\pi|; \mathcal{H}) \cap L^2(|\nu|; \mathcal{H})$. Then the flowing Gateaux derivative is well defined:*

$$\begin{aligned} D_\pi \left[\int \left\| \sum_{i=1}^m g(w_i) \right\|^2 d\pi^m(\mathbf{w}) \right] \nu \\ = m\pi(1)^{m-1} \nu \left(\|g\|^2 \right) + 2m(m-1)\pi(1)^{m-2} \langle \pi(g), \nu(g) \rangle + \lambda \nu(1) \end{aligned}$$

for some constant $\lambda \in \mathbb{R}$.

Proof. Let D be the derivative in the lemma. It follows from Lemma D.1 that the derivative exists with

$$D = \sum_{k=1}^m \int \left\| \sum_{i=1}^m g(w_i) \right\|^2 \prod_{\substack{j=1 \\ j \neq k}}^m d\pi(w_j) d\nu(w_k).$$

By symmetry, we can permute each summand so that k is replaced by 1 and the expression simplifies to

$$D = m \int \int \left\| \sum_{i=1}^m g(w_i) \right\|^2 \prod_{j=2}^m d\pi(w_j) d\nu(w_1) = m(\nu \otimes \pi^{m-1}) \left(\left\| \sum_{i=1}^m g(w_i) \right\|^2 \right).$$

We abbreviate $G(\mathbf{w}_{\setminus 1}) := \sum_{i=2}^m g(w_i)$ and split the norm accordingly

$$\begin{aligned} D &= m(\nu \otimes \pi^{m-1}) \left(\|g(w_1)\|^2 \right) \\ &\quad + 2m(\nu \otimes \pi^{m-1}) \left(\langle g(w_1), G(\mathbf{w}_{\setminus 1}) \rangle \right) \\ &\quad + m(\nu \otimes \pi^{m-1}) \left(\|G(\mathbf{w}_{\setminus 1})\|^2 \right) \\ &=: I + II + III. \end{aligned}$$

The first summand simplifies to

$$I = m\nu \left(\|g\|^2 \right) \pi^{m-1}(1)$$

and the second to

$$II = 2m \langle \nu(g), \pi^{m-1}(G) \rangle = 2m(m-1) \langle \nu(g), \pi(g) \rangle \pi^{m-2}(1),$$

where in the last step we have used that $\pi^{m-1}(G) = \sum_{i=2}^m \pi(g) \pi^{m-2}(1) = (m-1) \pi(g) \pi^{m-2}(1)$. Finally, the last term is

$$III = m\nu(1) \pi^{m-1} \left(\|G(\mathbf{w}_{\setminus 1})\|^2 \right) = \lambda \nu(1),$$

for some constant $\lambda \in \mathbb{R}$ independent of ν . Combining all summands shows the lemma. \square

Lemma D.3. For $\pi, \nu \in \mathcal{M}$, let $g_\pi \in L^2(|\pi|, \mathcal{H})$ be Gateaux differentiable with respect to $\pi \in \mathcal{M}$ and direction $\nu \in \mathcal{M}$ so that $w \rightarrow D_\pi[g_\pi(w)]\nu \in L^2(|\pi|, \mathcal{H})$. Then the flowing Gateaux derivative is well defined:

$$\begin{aligned} & \int D_\pi \left[\left\| \sum_{i=1}^m g_\pi(w_i) \right\|^2 \right] \nu d\pi^n(\mathbf{w}) \\ &= m\pi(1)^{m-1} \int D_\pi \left[\|g_\pi(w)\|^2 \right] \nu d\pi(w) + 2m(m-1)\pi(1)^{m-2} \langle \pi(g_\pi), \pi(D_\pi[g_\pi]\nu) \rangle. \end{aligned}$$

Proof. Since the outer function $\|\cdot\|^2 : \mathcal{H} \rightarrow \mathbb{R}$ is Fréchet differentiable and the inner functions $\pi \rightarrow g_\pi(w_i)$ are Gateaux differentiable, by the chain rule we have

$$D_\pi \left[\left\| \sum_{i=1}^m g_\pi(w_i) \right\|^2 \right] \nu = \sum_{\substack{i=1 \\ j=1}}^m 2 \langle g_\pi(w_i), D_\pi[g_\pi(w_j)]\nu \rangle =: \sum_{\substack{i=1 \\ j=1}}^m D_{ij}.$$

Undoing the derivative for the $i = j$ terms, we obtain

$$D_{ii} = D_\pi \left[\|g_\pi(w_i)\|^2 \right] \nu$$

and thus by integrating and renaming the integration variable $w_i \rightarrow w$

$$\pi^m(D_{ii}) = \int D_\pi \left[\|g_\pi(w)\|^2 \right] \nu d\pi(w) \pi(1)^{m-1}.$$

For $i \neq j$, the variables w_i and w_j in the inner product are independent so that

$$\pi^m(D_{ij}) = 2 \langle \pi(g_\pi), \pi(D_\pi[g_\pi]\nu) \rangle \pi(1)^{m-2}.$$

Noting that the last two right hand sides are independent of i and j , we can eliminate the i, j sum above and obtain the statement of the lemma. \square

Lemma D.4. Let $\pi, \nu \in \mathcal{M}$ with $\pi(1) = 1$ and $\|\sigma(\cdot)\| \in L^\infty(|\pi|)$. Let $g_\pi(w) = a_\pi(w)\sigma(w) - f$ for some $a_\pi \in L^2(|\pi|; \mathcal{H}) \cap L^2(|\nu|; \mathcal{H})$, Gateaux differentiable with respect to $\pi \in \mathcal{M}$ and direction $\nu \in \mathcal{M}$ so that $w \rightarrow D_\pi[a_\pi(w)]\nu \in L^2(|\pi|) \cap L^2(|\nu|)$. Let π, ν , both be absolutely continuous with respect to some measure $\mu \in \mathcal{M}$. Then there is a $\lambda \in \mathbb{R}$ with

$$D_\pi \left[\int \left\| \sum_{i=1}^m g_\pi(w_i) \right\|^2 d\pi^n(\mathbf{w}) \right] \nu = I + II + III + IV + V,$$

with

$$\begin{aligned} I &= -m \int \|a_\pi(w)\sigma(w)\|^2 d\nu(w) \\ II &= 2m \int \left\langle g_\pi(w), D_\pi \left[a_\pi(w)\sigma(w) \frac{d\pi}{d\mu} \right] \nu \right\rangle d\mu(w) \\ III &= 2m(m-1) \langle \pi(g_\pi), \nu(g_\pi) \rangle \\ IV &= 2m(m-1) \langle \pi(g_\pi), \pi(D_\pi[g_\pi]\nu) \rangle \\ V &= \lambda \nu(1). \end{aligned}$$

for some $\lambda \in \mathbb{R}$.

1134 *Proof.* Let D be the derivative in the lemma. The two terms dependent on the argument π are g_π
 1135 and $d\pi^n(\mathbf{w})$ so that, by the product rule we have

$$1136 \quad D = \int D_\pi \left[\left\| \sum_{i=1}^m g_\pi(w_i) \right\|^2 \right] \nu d\pi^n(\mathbf{w}) + D_\pi \left[\int \left\| \sum_{i=1}^m g_\mu(w_i) \right\|^2 d\pi^n(\mathbf{w}) \right] \nu \Big|_{\mu=\pi},$$

1140 where in the second term the inner g_π is replaced by g_μ for some temporary $\mu = \pi$ so that it is not
 1141 differentiated. The derivatives in the right hand side are given by Lemmas D.3 and D.2. Together
 1142 with the normalization $\pi(1) = 1$, they yield

$$1143 \quad D = (i + IV) + (ii + III + V),$$

1144 with

$$1145 \quad \begin{aligned} i &= m \int D_\pi \left[\|g_\pi(w)\|^2 \right] \nu d\pi(w), \\ IV &= 2m(m-1) \langle \pi(g_\pi), \pi(D_\pi[g_\pi]\nu) \rangle, \\ ii &= m\nu \left(\|g_\pi\|^2 \right), \\ III &= 2m(m-1) \langle \pi(g_\pi), \nu(g_\pi) \rangle, \\ V &= \lambda\nu(1). \end{aligned}$$

1155 The terms III , IV and V match the corresponding terms in the conclusion of the lemma. The terms
 1156 i and ii require some attention. For the former, using the chain rule and $D_\pi f = 0$, we obtain

$$1157 \quad i = 2m \int \langle g_\pi(w), D_\pi[g_\pi(w)]\nu \rangle d\pi(w) = 2m \int \langle g_\pi(w), D_\pi[a_\pi(w)\sigma(w)]\nu \rangle d\pi(w).$$

1160 Since π is absolutely continuous with respect to μ , there is a Radon-Nikodym derivative $d\pi = \frac{d\pi}{d\mu} d\mu$
 1161 so that

$$1162 \quad i = 2m \int \langle g_\pi(w), D_\pi[a_\pi(w)\sigma(w)]\nu \rangle \frac{d\pi}{d\mu} d\mu(w).$$

1163 Then, with $D_\pi \left[\frac{d\pi}{d\mu} \right] \nu d\mu = \frac{d\nu}{d\mu} d\mu = d\nu$ the product rule entails

$$1164 \quad \begin{aligned} i &= 2m \int \left\langle g_\pi(w), D_\pi \left[a_\pi(w)\sigma(w) \frac{d\pi}{d\mu} \right] \nu \right\rangle d\mu(w) \\ &\quad - 2m \int \langle g_\pi(w), a_\pi(w)\sigma(w) \rangle d\nu(w) \\ &=: II + i.ii. \end{aligned}$$

1173 Finally, invoking the definition of $g_\pi = a_\pi\sigma - f$, we conclude that

$$1174 \quad i.ii = -2m \int \|a_\pi(w)\sigma(w)\|^2 d\nu(w) + 2m \int \langle f, a_\pi(w)\sigma(w) \rangle d\nu(w).$$

1175 Likewise, the term ii expands to

$$1176 \quad ii = m \int \|a_\pi(w)\sigma(w)\|^2 d\nu(w) - 2m \int \langle f, a_\pi(w)\sigma(w) \rangle d\nu(w) + m \int \|f\|^2 d\nu(w).$$

1177 Combining $i.ii$ and ii and cancelling terms, we obtain

$$1178 \quad i.ii + ii = -m \int \|a_\pi(w)\sigma(w)\|^2 d\nu(w) + m \int \|f\|^2 d\nu(w) =: I + V.i.$$

1183 Joining $V + V.i$ by adjusting the constant $\lambda \rightarrow \lambda + m \|f\|^2$, and adding all terms proves the lemma.

1184 \square

1188

1189

D.2 FIRST ORDER OPTIMALITY CRITERIA

1190

1191

1192

In this section, we briefly review first order optimality criteria for the optimization of probability measures. To this end, we optimize a function $\ell: \mathcal{M} \rightarrow \mathbb{R}$, confined to positive and normalized measures:

1193

$$\min_{\pi \in \mathcal{M}_{+,1}} \ell(\pi).$$

1194

1195

For positive measures $\nu \in \mathcal{M}_+$, by a short computation and rescaling, the one-sided first order optimality criteria for the convex combinations $t \rightarrow \ell(\pi + t(\nu - \pi))$ at $t = 0$ yield

1196

1197

$$D\ell(\pi)(\nu) + \lambda\nu(1) \geq 0, \quad D\ell(\pi)(\pi) + \lambda\pi(1) = 0, \quad (32)$$

1198

1199

for some constant $\lambda \in \mathbb{R}$. With the definition $-h(\nu) := D\ell(\pi)(\nu) + \lambda\nu(1)$, we can extend the criteria to arbitrary $\nu \in \mathcal{M}$ and obtain standard KKT conditions:

1200

1201

1202

1203

$$\begin{aligned} \text{stationarity:} & \quad D\ell(\pi)(\nu) + \lambda\nu(1) + h(\nu) = 0, \quad \nu \in \mathcal{M} \\ \text{dual feasibility:} & \quad h(\nu) \leq 0, \quad \nu \in \mathcal{M}_+ \\ \text{complementary slackness:} & \quad h(\pi) = 0. \end{aligned}$$

1204

1205

E PROOF OF MAIN RESULTS: EXACT OUTER WEIGHTS

1206

1207

E.1 PROOF OF THEOREM 3.1: CONVEXITY

1208

1209

To show convexity of the loss

1210

1211

$$\ell(\pi) := \mathbb{E}_{\mathbf{w} \sim \pi^m} \left\| \frac{1}{m} \sum_{i=1}^m a_\pi(w_i) \sigma(w_i) - f \right\|^2,$$

1212

1213

we first expand the squared norm

1214

1215

1216

1217

$$\begin{aligned} \ell(\pi) &= \mathbb{E}_{\mathbf{w} \sim \pi^m} \left\| \frac{1}{m} \sum_{i=1}^m a_\pi(w_i) \sigma(w_i) \right\|^2 - 2 \mathbb{E}_{\mathbf{w} \sim \pi^m} \left\langle \frac{1}{m} \sum_{i=1}^m a_\pi(w_i) \sigma(w_i), f \right\rangle + \|f\|^2. \\ &=: I + II + III. \end{aligned}$$

1218

1219

The third term III is constant and therefore convex. For the second term II , we plug in the definition of $a_\pi = d\phi/d\pi$ to obtain

1220

1221

1222

1223

1224

1225

$$\begin{aligned} II &= -2 \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\mathbf{w}_{\setminus i}} \int \langle \sigma(w_i), f \rangle \frac{d\phi}{d\pi}(w_i) d\pi(w_i) \\ &= -2 \frac{1}{m} \sum_{i=1}^m \int \langle \sigma(w_i), f \rangle d\phi(w_i), \end{aligned}$$

1226

1227

where we have dropped the expectation $\mathbb{E}_{\mathbf{w}_{\setminus i}}$ because its argument is independent of $w_{\setminus 1}$. This term is independent of π and thus convex.

1228

1229

1230

It remains to show that the first term is also convex. To this end, we first expand the inner sum and split indices $i = j$ and $i \neq j$

1231

1232

1233

1234

$$\begin{aligned} I &= \frac{1}{m^2} \sum_{i=1}^m \mathbb{E}_{\mathbf{w} \sim \pi^m} \|a_\pi(w_i) \sigma(w_i)\|^2 + \frac{1}{m^2} \sum_{i \neq j} \mathbb{E}_{\mathbf{w} \sim \pi^m} \langle a_\pi(w_i) \sigma(w_i), a_\pi(w_j) \sigma(w_j) \rangle \\ &=: I.1 + I.2. \end{aligned}$$

1235

1236

The term $I.2$ is analogous to II : We plug in the definition of a_π , and remove unnecessary expectations \mathbb{E}_{w_k} for $k \neq i, j$ to obtain

1237

1238

1239

1240

1241

$$\begin{aligned} I.2 &= \frac{1}{m^2} \sum_{i \neq j} \iint \langle \sigma(w_i), \sigma(w_j) \rangle \frac{d\phi}{d\pi}(w_i) \frac{d\phi}{d\pi}(w_j) d\pi(w_i) d\pi(w_j) \\ &= \frac{1}{m^2} \sum_{i \neq j} \iint \langle \sigma(w_i), \sigma(w_j) \rangle d\phi(w_i) d\phi(w_j), \end{aligned}$$

which is independent of π and therefore convex. Noting that all expectations in *I.1* are identical, we simplify to

$$\begin{aligned} mI.1 &= \int \|\sigma(w)\|^2 \left| \frac{d\phi}{d\pi} \right|^2 (w) d\pi(w) \\ &= \int \|\sigma(w)\|^2 \left(\frac{d|\phi|}{d\pi} \right)^2 (w) d\pi(w) = \int \|\sigma(w)\|^2 \frac{d|\phi|}{d\pi} (w) d|\phi|(w). \end{aligned}$$

By the Lebesgue decomposition theorem, we can split $\pi = \pi_0 + \pi_1$ into a part that is absolutely continuous $\pi_0 \ll |\phi|$ and a part that is singular $\pi_1 \perp |\phi|$ with respect to $|\phi|$. Since by assumption $\phi_\pm \ll \phi \ll \pi$, it is not hard to see that $\pi_0 \ll |\phi| \ll \pi_0$ and $d|\phi|/d\pi = d|\phi|/d\pi_0 = (d\pi_0/d|\phi|)^{-1}$ on the support of $|\phi|$, so that

$$mI.1 = \int \|\sigma(w)\|^2 \left(\frac{d\pi_0}{d|\phi|} \right)^{-1} (w) d|\phi|(w).$$

The functions $\pi \rightarrow \pi_0$ and $\pi_0 \rightarrow d\pi_0/d|\phi|$ are linear and since π_0 and $|\phi|$ are non-negative measures, $d\pi_0/d|\phi|$ is non-negative, as well, so that the inverse $x \rightarrow x^{-1}$ is convex. Since all remaining terms in the integral of *I.1* are also non-negative, it follows that *I.1* is convex in π . Together with the convexity of *I.2*, *II* and *III*, this completes the proof.

E.2 PROOF OF LEMMA 3.2: EQUIDISTRIBUTION

We first compute the derivative of the loss

Lemma E.1. *Let $\nu \ll \pi \in \mathcal{M}$ with $\pi(1) = 1$ and $\|\sigma(\cdot)\| \in L^\infty(|\pi|)$, and bounded Radon-Nikodym derivative $d\nu/d\pi$. Let $a_\pi \in L^2(|\pi|)$ be given by (14). Then with $g_\pi(w) = a_\pi(w)\sigma(w) - f$, there is a $\lambda \in \mathbb{R}$ so that*

$$D_\pi \left[\int \left\| \sum_{i=1}^m g_\pi(w_i) \right\|^2 d\pi^n(w) \right] \nu = -m \int \|a_\pi(w)\sigma(w)\|^2 d\nu(w) + \nu(\lambda),$$

Proof. Let D be the derivative in the lemma, which we compute by Lemma D.4. To this end, we first establish its assumptions. Since $d\nu/d\pi$ is bounded, we have

$$\int |a_\pi(w)|^2 d\nu(w) = \int |a_\pi(w)|^2 \frac{d\nu}{d\pi}(w) d\pi(w) \lesssim \int |a_\pi(w)|^2 d\pi(w) < \infty$$

so that $a_\pi \in L^2(|\nu|; \mathcal{H})$. After eventually splitting of singular parts of π with respect to ϕ as in the proof of Theorem 3.1, the derivative $a_\pi = d\phi/d\pi$ is invertible and Gateaux differentiable with

$$D[a_\pi]\nu = D \left[\left(\frac{d\pi}{d\phi} \right)^{-1} \right] \nu = - \left(\frac{d\pi}{d\phi} \right)^{-2} \frac{d\nu}{d\phi} = - \left(\frac{d\phi}{d\pi} \right)^2 \frac{d\nu}{d\phi} = - \frac{d\phi}{d\pi} \frac{d\nu}{d\pi} = -a_\pi \frac{d\nu}{d\pi}.$$

Again using that $d\nu/d\pi$ is bounded, we find $D[a_\pi]\nu \in L^2(|\pi|)$ so that all assumptions of Lemma D.4 are satisfied. The lemma shows that the derivative is of the form $D = I + II + III + IV + V$ and provides formulas for the left hand side. The terms *I* and *V* are unchanged in the formula we wish to prove so that it remains to show that *II*, *III* and *IV* are zero.

III and *IV* vanish because both contain the factor $\pi(g_\pi)$, which is zero. Indeed, since ϕ is absolutely continuous with respect to π , the Radon-Nikodym derivative exists and $a_\pi = \frac{d\phi}{d\pi}$. Hence

$$f = \int \sigma(w) d\phi(w) = \int \sigma(w) \frac{d\phi}{d\pi}(w) d\pi(w) = \int a_\pi(w)\sigma(w) d\pi(w) = \pi(a\sigma)$$

and thus in particular $\pi(g_\pi) = 0$.

The remaining term *II* contains the factor $D_\pi \left[a_\pi(w)\sigma(w) \frac{d\pi}{d\mu} \right] \nu$ for an arbitrary measure μ with respect to which π and ν are absolutely continuous. The particular choice is not important, e.g. $\mu := |\pi| + |\nu|$ is one possibility. Plugging in $a_\pi = d\phi/d\pi$, we obtain

$$D_\pi \left[a_\pi(w)\sigma(w) \frac{d\pi}{d\mu} \right] \nu = D_\pi \left[\sigma(w) \frac{d\phi}{d\pi} \frac{d\pi}{d\mu} \right] \nu = D_\pi \left[\sigma(w) \frac{d\phi}{d\mu} \right] \nu = 0,$$

1296 where in the last step we have used that the term in the bracket is independent of π . Hence also
 1297 $III = 0$, which completes the proof.

1298 □

1300 *Proof of Lemma 3.2.* Let $\nu \in \mathcal{M}_+$ be a positive measure with bounded Radon-Nikodym derivative
 1301 $d\nu/d\pi$. Then, the local minimizer π necessarily satisfies the first order optimality criteria (32)

$$1302 \quad D[\ell(\pi)]\nu + \nu(\lambda) \geq 0,$$

1303 with equality for $\nu = \pi$. Together with the derivative from Lemma E.1, this yields

$$1304 \quad -m \int \|a_\pi(w)\sigma(w)\|^2 d\nu(w) + \nu(\lambda) \geq 0,$$

1305 or equivalently

$$1306 \quad \int \left[-m \|a_\pi(w)\sigma(w)\|^2 + \lambda \right] \frac{d\nu}{d\pi} d\pi(w) \geq 0$$

1307 for all non-negative and bounded densities $d\nu/d\pi$, with equality if $d\nu/d\pi = 1$. Thus, it follows that

$$1308 \quad -m \|a_\pi(w)\sigma(w)\|^2 + \lambda = 0,$$

1309 π -a.e. Dividing by m and redefining λ , concludes the proof.

1310 □

1311

1312 E.3 PROOF OF THEOREM 3.3: APPROXIMATION ERROR

1313 By Theorem 3.1, π is a global minimum. Let π_B be the probability measure from Corollary 2.2.
 1314 Then, we have

$$1315 \quad \ell(\pi) \leq \ell(\pi_B) \leq \frac{2}{m} |f|_B^2.$$

1316

1317 F PROOF OF MAIN RESULTS: AVERAGED OUTER WEIGHTS

1318

1319 The results in this section make heavy use of our conventions on bound and unbound variables
 1320 summarized in Section A.

1321

1322 F.1 SYMMETRY

1323

1324 The coefficients $\mathbf{a}(w)$ in (16) are only confined to a symmetry condition. In this section, we show
 1325 that this condition is satisfied for least squares minimizers as well as several useful properties for the
 1326 proofs of the main results.

1327 **Lemma F.1.** For $f \in \mathcal{H}$, $\sigma: \mathcal{W} \rightarrow \mathcal{H}$ and $w_i \in \mathcal{W}$, $i \in [M]$, define

1328

$$1329 \quad \mathbf{a} = \mathbf{a}(w) = \arg \min^+_{\mathbf{a} \in \mathbb{R}^M} \left\| \frac{1}{M} \sum_{i=1}^M a_i \sigma(w_i) - f \right\|^2,$$

1330 where $\arg \min^+$ picks the candidate with minimal Euclidean norm, in case the minimizer is not
 1331 unique. Then for any permutation matrix $P \in \mathbb{R}^{M \times M}$, we have $P\mathbf{a}(w) = \mathbf{a}(Pw)$.

1332

1333 *Proof.* Abbreviating $\sigma_i = \sigma(w_i)$ and solving the least squares problem, \mathbf{a} satisfies the linear system
 1334 of equations

1335

$$1336 \quad A_w \mathbf{a}(w) = b_w, \quad A_w := \left[\frac{1}{M} \langle \sigma_i, \sigma_j \rangle \right]_{i,j=1}^M, \quad b_w := [\langle \sigma_i, f \rangle]_{i=1}^M.$$

1337

1338 Its solution is invariant under permutations. Indeed, for permutation matrix $P \in \mathbb{R}^{M \times M}$, replacing
 1339 w by Pw , we have $A_{Pw} \mathbf{a}(Pw) = b_{Pw}$. But since $P^T = P^{-1}$ we also have $(PA_w P^T) P\mathbf{a}(w) =$
 1340 Pb_w , or equivalently $A_{Pw} P\mathbf{a}(w) = b_{Pw}$. Hence

1341

$$1342 \quad P\mathbf{a}(w) = \mathbf{a}(Pw) = A_{Pw}^+ b_{Pw},$$

1343

1344 which concludes the proof.

1345 □

1350 **Lemma F.2.** Let $p: [M] \rightarrow [M]$ be a permutation that leaves the complement $[M] \setminus \Lambda$ of a set
 1351 $\Lambda \subset [M]$ invariant. Assume that the corresponding permutation matrix satisfies $P\mathbf{a}(\mathbf{w}) = \mathbf{a}(P\mathbf{w})$
 1352 for some function $\mathbf{a}: \mathcal{W}^M \rightarrow \mathbb{R}$. Let $h: \mathcal{W} \rightarrow \mathcal{H}$ be a function and $\pi_k \in \mathcal{M}_{+,1}$, $k \in \Lambda$ be
 1353 probability measures. Then we have

$$1354 \mathbb{E}_{w_{p_k} \sim \pi_k; k \in \Lambda} [\mathbf{a}_{p_i}(\mathbf{w})h(w_{p_i})] = \mathbb{E}_{w_k \sim \pi_k; k \in \Lambda} [\mathbf{a}_i(\mathbf{w})h(w_i)] \quad \text{for all } i \in \Lambda.$$

1356 *Proof.* Let E be the expectation on the left hand side in the lemma. By the given properties of the
 1357 permutation, we have

$$1358 \mathbf{a}_{p_i}(w_1, \dots, w_M) = \mathbf{a}_i(w_{p_1}, \dots, w_{p_M}),$$

1359 so that

$$1360 E = \int \mathbf{a}_{p_i}(w_1, \dots, w_M)h(w_{p_i}) \prod_{k \in \Lambda} d\pi_k(w_{p_k})$$

$$1361 = \int \mathbf{a}_i(w_{p_1}, \dots, w_{p_M})h(w_{p_i}) \prod_{k \in \Lambda} d\pi_k(w_{p_k}).$$

1362 Since p leaves the complement of Λ invariant, we can change the integration variables names $w_{p_i} \rightarrow$
 1363 w_i to obtain

$$1364 E = \int \mathbf{a}_i(w_1, \dots, w_M)h(w_i) \prod_{k \in \Lambda} d\pi_k(w_k),$$

1365 which concludes the proof. □

1366 **Lemma F.3.** Let $\Lambda \subset [M]$. Let $w_i \in \mathcal{W}$ for $i \notin \Lambda$ and sample $w_i \in \mathcal{W}$, $i \in \Lambda$ i.i.d. from
 1367 probability measure π . Assume $\mathbf{a}: \mathcal{W}^M \rightarrow \mathbb{R}$ satisfies $P\mathbf{a}(\mathbf{w}) = \mathbf{a}(P\mathbf{w})$ for all permutation
 1368 matrices $P \in \mathbb{R}^{M \times M}$ and define $f_{\mathbf{a}(\mathbf{w}), \mathbf{w}} = \frac{1}{M} \sum_{i=1}^M \mathbf{a}_i(\mathbf{w})\sigma(w_i)$. Then for $\mathbb{E}_\Lambda = \mathbb{E}_{w_i \sim \pi, i \in \Lambda}$, we
 1369 have

$$1370 \mathbb{E}_\Lambda [\mathbf{a}_i(\mathbf{w})\sigma(w_i)] = \mathbb{E}_\Lambda [\mathbf{a}_j(\mathbf{w})\sigma(w_j)] \quad \text{for all } i, j \in \Lambda \quad (33)$$

1371 and

$$1372 \mathbb{E}_\Lambda [\mathbf{a}_i(\mathbf{w})\sigma(w_i)] = \frac{M}{|\Lambda|} \mathbb{E}_\Lambda [f_{\mathbf{a}(\mathbf{w}), \mathbf{w}}] - \frac{1}{|\Lambda|} \sum_{j \notin \Lambda} \mathbb{E}_\Lambda [\mathbf{a}_j(\mathbf{w})\sigma(w_j)] \quad \text{for all } i \in \Lambda. \quad (34)$$

1373 *Proof.* We abbreviate $\mathbf{a}_i = \mathbf{a}_i(\mathbf{w})$ and $\sigma_i = \sigma(w_i)$. The first identity (33) follows directly from
 1374 Lemma F.2 with $\pi_k = \pi$, permutation p that swaps i and j and $h = \sigma$. To show (34), we split the
 1375 indices across Λ :

$$1376 \mathbb{E}_\Lambda [f_{\mathbf{a}(\mathbf{w}), \mathbf{w}}] = \frac{1}{M} \sum_{i \in \Lambda} \mathbb{E}_\Lambda [\mathbf{a}_i \sigma_i] + \frac{1}{M} \sum_{i \notin \Lambda} \mathbb{E}_\Lambda [\mathbf{a}_i \sigma_i].$$

1377 By (33) all summands in the sum over $i \in \Lambda$ are equal so that we can choose one index $i \in \Lambda$ and
 1378 collapse the sum:

$$1379 \mathbb{E}_\Lambda [f_{\mathbf{a}(\mathbf{w}), \mathbf{w}}] = \frac{|\Lambda|}{M} \mathbb{E}_\Lambda [\mathbf{a}_i \sigma_i] + \frac{1}{M} \sum_{i \notin \Lambda} \mathbb{E}_\Lambda [\mathbf{a}_i \sigma_i], \quad \text{for all } i \in \Lambda.$$

1380 Rearranging terms shows (34) and concludes the proof. □

1381 **Corollary F.4.** Assume $w_i \in \mathcal{W}$, $i \in [M]$ are sampled i.i.d. from a probability measure π . Assume
 1382 $\mathbf{a}: \mathcal{W}^M \rightarrow \mathbb{R}$ satisfies $P\mathbf{a}(\mathbf{w}) = \mathbf{a}(P\mathbf{w})$ for all permutation matrices $P \in \mathbb{R}^{M \times M}$ and define
 1383 $f_{\mathbf{a}(\mathbf{w}), \mathbf{w}} = \frac{1}{M} \sum_{i=1}^M \mathbf{a}_i(\mathbf{w})\sigma(w_i)$. Then we have

$$1384 \mathbb{E} [\mathbf{a}_i(\mathbf{w})\sigma(w_i)] = \mathbb{E} [\mathbf{a}_j(\mathbf{w})\sigma(w_j)] = \mathbb{E} [f_{\mathbf{a}(\mathbf{w}), \mathbf{w}}] \quad \text{for all } i, j \in \Lambda.$$

1385 *Proof.* The result follows directly from Lemma F.3 with $\Lambda = [M]$. Then $M/|\Lambda| = 1$ and the sum
 1386 over $j \notin \Lambda = [M]$ in (34) vanishes. □

1404 F.2 DERIVATIVES
1405

1406 This section provides several derivatives for the main results. We abbreviate

$$1407 g_\pi(w) := \bar{a}_\pi(w)\sigma(w) - f, \quad (35)$$

1408 with \bar{a}_π defined in (17).

1409 **Lemma F.5.** *Let $\pi \in \mathcal{M}_{+,1}$, g_π be given by (35) and f_w be given by (17). Then*

$$1410 \pi(g_\pi) = \mathbb{E}_{\mathbf{w} \sim \pi^M} [f_w - f].$$

1411 *Proof.* Since all $w_i, i \in [M]$ are distributed by π , we have

$$\begin{aligned} 1412 \pi(g_\pi) &= \mathbb{E}_{w_1} [\bar{a}(w_1)\sigma(w_1) - f] \\ 1413 &= \mathbb{E}_{w_1} [\mathbb{E}_{\mathbf{w}_{\setminus 1}} [\mathbf{a}_1(\mathbf{w})] \sigma(w_1) - f] \\ 1414 &= \mathbb{E}_{w_1} \mathbb{E}_{\mathbf{w}_{\setminus 1}} [\mathbf{a}_1(\mathbf{w})\sigma(w_1) - f] \\ 1415 &= \mathbb{E}_{\mathbf{w}} [\mathbf{a}_1(\mathbf{w})\sigma(w_1) - f] \\ 1416 &= \mathbb{E}_{\mathbf{w}} [f_{\mathbf{a}(\mathbf{w}), \mathbf{w}} - f] \\ 1417 &= \mathbb{E}_{\mathbf{w}} [f_w - f], \end{aligned}$$

1418 where in the first equality we use the definition (35) of g_π , in the second the definition (17) of $\bar{a}(w_1)$,
1419 in the fourth we join expectations, in the fifth we use Corollary F.4 and in the last the definition (17)
1420 of f_w .

1421 \square

1422 **Lemma F.6.** *Let $\pi \in \mathcal{M}_{+,1}$, $\mathbf{a}_1 \in L^2(\pi^M)$ be defined by (16) and \bar{a}_π be defined by (17). Then for
1423 any $\nu \ll \pi$ with bounded $d\nu/d\pi$*

$$1424 D_\pi[\bar{a}_\pi(w_1)]\nu = (M-1) \int \mathbb{E}_{\mathbf{w}_{\setminus 1,j}} [\mathbf{a}_1(\mathbf{w})] d\nu(w_j), \quad \text{for all } j \neq 1.$$

1425 *Proof.* With definition (17) of \bar{a}_π applied to input w_1 , we compute the derivative

$$\begin{aligned} 1426 D_\pi[\bar{a}_\pi(w_1)]\nu &= D_\pi \left[\int \mathbf{a}_1(\mathbf{w}) \prod_{i=2}^M d\pi(w_i) \right] \\ 1427 &= \sum_{j=2}^M \int \mathbf{a}_1(\mathbf{w}) \prod_{\substack{i=2 \\ i \neq j}}^M d\pi(w_i) d\nu(w_j) \\ 1428 &= \sum_{j=2}^M \int \mathbb{E}_{\mathbf{w}_{\setminus 1,j}} [\mathbf{a}_1(\mathbf{w})] d\nu(w_j). \end{aligned}$$

1429 Next, we show that the summands are independent of j . Indeed, by Lemma F.2, with $\Lambda = 2, \dots, M$,
1430 permutation p that swaps indices j and arbitrary $k \neq h = 1$ and $\pi_k = \nu$, we have

$$1431 \int \mathbb{E}_{\mathbf{w}_{\setminus 1,j}} [\mathbf{a}_1(\mathbf{w})] d\nu(w_j) = \int \mathbb{E}_{\mathbf{w}_{\setminus 1,k}} [\mathbf{a}_1(\mathbf{w})] d\nu(w_k).$$

1432 Therefore, in the sum above, we can choose an arbitrary index $j \neq 1$ and collapse the sum to the
1433 factor $M-1$, which concludes the proof.

1434 \square

1435 **Lemma F.7.** *Let $\pi \in \mathcal{M}_{+,1}$, $\mathbf{a}_1 \in L^2(\pi^M)$ be defined by (16), \bar{a}_π by (17) and $g_\pi = \bar{a}_\pi\sigma - f$. Then
1436 for any $\nu \ll \pi$ with bounded $d\nu/d\pi$ and $i \in [M]$*

$$1437 \pi(D_\pi[g_\pi])\nu = M\nu(\mathbb{E}_{w_i} [f_w]) - \nu(g_\pi) - \nu(1)f.$$

1458 *Proof.* Using the definition of g_π and Lemma F.6, for an arbitrary index $j \neq 1$, we have

$$\begin{aligned}
1459 \quad \pi(D_\pi[g_\pi]\nu) &= \mathbb{E}_{w_1} [D_\pi[\bar{a}_\pi(w_1)]\nu\sigma(w_1)] \\
1460 &= (M-1) \int \mathbb{E}_{w_1} \mathbb{E}_{\mathbf{w}_{\setminus 1, j}} [\mathbf{a}_1(\mathbf{w})\sigma(w_1)] d\nu(w_j) \\
1461 &= (M-1) \int \mathbb{E}_{\mathbf{w}_{\setminus j}} [\mathbf{a}_1(\mathbf{w})\sigma(w_1)] d\nu(w_j).
\end{aligned}$$

1465 Invoking Lemma F.3 with $\Lambda = [M] \setminus \{j\}$, this implies

$$\begin{aligned}
1466 \quad \pi(D_\pi[g_\pi]\nu) &= M \int \mathbb{E}_{\mathbf{w}_{\setminus j}} [f_{\mathbf{a}(\mathbf{w}), \mathbf{w}}] d\nu(w_j) - \int \mathbb{E}_{\mathbf{w}_{\setminus j}} [\mathbf{a}_j(\mathbf{w})\sigma(w_j)] d\nu(w_j) \\
1467 &= I - II.
\end{aligned}$$

1470 Using the abbreviation $f_{\mathbf{w}} = f_{\mathbf{a}(\mathbf{w}), \mathbf{w}}$, we simplify the first term to

$$1471 \quad I = M \int \mathbb{E}_{\mathbf{w}_{\setminus j}} [f_{\mathbf{w}}] d\nu(w_j).$$

1474 Since the $f_{\mathbf{w}}$ is invariant under permutations of \mathbf{w} , we can replace j by an arbitrary index $i \in [M]$, including 1, which was not allowed in the definition of j :

$$1475 \quad I = M \int \mathbb{E}_{\mathbf{w}_{\setminus i}} [f_{\mathbf{w}}] d\nu(w_i) = M\nu(\mathbb{E}_{\mathbf{w}_{\setminus i}} [f_{\mathbf{w}}]).$$

1478 To simplify II , we use symmetry from Lemma F.2, to rename $1 \leftrightarrow j$ so that

$$\begin{aligned}
1480 \quad II &= \int \mathbb{E}_{\mathbf{w}_{\setminus 1}} [\mathbf{a}_1(\mathbf{w})\sigma(w_1)] d\nu(w_1) \\
1481 &= \int \mathbb{E}_{\mathbf{w}_{\setminus 1}} [\mathbf{a}_1(\mathbf{w})] \sigma(w_1) d\nu(w_1) = \int \bar{a}_\pi(w_1)\sigma(w_1) d\nu(w_1),
\end{aligned}$$

1485 where the last step is the definition (17) of \bar{a}_π . Inserting f , we can express this in terms of g_π :

$$\begin{aligned}
1486 \quad II &= \int \bar{a}_\pi(w_1)\sigma(w_1) - f d\nu(w_1) + \int f d\nu(w_1), \\
1487 &= \int g_\pi(w_1) d\nu(w_1) + \int f d\nu(w_1), \\
1488 &= \nu(g_\pi) + \nu(1)f.
\end{aligned}$$

1492 Combining all terms

$$1493 \quad \pi(D_\pi[g_\pi]\nu) = M\nu(\mathbb{E}_{\mathbf{w}_{\setminus i}} [f_{\mathbf{w}}]) - \nu(g_\pi) - \nu(1)f$$

1495 concludes the proof. □

1498 The following corollary is a minor rearrangement of Lemma F.7, in the exact form used below.

1499 **Corollary F.8.** *Let $\pi \in \mathcal{M}_{+,1}$, $\mathbf{a}_1 \in L^2(\pi^M)$ be defined by (16), \bar{a}_π by (17) and $g_\pi = \bar{a}_\pi\sigma - f$. Then for any $\nu \ll \pi$ with bounded $d\nu/d\pi$ and $i \in [M]$*

$$1500 \quad \pi(D_\pi[g_\pi]\nu) = M\nu(\mathbb{E}_{\mathbf{w}_{\setminus i}} [f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus i}}]) + M\nu(\mathbb{E}_{\mathbf{w}_{\setminus i}} [f_{\mathbf{w}_{\setminus i}}]) - \nu(g_\pi) - \nu(1)f.$$

1504 *Proof.* Follows directly from Lemma F.7 with

$$1505 \quad \mathbb{E}_{\mathbf{w}_{\setminus i}} [f_{\mathbf{w}}] = \mathbb{E}_{\mathbf{w}_{\setminus i}} [f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus i}}] + \mathbb{E}_{\mathbf{w}_{\setminus i}} [f_{\mathbf{w}_{\setminus i}}]$$

1508 **Lemma F.9.** *Let $\pi \in \mathcal{M}_{+,1}$, $\mathbf{a}_1 \in L^2(\pi^M)$ be defined by (16), \bar{a}_π by (17), Δ_1 by (19) and $g_\pi = \bar{a}_\pi\sigma - f$. Then for any $\nu \ll \pi$ with bounded $d\nu/d\pi$*

$$1509 \quad \int D_\pi [\|g_\pi(w)\|^2] \nu d\pi(w) = \nu(\Delta_1) - 2\nu(\|g_\pi\|^2) - 2\langle \nu(g_\pi), f \rangle.$$

1512 We motivate the lemma with a simplified analogy, where:
1513

- 1514 • $\sigma(w) = 1$ and $f = 0$,
- 1515 • the measures $d\pi(w) = \rho(w) dw$ and $d\nu(w) = v(w) dw$ and $d\phi(w) = F(w) dw$ are given
1516 by densities
- 1517 • and the coefficients by $a_\rho = d\phi/d\pi = F/\rho$ are defined by Radon-Nikodym derivatives as
1518 for the exact outer weights (14).
1519

1520 Then the derivative $D_\rho[a_\rho\rho] = D_\rho[F] = 0$ vanishes because the inner function is independent of ρ .
1521 Thus by the product rule, we have
1522

$$\begin{aligned}
1523 \quad 0 &= 2 \int \langle a_\rho(w), D_\rho[a_\rho(w)\rho(w)]v \rangle \rho(w) dw \\
1524 &= 2 \int \langle a_\rho(w), D_\rho[a_\rho(w)]v \rangle \rho(w) dw + 2 \int \langle a_\rho(w), a_\rho(w) \rangle D_\rho[\rho(w)]v dw \\
1525 &= \int D_\rho[\|a_\rho(w)\|^2 v] \rho(w) dw + 2 \int \|a_\rho(w)\|^2 v(w) dw. \\
1526 &=: d - ii.
\end{aligned}$$

1531 We can now compare the term in this informal formula with terms in Lemma F.9, split into $D =$
1532 $I + II + III$. Then, d is analogous to D and ii to II . The remaining term I is a perturbation and
1533 III is a constant, which will be joined with Lagrange multipliers for normalization $\pi(1) = 1$.

1534 The above motivation also allows a direct comparison with the main results for exact outer weights
1535 a_π in (14). Then, $D_\rho[a_\rho\rho]$ corresponds to the derivative in term II in Lemma D.4. Here, we have
1536 used that $D[a_\rho\rho]$ is zero and likewise Lemma 3.2 shows that II is zero, which is a crucial step in
1537 the proof of equidistribution Lemma 3.2.
1538

1539 *Proof.* Let D denote the left hand side in the lemma's conclusion. Computing the derivative, plug-
1540 ging in the definition of g_π and renaming $w \rightarrow w_1$ for convenience below, we obtain
1541

$$1542 \quad D = 2 \int \langle g_\pi(w), D_\pi[g_\pi(w)]v \rangle d\pi(w) = 2 \int \langle g_\pi(w_1), D_\pi[\bar{a}_\pi(w_1)]\sigma(w_1)v \rangle d\pi(w_1).$$

1543 We compute the derivative of $\bar{a}_\pi(w)$ by Lemma F.6, with an arbitrary $j \neq 1$, and rearrange terms to
1544 conclude that
1545

$$1546 \quad D = 2(M-1) \iint \langle g_\pi(w_1), \mathbb{E}_{\mathbf{w}_{\setminus 1, j}}[\mathbf{a}_1(\mathbf{w})\sigma(w_1)] \rangle d\nu(w_j) d\pi(w_1).$$

1547 Next, we swap the variables w_1 and w_j and apply symmetry, Lemma F.2 with $h(w) =$
1548 $\langle g_\pi(w), \sigma(w) \rangle$, to obtain
1549

$$1550 \quad D = 2(M-1) \iint \langle g_\pi(w_j), \mathbb{E}_{\mathbf{w}_{\setminus 1, j}}[\mathbf{a}_j(\mathbf{w})\sigma(w_j)] \rangle d\nu(w_1) d\pi(w_j).$$

1551 Since D is independent of $j \neq 1$, so is the right hand side and we can average over all choices:
1552

$$1553 \quad D = 2 \sum_{j=2}^M \iint \langle g_\pi(w_j), \mathbb{E}_{\mathbf{w}_{\setminus 1, j}}[\mathbf{a}_j(\mathbf{w})\sigma(w_j)] \rangle d\nu(w_1) d\pi(w_j).$$

1554 Next, we add and subtract the missing $j = 1$ term in the sum:
1555

$$\begin{aligned}
1556 \quad D &= 2 \sum_{j=1}^M \iint \langle g_\pi(w_j), \mathbb{E}_{\mathbf{w}_{\setminus 1, j}}[\mathbf{a}_j(\mathbf{w})\sigma(w_j)] \rangle d\nu(w_1) d\pi(w_j). \\
1557 &\quad - 2 \int \langle g_\pi(w_1), \mathbb{E}_{\mathbf{w}_{\setminus 1}}[\mathbf{a}_1(\mathbf{w})\sigma(w_1)] \rangle d\nu(w_1) \\
1558 &= I - II.
\end{aligned}$$

1566 Note that II does not contain a π integral because w_1 is already integrated with respect to ν and
 1567 $\pi(1) = 1$. In order to simplify II , we first replace the expectation with the definition of \bar{a}_π
 1568

$$1569 \quad II = 2 \int \langle g_\pi(w_1), \bar{a}_\pi(\mathbf{w}_1) \sigma(w_1) \rangle d\nu(w_1).
 1570$$

1571 We add and subtract f so that we can replace the second component of the inner product with g_π :
 1572

$$1573 \quad II = 2 \int \langle g_\pi(w_1), \bar{a}_\pi(w_1) \sigma(w_1) - f \rangle d\nu(w_1) + 2 \int \langle g_\pi(w_1), f \rangle d\nu(w_1)
 1574 \quad = 2 \int \|g_\pi(w_1)\|^2 d\nu(w_1) + 2 \int \langle g_\pi(w_1), f \rangle d\nu(w_1)
 1575 \quad = 2\nu(\|g_\pi\|^2) + 2 \langle \nu(g_\pi), f \rangle.
 1576
 1577
 1578$$

1579 Collecting all terms $D = I - II$ and noting that $I = \nu(\Delta_1)$, given by (19), completes the proof.
 1580

□

1581 **Lemma F.10.** Let $\pi \in \mathcal{M}_{+,1}$, $\mathbf{a}_1 \in L^2(\pi^M)$ be defined by (16), \bar{a}_π by (17), Δ_1 by (19) and $G_\pi \in \mathcal{H}$
 1582 by (22). Then for any $\nu \ll \pi$ with bounded $d\nu/d\pi$
 1583

$$1584 \quad \nu(\Delta_1) = 2M \langle G_\pi, \nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}}]) \rangle.
 1585$$

1586 *Proof.* Recall that $\nu(\Delta_1)$ is defined by
 1587

$$1588 \quad I = 2 \sum_{j=1}^M \iint \langle g_\pi(w_j), \mathbb{E}_{\mathbf{w}_{\setminus 1,j}} [\mathbf{a}_j(\mathbf{w})] \sigma(w_j) \rangle d\nu(w_1) d\pi(w_j).
 1589
 1590$$

1591 Since $\mathbb{E}_{\mathbf{w}_{\setminus 1,j}} [\mathbf{a}_j(\mathbf{w})]$ is scalar, we can factor it out of the inner product to obtain
 1592

$$1593 \quad I = 2 \sum_{j=1}^M \iint \langle g_\pi(w_j), \sigma(w_j) \rangle \mathbb{E}_{\mathbf{w}_{\setminus 1,j}} [\mathbf{a}_j(\mathbf{w})] d\nu(w_1) d\pi(w_j).
 1594
 1595$$

1596 and invoke the given assumption of G_π :
 1597

$$1598 \quad I = 2 \sum_{j=1}^M \iint \langle G_\pi, \sigma(w_j) \rangle \mathbb{E}_{\mathbf{w}_{\setminus 1,j}} [\mathbf{a}_j(\mathbf{w})] d\nu(w_1) d\pi(w_j).
 1599
 1600$$

1601 For $j = 1$, the variable $w_j = w_1$ is bound to the ν integral and the outer π integral can be removed.
 1602 For $j \neq 1$, note that unlike $g_\pi(w_j)$, the new quantity G_π is independent of \mathbf{w} so that we can move
 1603 the outer π integral inside to join it with the expectation:
 1604

$$1605 \quad I = 2 \sum_{j=1}^M \int \langle G_\pi, \sigma(w_j) \rangle \mathbb{E}_{\mathbf{w}_{\setminus 1}} [\mathbf{a}_j(\mathbf{w})] d\nu(w_1).
 1606
 1607$$

1608 Next, we move the sum inside and abbreviate it by $f_{\mathbf{w}}$
 1609

$$1610 \quad I = 2M \int \left\langle G_\pi, \mathbb{E}_{\mathbf{w}_{\setminus 1}} \left[\frac{1}{M} \sum_{j=1}^M \mathbf{a}_j(\mathbf{w}) \sigma(w_j) \right] \right\rangle d\nu(w_1)
 1611 \quad = 2M \int \langle G_\pi, \mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}}] \rangle d\nu(w_1)
 1612 \quad = 2M \langle G_\pi, \nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}}]) \rangle,
 1613
 1614
 1615$$

1616 which proves the lemma.
 1617

□

1618 The following is a slight rearrangement of Lemma F.9 in the exact format that will be used below.
 1619

1620 **Corollary F.11.** Let $\pi \in \mathcal{M}_{+,1}$, $\mathbf{a}_1 \in L^2(\pi^M)$ be defined by (16), \bar{a}_π by (17) and $g_\pi = \bar{a}_\pi \sigma - f$.
 1621 Then for any $\nu \ll \pi$ with bounded $d\nu/d\pi$
 1622

$$1623 \int D_\pi \left[\|g_\pi(w)\|^2 \right] \nu d\pi(w) = \nu(\Delta_1) + \nu(\bar{\lambda}) - 2\nu(\|g_\pi\|^2) - 2 \langle \nu(g_\pi), f \rangle.$$

1624
 1625 with

$$1626 \bar{\lambda} = 0, \quad \text{if } \Delta_1 \text{ is defined by (19),}$$

$$1627 \bar{\lambda} = 2M \langle G_\pi, \mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}_{\setminus 1}}] \rangle, \quad \text{if } \Delta_1 \text{ is defined by (20).}$$

1628 *Proof.* In case Δ_1 is defined by (19), there is nothing to show. If it is defined by (20), by Lemma
 1629 F.10, we have

$$1630 \nu(\Delta_1) = 2M \langle G_\pi, \nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}}]) \rangle$$

$$1631 = 2M \langle G_\pi, \nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}]) \rangle + 2M \langle G_\pi, \nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}_{\setminus 1}}]) \rangle.$$

1632 On the right hand side, the first term matches the alternative definition (20) of Δ_1 . The second
 1633 simplifies to

$$1634 2M \langle G_\pi, \mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}_{\setminus 1}}] \rangle \nu(1) = \bar{\lambda} \nu(1),$$

1635 because the inner expectation $\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}_{\setminus 1}}]$ does not depend on the integration variable w_1 of the ν
 1636 integral. This concludes the proof.
 1637

1638 \square

1639 **Lemma F.12.** Let $\pi \in \mathcal{M}_{+,1}$, $\mathbf{a}_1 \in L^2(\pi^M)$ be defined by (16), \bar{a}_π by (17), Δ_1 by (19) or (20) and
 1640 Δ_2 by (21). Then there is a constant $\lambda \in \mathbb{R}$ so that for all $\nu \ll \pi$ with bounded $d\nu/d\pi$ the derivative
 1641 of the loss (18) is

$$1642 D_\pi [\ell(\pi)] \nu = -m\nu(\|\bar{a}_\pi \sigma\|^2) + m\lambda\nu(1) + m\nu(\Delta_1) + m\nu(\Delta_2).$$

1643 *Proof.* With $g_\pi = \bar{a}_\pi \sigma - f$ and $\mathbf{g}_\pi = \sum_{i=1}^m g_\pi(w_i)$ by the product rule, we have

$$1644 D_\pi [\ell(\pi)] \nu = D_\pi \left[\int \|g_\pi\|^2 d\pi^n(\mathbf{w}) \right] \nu$$

$$1645 = D_\pi \left[\int \|g_\eta\|^2 d\pi^n(\mathbf{w}) \right] \nu \Big|_{\eta=\pi} + \int D_\pi [\|g_\pi\|^2] \nu d\pi^n(\mathbf{w}).$$

1646 The first term, for fixed inner $\eta = \pi$, not differentiated, is given by Lemma D.2 and the second by
 1647 Lemma D.3. With normalization $\pi(1) = 1$, this yields

$$1648 D_\pi [\ell(\pi)] \nu = I + II + III + IV + V$$

1649 with

$$1650 I = m\nu(\|g_\pi\|^2),$$

$$1651 II = 2m(m-1) \langle \pi(g_\pi), \nu(g_\pi) \rangle,$$

$$1652 III = \lambda\nu(1),$$

$$1653 IV = m \int D_\pi [\|g_\pi(w)\|^2] \nu d\pi(w),$$

$$1654 V = 2m(m-1) \langle \pi(g_\pi), \pi(D_\pi[g_\pi]\nu) \rangle.$$

1655 By Corollary F.11, we have $IV = IV.1 + IV.2 + IV.3 + IV.4$ with

$$1656 IV.1 = m\nu(\Delta_1),$$

$$1657 IV.2 = m\nu(\bar{\lambda}),$$

$$1658 IV.3 = -2m\nu(\|g_\pi\|^2),$$

$$1659 IV.4 = -2m \langle \nu(g_\pi), f \rangle.$$

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

By Lemma F.5 and Corollary F.8, we have $V = V.1 + V.2 + V.3 + V.4$ with

$$V.1 = 2m(m-1)M \langle \mathbb{E}_{\mathbf{w}} [f_{\mathbf{w}} - f], \nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}]) \rangle,$$

$$V.2 = 2m(m-1)M \langle \pi(g_{\pi}), \nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}_{\setminus 1}}]) \rangle,$$

$$V.3 = -2m(m-1) \langle \pi(g_{\pi}), \nu(g_{\pi}) \rangle,$$

$$V.4 = -2m(m-1) \langle \pi(g_{\pi}), \nu(1)f \rangle.$$

The parts of III, IV.2, V.2, V.4 inside of $\nu(\cdot)$ are independent of w_1 and therefore can be joined into one constant λ_1 :

$$III + IV.2 + V.3 = \lambda_1 \nu(1).$$

The terms II and V.3 cancel:

$$II + V.3 = 0.$$

Next, we join the terms I, IV.3 and IV.4. Using $g_{\pi} = \bar{a}_{\pi}\sigma - f$ and defining $\lambda_2 := m\|f\|^2$, we have:

$$\begin{aligned} I + IV.3 + IV.4 &= -m\nu(\|g_{\pi}\|^2) - 2m \langle \nu(g_{\pi}), f \rangle \\ &= -m\nu(\|\bar{a}_{\pi}\sigma\|^2) + 2m\nu(\langle \bar{a}_{\pi}\sigma, f \rangle) - m\nu(\|f\|^2) - 2m \langle \nu(\bar{a}_{\pi}\sigma), f \rangle + 2m\nu(\|f\|^2) \\ &= -m\nu(\|\bar{a}_{\pi}\sigma\|^2) + m\|f\|^2 \nu(1) \\ &= -m\nu(\|\bar{a}_{\pi}\sigma\|^2) + \lambda_2 \nu(1). \end{aligned}$$

The only remaining terms are IV.1 $= m\nu(\Delta_1)$ and V.1 $= m\nu(\Delta_2)$, which are contained unchanged in the statement of the lemma. Defining $m\lambda = \lambda_1 + \lambda_2$ and combining all terms concludes the proof. \square

F.3 STABLE BARRON NORMS AND MAUREY SAMPLING

The following corollary is a variant of standard Maurey sampling arguments in Corollary 2.2, applied to stable Barron norms.

Corollary F.13. *Assume probability measure $\pi \in \mathcal{M}_{+,1}$ and $a_{\pi} \in L^2(\pi)$ satisfy the bounds (23) (24) for some π integrable function δ and $\epsilon \in \mathbb{R}$. Then for all $m \in \mathbb{N}$*

$$\mathbb{E}_{\mathbf{w} \sim \pi^m} \left\| \frac{1}{m} \sum_{i=1}^m a_{\pi}(w_i) \sigma(w_i) - f \right\|^2 \leq \frac{4}{m} [\|f\|_{B(\delta, \epsilon)}^2 + \pi(\delta)] + 2\epsilon^2.$$

Proof. Throughout this proof, let $g_{\pi}(w) = a_{\pi}(w)\sigma(w) - f$, for arbitrary $a_{\pi}: \mathcal{W} \rightarrow \mathbb{R}$, not only $a_{\pi} = \bar{a}_{\pi}$ as in our usual convention. Let w, w_i and $w'_i, i \in [m]$ be i.i.d. sampled from π . Since $g_{\pi}(w) - \pi(g_{\pi})$ has zero mean, by Lemma C.1 we have

$$\begin{aligned} \mathbb{E}_{\mathbf{w}} \left\| \frac{1}{m} \sum_{i=1}^m g_{\pi}(w_i) - \pi(g_{\pi}) \right\|^2 &= \frac{1}{m^2} \sum_{i=1}^m \mathbb{E}_{\mathbf{w}, \mathbf{w}'} \|g_{\pi}(w_i) - g_{\pi}(w'_i)\|^2 \\ &\leq \frac{2}{m^2} \sum_{i=1}^m \mathbb{E}_{\mathbf{w}} \|a_{\pi}(w_i)\sigma(w_i)\|^2 = \frac{2}{m} \mathbb{E}_{\mathbf{w}} \|a_{\pi}(w)\sigma(w)\|^2, \end{aligned}$$

which implies

$$\mathbb{E}_{\mathbf{w}} \left\| \frac{1}{m} \sum_{i=1}^m g_{\pi}(w_i) \right\|^2 \leq \frac{4}{m} \mathbb{E}_{\mathbf{w}} \|a_{\pi}(w)\sigma(w)\|^2 + 2\|\pi(g_{\pi})\|^2.$$

From (23) and (24), we have $|a_{\pi}(w)|^2 \|\sigma(w)\|^2 \leq \lambda + \delta(w)$ and $\|\pi(g_{\pi})\|^2 \leq \epsilon^2$ so that

$$\mathbb{E}_{\mathbf{w}} \left\| \frac{1}{m} \sum_{i=1}^m g_{\pi}(w_i) \right\|^2 \leq \frac{4}{m} [\lambda + \pi(\delta)] + 2\epsilon^2.$$

Since a_{π} and π are eligible candidates in the definition of the stable Barron norm, we have $\lambda \leq \|f\|_{B(\delta, \epsilon)}$, which concludes the proof. \square

1728 F.4 PROOF OF LEMMA 3.4: EQUIDISTRIBUTION

1729
1730 We heavily rely on the notational conventions in Appendix A.2. Let $\nu \in \mathcal{M}_+$ be a positive measure
1731 with bounded Radon-Nikodym derivative $d\nu/d\pi$. Then, the local minimizer π necessarily satisfies
1732 the first order optimality criteria (32)

$$1733 \quad D_\pi[\ell(\pi)]\nu + \nu(\lambda) \geq 0,$$

1734
1735 with equality for $\nu = \pi$. With Δ_1 and Δ_2 defined in (19) or (20) and (21), respectively, and
1736 $D_\pi[\ell(\pi)]\nu$ computed by Lemma F.12, we obtain upon a redefinition of the constant $\lambda \in \mathbb{R}$

$$1737 \quad -m\nu(\|\bar{a}_\pi\sigma\|^2) + m\lambda\nu(1) + m\nu(\Delta_1) + m\nu(\Delta_2) \geq 0.$$

1738
1739 Equivalently, we have

$$1740 \quad \int \left[-m \|\bar{a}_\pi(w)\sigma(w)\|^2 + m\lambda + m\Delta_1(w) + m\Delta_2(w) \right] \frac{d\nu}{d\pi} d\pi(w) \geq 0,$$

1741
1742 for all non-negative and bounded densities $d\nu/d\pi$, with equality if $d\nu/d\pi = 1$. Thus, it follows that

$$1743 \quad \|\bar{a}_\pi\sigma\|^2 = \lambda + \Delta_1 + \Delta_2$$

1744
1745 π -almost surely, as stated in the lemma.

1746 1747 F.5 PROOF OF THEOREM 3.5: APPROXIMATION ERROR

1748
1749 We heavily rely on the notational conventions in Appendix A.2. For arbitrary $\bar{\lambda}$ and $\delta := \bar{\lambda} + \Delta_1 +$
1750 Δ_2 , by Lemma 3.4, we have

$$1751 \quad |\bar{a}(w)|^2 \|\sigma(w)\|^2 = \lambda + \delta(w) \quad \pi - \text{a.s.}$$

1752
1753 for some $\lambda \in \mathbb{R}$ and thus

$$1754 \quad \left| \lambda - |\bar{a}(w)|^2 \|\sigma(w)\|^2 \right| \leq |\delta(w)|, \quad \pi - \text{a.s.}$$

1755
1756 This establishes the first perturbation bound (23) in the definition of stable Barron norms. By
1757 Lemma F.5 the second perturbation bound (24) is satisfied with $\|\mathbb{E}_{w \sim \pi} [g_\pi]\| = \|\pi(g_\pi)\| =$
1758 $\|\mathbb{E}_{w \sim \pi} [f_w - f]\| = \epsilon$. Hence, the result follows from Corollary F.13.

1759 1760 G TRAINING

1761 1762 G.1 PROOF OF LEMMA 4.1: PARTICLE APPROXIMATION OF WGF

1763
1764 The proof is standard (Chizat and Bach, 2018) and only included to make the exposition self con-
1765 tained. We first compute the gradient in the gradient flow of $w(t)$:

$$1766 \quad \begin{aligned} 1767 \quad \nabla_{w_i} \ell(\pi_w) &= \frac{1}{\mathfrak{m}} D[\ell(\pi_w)] \nabla_{w_i} \delta_{w_i} \\ 1768 \quad &= \frac{1}{\mathfrak{m}} \int \nabla_\pi \ell(\pi_w)(w) \nabla_{w_i} \delta(w - w_i) dw \\ 1769 \quad &= \frac{1}{\mathfrak{m}} \int \nabla_w \nabla_\pi \ell(\pi_w)(w) \delta(w - w_i) dw \\ 1770 \quad &= \frac{1}{\mathfrak{m}} \nabla_w \nabla_\pi \ell(\pi_w)(w_i), \end{aligned}$$

1771
1772 where in the third equality we have used integration by parts and the compact support of δ . Hence,
1773 the gradient flow of $w(t)$ is equivalent to

$$1774 \quad \dot{w}_i(t) = -\mathfrak{m} \nabla_{w_i} \ell(\pi_{w(t)}) = -\nabla_w \nabla_\pi \ell(\pi_{w(t)})(w_i(t)).$$

1775
1776 Next, we show that $\pi_{w(t)}$ satisfies Wasserstein gradient flow, i.e. that

$$1777 \quad \iint \varphi_t - \nabla_w \varphi(t, w) \nabla_w \nabla_\pi \ell(\pi_{w(t)})(w_i) d\pi_{w(t)} dt = 0$$

for all $\varphi(t, w)$ with compact support. Using that $\pi_{w(t)}$ consists of Dirac deltas, we can remove the inner integral to obtain

$$\frac{1}{m} \sum_{i=1}^m \int \varphi_t - \nabla_w \varphi(t, w_i(t)) \nabla_w \nabla_\pi \ell(\pi_{w(t)})(w_i(t)) dt = 0.$$

Plugging in the identity for $\dot{w}_i(t)$ above, this simplifies to

$$\frac{1}{m} \sum_{i=1}^m \int \varphi_t + \partial_w \varphi(t, w_i(t)) \dot{w}_i(t) dt = \frac{1}{m} \sum_{i=1}^m \int \frac{d}{dt} \varphi(t, w_i(t)) dt = 0,$$

where the last term is indeed zero because φ is compactly supported. This concludes the proof.

G.2 PROOF OF LEMMA 4.2

We first characterize stationary points of Wasserstein gradient flow.

Lemma G.1. *Let π be a stationary point of the Wasserstein gradient flow (25) with bounded support. Then $\nabla_w \nabla_\pi \ell(\pi) = 0$ π -almost surely.*

Proof. Let $\phi(t) \geq 0$ be compactly supported and normalized $\int \phi(t) dt = 1$ and $\psi(w)$ be compactly supported, and equal to one on the support of π . Then, plugging $\varphi(t, w) = \phi(t)\psi(w)\nabla_\pi \ell(\pi)(w)$ into the distributional definition of WGF and using that π is stationary and therefore independent of t , we obtain

$$\begin{aligned} 0 &= \int \dot{\phi}(t) dt \int \psi(w) \nabla_\pi \ell(\pi)(w) d\pi(w) \\ &\quad - \int \phi(t) dt \int \nabla_w [\psi(w) \nabla_\pi \ell(\pi)(w)] \cdot \nabla_w \nabla_\pi \ell(\pi)(w) d\pi(w) dt. \end{aligned}$$

The first summand vanishes because ϕ has compact support so that $\int \dot{\phi}(t) dt = 0$. For the second, we simplify $\psi(w) = 1$ on the support of π , and with $\int \phi(t) dt = 1$, we obtain

$$\int |\nabla_w \nabla_\pi \ell(\pi)(w)|^2 d\pi(w) dt = 0.$$

It follows that $\nabla_w \nabla_\pi \ell(\pi)(w) = 0$, π -almost surely. □

Proof of Lemma 4.2. Recall the notational conventions in Appendix A.2. By Lemma F.12, the directional derivatives of the loss are given by

$$D_\pi [\ell(\pi)] \nu = \int -m \|\bar{a}_\pi(w)\sigma(w)\|^2 + m\lambda + m\Delta_1(w) + m\Delta_2(w) d\nu(w)$$

so that by our conventions the gradient is

$$\nabla_\pi \ell(\pi)(w) = -m \|\bar{a}_\pi(w)\sigma(w)\|^2 + m\lambda + m\Delta_1(w) + m\Delta_2(w),$$

π -a.e. By Lemma 4.1 the discrete measure π_w is a stationary point of Wasserstein gradient flow. Then by Lemma G.1, we have $\nabla_w \nabla_\pi \ell(\pi) = 0$ π -almost surely and thus

$$\nabla_w \left[-\|\bar{a}_\pi(w)\sigma(w)\|^2 + \Delta_1(w) + \Delta_2(w) \right] = 0,$$

π -a.e. Since π_w is a sum of Dirac deltas at location w_i , this implies the lemma. □

1836 G.3 PROOF OF (29)

1837
1838 In this section, we compute (29). To this end, we abbreviate

$$1839 \tilde{\ell}(\pi) = \frac{1}{N_\ell} \sum_{j=1}^{N_\ell} L(\mathbf{v}^j) := \frac{1}{N_\ell} \sum_{j=1}^{N_\ell} \left\| \frac{1}{m} \sum_{i=1}^m \tilde{a}_\pi(v_i^j) \sigma(v_i^j) - f \right\|^2, \quad \mathbf{v}^j \sim \pi^m$$

1842 and recall that gradient descent is given by

$$1843 \mathbf{w}^n = \mathbf{w}^{n-1} - \gamma \nabla_{\mathbf{w}} \tilde{\ell}(\pi_{\mathbf{w}^{n-1}}),$$

1844 or plugging in all intermediate steps as well as $\tilde{\ell}(\pi)$ by

$$1845 \mathbf{w}^n = \mathbf{w}^0 - \gamma \sum_{k=0}^{n-1} \nabla_{\mathbf{w}} \tilde{\ell}(\pi_{\mathbf{w}^k}) = \mathbf{w}^0 - \frac{\gamma}{N_\ell} \sum_{k=0}^{n-1} \sum_{j=1}^{N_\ell} \nabla_{\mathbf{w}} L(\mathbf{v}^{jk}), \quad \mathbf{v}^{jk} \sim \pi_{\mathbf{w}^k}^m.$$

1846
1847 Instead of sampling $\mathbf{v}^{jk} \sim \pi_{\mathbf{w}^k}^m$ for all j , repeatedly from the same distribution, we always sample
1848 from the latest one available. To this end, we denote by l the lexicographic ordering the index pairs
1849 jk and by l the upper bound $k = n$ and $j = N_\ell$ so that

$$1850 \mathbf{w}^l = \mathbf{w}^0 - \frac{\gamma}{N_\ell} \sum_{l=0}^{l-1} \nabla_{\mathbf{w}} L(\mathbf{v}^l) \quad \mathbf{v}^l \sim \pi_{\mathbf{w}^l}^m.$$

1851
1852 Rewriting as a recursive formula

$$1853 \mathbf{w}^l = \mathbf{w}^{l-1} - \frac{\gamma}{N_\ell} \nabla_{\mathbf{w}} L(\mathbf{v}^{l-1}) \quad \mathbf{v}^{l-1} \sim \pi_{\mathbf{w}^{l-1}}^m.$$

1854
1855 this yields the gradient descent iteration in (29).

1860 G.4 PROOF OF (30)

1861
1862 In this section, we show (30). To this end, we abbreviate

$$1863 L(\mathbf{a}, \mathbf{w}) := \left\| \frac{1}{M} \sum_{i=1}^M a_i^n \sigma(w_i) - f \right\|^2$$

1864
1865 and compute the best approximation coefficients $\mathbf{a}(\mathbf{w})$ in (15) by gradient descent

$$1866 \mathbf{a}^{n+1}(\mathbf{w}) = \mathbf{a}^n(\mathbf{w}) - \lambda \nabla_{\mathbf{a}} L(\mathbf{a}^n, \mathbf{w}).$$

1867
1868 Summing the first component \mathbf{a}_1 over samples $\mathbf{w} = (v, \mathbf{v}_{\setminus 1}^j)$ with $\mathbf{v}^j \sim \pi^M$, we obtain

$$1869 \frac{1}{N_a} \sum_{j=1}^{N_a} \mathbf{a}_1^{n+1}(v, \mathbf{v}_{\setminus 1}^j) = \frac{1}{N_a} \sum_{j=1}^{N_a} \mathbf{a}_1^n(v, \mathbf{v}_{\setminus 1}^j) - \lambda \nabla_{\mathbf{a}_1} \frac{1}{N_a} \sum_{j=1}^{N_a} L(\mathbf{a}^n, v, \mathbf{v}_{\setminus 1}^j)$$

1870
1871 Abbreviating the left hand side by $\mathbf{a}_\pi^n(v)$, it converges to the mean

$$1872 \tilde{\mathbf{a}}_\pi^n(v) \rightarrow \tilde{\mathbf{a}}_\pi(v) = \frac{1}{N_a} \sum_{j=1}^{N_a} \mathbf{a}_1^\infty(v, \mathbf{v}_{\setminus 1}^j),$$

1873
1874 defined in (27). To ease the computation, we replace the argument \mathbf{a}_i^n of L with $\tilde{\mathbf{a}}_\pi^n(v_i^j)$ to obtain

$$1875 \tilde{\mathbf{a}}_\pi^n(v) = \tilde{\mathbf{a}}_\pi^{n-1}(v) - \lambda \nabla_{\mathbf{a}} \frac{1}{N_a} \sum_{j=1}^{N_a} L(\tilde{\mathbf{a}}_\pi^{n-1}(v, \mathbf{v}_{\setminus 1}^j), v, \mathbf{v}_{\setminus 1}^j).$$

1876
1877 Note that by our notational conventions the application of $\tilde{\mathbf{a}}_\pi^n$ to the vector $(v, \mathbf{v}_{\setminus 1}^j)$ is component
1878 wise and that v belongs to the discrete set w_1, \dots, w_m of all particles in the particle approximation
1879 (26). Analogous to Appendix G.3, we can always use the newest possible $\tilde{\mathbf{a}}_\pi^n$ in the iteration and
1880 relabel n to obtain

$$1881 \tilde{\mathbf{a}}_\pi^n(v) = \tilde{\mathbf{a}}_\pi^{n-1}(v) - \frac{\lambda}{N_a} \nabla_{\mathbf{a}} L(\tilde{\mathbf{a}}_\pi^{n-1}(v, \mathbf{v}_{\setminus 1}), v, \mathbf{v}_{\setminus 1}), \quad \mathbf{v} \sim \pi^M,$$

1882
1883 which provides (30).

1890 H PROOFS: PERTURBATIONS

1891 H.1 PRELIMINARIES

1892 **Lemma H.1.** *Let \mathbf{a} and $f_{\mathbf{w}}$ be given by (15) and $\mathcal{H}(\mathbf{w})$ by (31). Then, for all $\mathbf{w} \in \mathcal{W}^M$ and*
 1893 *$h \in \mathcal{H}(\mathbf{w}_{\setminus 1})$, we have*

$$1894 \langle f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}, h \rangle = 0.$$

1895 *Proof.* Since $f_{\mathbf{w}}$ and $f_{\mathbf{w}_{\setminus 1}}$ are best approximations, we have

$$1896 \begin{aligned} 1897 \langle f_{\mathbf{w}}, h \rangle &= \langle f, h \rangle, & \text{for all } h \in \mathcal{H}(\mathbf{w}), \\ 1898 \langle f_{\mathbf{w}_{\setminus 1}}, h \rangle &= \langle f, h \rangle, & \text{for all } h \in \mathcal{H}(\mathbf{w}_{\setminus 1}) \end{aligned}$$

1899 so that in particular $\langle f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}, h \rangle = 0$ for all $h \in \mathcal{H}(\mathbf{w}_{\setminus 1})$.

1900 \square

1901 H.2 PROOF OF LEMMA B.1

1902 From Lemma H.1, we have

$$1903 \langle f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}, h \rangle = 0 \quad \text{for all } h \in \mathcal{H}(\mathbf{w}_{\setminus 1})$$

1904 so that $f_{\mathbf{w}_{\setminus 1}}$ is not only a best approximation of f but also of $f_{\mathbf{w}}$, i.e.

$$1905 f_{\mathbf{w}_{\setminus 1}} = \arg \min_{\varphi \in \mathcal{H}(\mathbf{w}_{\setminus 1})} \|\varphi - f_{\mathbf{w}}\|^2.$$

1906 Since $M^{-1}\mathbf{a}_1(\mathbf{w})$ is defined as a coefficient of the best approximation, we have $f_{\mathbf{w}} -$
 1907 $M^{-1}\mathbf{a}_1(\mathbf{w})\sigma(w_1) \in \mathcal{H}(\mathbf{w}_{\setminus 1})$ and thus

$$1908 \|f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}\| \leq \|f_{\mathbf{w}} - [f_{\mathbf{w}} - M^{-1}\mathbf{a}_1(\mathbf{w})\sigma(w_1)]\| = M^{-1} \|\mathbf{a}_1(\mathbf{w})\sigma(w_1)\|.$$

1909 Applying the expectation on both sides and factoring out terms independent of $\mathbf{w}_{\setminus 1}$, proves the
 1910 lemma.

1911 H.3 PROOF OF COROLLARY B.2

1912 The result is a direct consequence of Corollary F.13. To this end, we first bound the perturbation
 1913 terms Δ_1 and Δ_2 given in (20) and (21), respectively. With the given assumptions, we have

$$1914 \mathbb{E}_{\mathbf{w}_{\setminus 1}} [\|\mathbf{a}_1(\mathbf{w})\|] \leq c \mathbb{E}_{\mathbf{w}_{\setminus 1}} [|\mathbf{a}_1(\mathbf{w})|] = |\bar{a}_{\pi}(w_1)|$$

1915 and therefore, with Lemma B.1

$$1916 \begin{aligned} 1917 \|\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}]\| &\leq \mathbb{E}_{\mathbf{w}_{\setminus 1}} \|f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}\| \leq M^{-1} \mathbb{E}_{\mathbf{w}_{\setminus 1}} [\|\mathbf{a}_1(\mathbf{w})\|] \|\sigma(w_1)\| \\ 1918 &\leq cM^{-1} |\bar{a}_{\pi}(w_1)| \|\sigma(w_1)\|. \end{aligned}$$

1919 Together with

$$1920 \langle h, \nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}]) \rangle = 0$$

1921 for arbitrary $h \in \mathcal{H}(\mathbf{w}_{\setminus 1})$ by Lemma H.1, we bound the first perturbation term (20) by

$$1922 \begin{aligned} 1923 \nu(\Delta_1) &= 2M \langle G_{\pi}, \nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}]) \rangle \\ 1924 &= 2M \langle G_{\pi} - h, \nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}]) \rangle \\ 1925 &\leq 2c \|G_{\pi} - h\| \nu(|\bar{a}_{\pi}(w_1)| \|\sigma(w_1)\|). \end{aligned}$$

1926 Similarly, the second perturbation term (21) is bounded by

$$1927 \begin{aligned} 1928 \nu(\Delta_2) &= 2(m-1)M \langle \mathbb{E}_{\mathbf{w}} [f_{\mathbf{w}} - f], \nu(\mathbb{E}_{\mathbf{w}_{\setminus 1}} [f_{\mathbf{w}} - f_{\mathbf{w}_{\setminus 1}}]) \rangle, \\ 1929 &\leq 2c(m-1) \|\mathbb{E}_{\mathbf{w}} [f_{\mathbf{w}} - f]\| \nu(|\bar{a}_{\pi}(w_1)| \|\sigma(w_1)\|). \end{aligned}$$

1944 Since $\nu \ll \pi$ is arbitrary, with $\delta := \Delta_1 + \Delta_2$ we conclude that
 1945

$$\begin{aligned}
 1946 \quad |\delta(w)| &\leq 2c \|G_\pi - h\| |\bar{a}_\pi(w)| \|\sigma(w)\| \\
 1947 \quad &\quad + 2c(m-1) \|\mathbb{E}_{\mathbf{w}} [f_{\mathbf{w}} - f]\| |\bar{a}_\pi(w)| \|\sigma(w)\|. \\
 1948 \quad &\leq 2c \left[\|G_\pi - h\| + (m-1) \|\mathbb{E}_{\mathbf{w}} [f_{\mathbf{w}} - f]\| \right] \|\bar{a}_\pi \sigma\|_{L^\infty(\pi; \mathcal{H})} \\
 1949 \quad &=: \Delta. \\
 1950
 \end{aligned}$$

1951 Next, we establish the bounds (23) and (24) for the definition of stable Barron norms. By Lemma
 1952 3.4, we have
 1953

$$1954 \quad |\bar{a}_\pi(w)|^2 \|\sigma(w)\|^2 = \lambda + \delta(w) \qquad \pi - \text{a.s.}$$

1955 for some $\lambda \in \mathbb{R}$. Combining the last two results yields
 1956

$$1957 \quad \left| \lambda - |\bar{a}_\pi(w)|^2 \|\sigma(w)\|^2 \right| \leq \Delta(w),$$

1958 which shows (23). By Lemma F.5 the second perturbation bound (24) is satisfied with
 1959 $\|\mathbb{E}_{\mathbf{w} \sim \pi} [g_\pi]\| = \|\pi(g_\pi)\| = \|\mathbb{E}_{\mathbf{w} \sim \pi} [f_{\mathbf{w}} - f]\| = \epsilon$. Hence, the result follows from Corollary
 1960 F.13.
 1961
 1962
 1963
 1964
 1965
 1966
 1967
 1968
 1969
 1970
 1971
 1972
 1973
 1974
 1975
 1976
 1977
 1978
 1979
 1980
 1981
 1982
 1983
 1984
 1985
 1986
 1987
 1988
 1989
 1990
 1991
 1992
 1993
 1994
 1995
 1996
 1997