
The Unexpected Deterministic and Universal Behavior of Large Softmax Classifiers

Mohamed El Amine Seddik^{1,2}, Cosme Louart^{1,3}, Romain Couillet^{2,3}, Mohamed Tamaazousti¹

¹CEA List, ²Centralesupélec, ³GIPSA Lab Grenoble-Alpes University

¹Palaiseau, France, ²Gif-sur-Yvette, France, ³Grenoble, France

firstname.lastname@{cea,centralesupelec}.fr

Abstract

In this paper, we provide a large dimensional analysis of the Softmax classifier, extensively used in modern neural networks. We discover and prove that, when the classifier is trained on data satisfying reasonable concentration assumptions, its weights become deterministic and solely depend on the statistical means and covariances of the data. As a striking consequence, despite the implicit and non-linear nature of the underlying optimization problem, the performance of the Softmax classifier is the same as if performed on a mere Gaussian mixture model, thereby disrupting the intuition that non-linearities inherently extract advanced statistical features from the data. Our findings are in particular theoretically sustained as well as numerically confirmed on CNN representations of images produced by GANs.

1 Introduction

The intricate nature of deep neural network training leaves little insight on the specific information encoded into the inter-layer connectivity weights of a fully trained network, thereby so far not allowing for particularly useful interpretation and control of their performances [YKYR18].

At the very source of these difficulties are the multiple non-linearities and the implicit optimization scheme involved in the network design: the activation functions in the intermediate layers as well as the soft or hard final decision layer [LWL⁺17]. For lack of a tractable comprehensive analysis, literature studies have mostly focused on individual components which, when isolated, become tractable. For instance, the effect of non-linearities in a single-hidden layer network was analysed in [PW17, LLC⁺18], the learning dynamics in elementary network designs in [SMG13, dCPS⁺18] and the overall understanding of the geometry of the loss surface in a largely approximated version of a deep neural net in [PB17, CHM⁺15].

These works are however restricted to the analysis either of the intermediate layers of practical neural nets, or oversimplify the network to an extent that makes the results rather impractical. The present article instead focuses on the training of the weights of the last decision layer, by specifically studying the widely used Softmax component in neural networks classifiers. The Softmax classifier has the property, which we will see to be of importance here, to be optimal for Gaussian mixture inputs with equal covariance [YW19]. Specifically, assuming the feature representations of the data fixed at the penultimate layer of the network, and modelling these features as *concentrated random vectors* [Led05] (which is a natural assumption as concentrated random vectors enjoy the property to be stable through Lipschitz maps, and thus through the action of intermediate neural network layers [SLTC20]), the article studies the statistical behavior of this last layer once trained (see Figure 1).

Our analysis leverages recent advances in random matrix theory by supposing the realistic setting where the number of data samples n (here their representations at the penultimate layer) and their

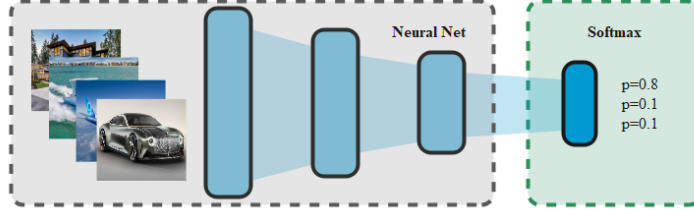


Figure 1: A neural network with fixed representations in the penultimate layer (in gray) and a trained Softmax output (in green).

dimension p (the size of this representation, i.e., the number of neurons in the one-before-last layer) are both large and comparable.

From a technical standpoint, as the Softmax classifier training corresponds to a (possibly non-convex) optimization problem, our analysis of the Softmax weights is performed by first expressing the optimization problem as a contracting fixed point equation, and then showing that the assumed *concentration properties* of the data features naturally transfer to the solution of the fixed-point equation, and thus to the Softmax weights. This has the major consequence that, as $n, p \rightarrow \infty$, the Softmax weights tend to have a deterministic behavior which we express explicitly as a function of the data statistics and the Softmax parameters.

Our most fundamental findings may be summarized as follows:

1. the above deterministic behavior exhibits a surprising *universality* of the Softmax classifier, in the sense that the large dimensional statistics of the weights solely depend on the statistical means and covariances of the input data features;
2. this suggests in turn that, quite counter-intuitively, at least as far as the last Softmax classification layer is concerned, no further discriminative feature of the data is extracted and, possibly most outstandingly, *the Softmax layer treats the input data as if they were Gaussian random vectors*; this, in passing, supports the Gaussianity assumption on the data representations commonly considered in the literature [HRU⁺17, PRU⁺18];
3. combined to the aforementioned optimality of the Softmax classifier on Gaussian mixture models with strongly discriminative class-wise means, this compellingly supports an overall classification optimality of the Softmax classifier on large dimensional representations of real data. A similar behavior was already pointed out, yet not well understood, by the authors in [MVPC13, GCM18];
4. Our findings are supported both theoretically and practically by considering the input data features as CNN-representations of images generated by the BigGAN model [BDS18].

In the remainder of the article, we introduce more precisely the present model of Softmax classification training and recall some concentration of measure tools necessary for best understanding (Section 2), before stating our main theoretical results (Section 3) along with supporting experiments and further concluding remarks (Section 4). The detailed derivations of our results are deferred to Section 5 and the Supplementary Material.

Notation: For $m \in \mathbb{N}$, $[m] \equiv \{1, \dots, m\}$. Vectors are denoted by boldface lowercase and matrices by boldface uppercase letters. The set of matrices of size $p \times n$ is denoted $\mathcal{M}_{p,n}$, the set of squared matrices and diagonal matrices of size n respectively \mathcal{M}_n and \mathcal{D}_n . $\|\cdot\|$ is the Euclidean (resp., spectral) norm for vectors (resp., matrices with $\|\cdot\| : \mathbf{M} \mapsto \sup_{\|u\| \leq 1} \|\mathbf{M}u\|$); $\|\cdot\|_F$ stands for the Frobenius norm $\|\cdot\|_F : \mathbf{M} \mapsto \sqrt{\text{Tr}(\mathbf{M}\mathbf{M}^T)}$ and $\|\cdot\|_*$ stands for the nuclear norm $\|\cdot\|_* : \mathbf{M} \mapsto \text{Tr}(\sqrt{\mathbf{M}\mathbf{M}^T})$ (which is the dual norm of the spectral norm for the scalar product $\mathbf{A}, \mathbf{B} \mapsto \text{Tr}(\mathbf{A}^T \mathbf{B})$). \otimes stands for the Kronecker product.

2 Model setting

2.1 The Softmax classifier

Let $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$ be a set of n labeled data associated to one of k classes $\mathcal{C}_1, \dots, \mathcal{C}_k$, where $\mathbf{x}_i \in \mathbb{R}^p$ and $\mathbf{y}_i \in \mathbb{R}^k$ is one-hot encoded vectors such that $y_{i\ell} = 1$ if $\mathbf{x}_i \in \mathcal{C}_\ell$. The \mathbf{x}_i 's are assumed

to be the input of an ℓ_2 -regularized Softmax classifier with regularization parameters $(\lambda_\ell)_{\ell \in [k]} \in \mathbb{R}^+$, which aims to determine the class-wise weight vectors $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^p$ minimizing the loss¹, for some real-valued function $\phi : \mathbb{R} \rightarrow \mathbb{R}$:

$$\mathcal{L}(\mathbf{w}_1, \dots, \mathbf{w}_k) = -\frac{1}{n} \sum_{i=1}^n \sum_{\ell=1}^k y_{i\ell} \log p_{i\ell} + \frac{1}{2} \sum_{\ell=1}^k \lambda_\ell \|\mathbf{w}_\ell\|^2 \quad \text{with} \quad p_{i\ell} = \frac{\phi(\mathbf{w}_\ell^\top \mathbf{x}_i)}{\sum_{j=1}^k \phi(\mathbf{w}_j^\top \mathbf{x}_i)}$$

In particular, the classical Softmax classifier corresponds to the case where $\phi(t) = e^t$ [GP17]. Cancelling the loss function gradient with respect to each weight vector \mathbf{w}_ℓ yields

$$\lambda_\ell \mathbf{w}_\ell = -\frac{1}{n} \sum_{i=1}^n \left(y_{i\ell} \psi(\mathbf{w}_\ell^\top \mathbf{x}_i) - \frac{\phi(\mathbf{w}_\ell^\top \mathbf{x}_i)}{\sum_{j=1}^k \phi(\mathbf{w}_j^\top \mathbf{x}_i)} \sum_{j=1}^k y_{ij} \psi(\mathbf{w}_j^\top \mathbf{x}_i) \right) \mathbf{x}_i, \quad \ell \in [k], \quad (1)$$

where $\psi \equiv \phi'/\phi$. Under appropriate statistical assumptions on the data matrix $\mathbf{X} \equiv [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathcal{M}_{p,n}$ and on ψ , and assuming p, n large, we subsequently show that the vector $\mathbf{W} \equiv [\mathbf{w}_1^\top, \dots, \mathbf{w}_k^\top]^\top \in \mathbb{R}^{pk}$ has a well defined behavior, which in turn allows us to accurately predict the performances of the Softmax classifier.

2.2 Mixture of concentrated random vectors

We first characterize the data classes: if $y_{i\ell} = 1$, then $\mathbf{x}_i \in \mathbb{R}^p$ is a random vector with

$$\mathbb{E}[\mathbf{x}_i] \equiv \boldsymbol{\mu}_\ell, \quad \mathbb{E}_{\mathbf{x}_i}[\mathbf{x}_i \mathbf{x}_i^\top] - \boldsymbol{\mu}_\ell \boldsymbol{\mu}_\ell^\top \equiv \boldsymbol{\Sigma}_\ell.$$

The vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ are further assumed to be independent and are such that $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathcal{M}_{p,n}$ satisfies a *concentration* property. To properly state this assumption (Assumption 1), which is central to our analysis, we first define the notion of random vector concentration.

Definition 1 (Concentrated vector). *Given a set of indices \mathbb{S} , a sequence of normed vector spaces $(E_s, \|\cdot\|_s)_{s \in \mathbb{S}}$, a sequence of random vectors $\mathbf{Z}_s \in E_s$, a sequence of positive numbers σ_s , we say that \mathbf{Z}_s is q -exponentially concentrated with an observable diameter of order $O(\sigma_s)$ if there exists two constants $C, c > 0$ such that for all sequence of 1-Lipschitz mappings $f_s : E_s \rightarrow \mathbb{R}$:*

$$\forall s \in \mathbb{S}, \forall t > 0 : \mathbb{P}(|f_s(\mathbf{Z}_s) - \mathbb{E}[f_s(\mathbf{Z}_s)]| \geq t) \leq C e^{-c(t/\sigma_s)^q}. \quad (2)$$

We note then $\mathbf{Z}_s \propto \mathcal{E}_q(\sigma_s)$, or simply $\mathbf{Z} \propto \mathcal{E}_q(\sigma)$; when $\sigma_s = O(1)$, we write $\mathbf{Z}_s \propto \mathcal{E}_q$.

The prototypical example of a concentrated random vector is the Gaussian random vector $\mathbf{Z}_s \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ for which $\mathbf{Z}_s \propto \mathcal{E}_2$ ([Led05]). But the richness of concentrated random vectors lies in their fundamental stability property through Lipschitz operations, which naturally generates wide families of concentrated random vectors.

Remark 2.1 (Stability through Lipschitz transformations). *It is easily deduced from Definition 1 that given a sequence of positive numbers $L_s > 0$ and a sequence of L_s -Lipschitz transformations $\phi_s : (E_s, \|\cdot\|_s) \rightarrow (F_s, \|\cdot\|'_s)$, if $\mathbf{Z}_s \propto \mathcal{E}_q(\sigma_s)$, then $\phi_s(\mathbf{Z}_s) \propto \mathcal{E}_q(L_s \sigma_s)$ (indeed, for all $f_s : F_s \rightarrow \mathbb{R}$, 1-Lipschitz, $\frac{1}{L_s} f_s \circ \phi_s$ is 1-Lipschitz, and one can employ inequality 2 to $\frac{t}{L_s}$).*

In particular, the concentration of Gaussian vectors combined with the stability through Lipschitz transformations as per Remark 2.1 provides a wide range of random vectors, among which random vectors with possibly quite complex dependence structures between their entries. A remarkable example of such random vectors are random vectors produced by generative adversarial networks (GANs) [GP14]: GAN random vectors notably satisfy that their outputs has the same concentration² as their inputs [SLTC20]; in particular, for Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I}_m)$ inputs (as traditionally assumed) whose observable diameter does not depend on the dimension m , the observable diameter of the GAN outputs does not increase with the data dimension. Besides, further operations through

¹Biases are not introduced in the present formulation as their effect is known to be negligible in practice [KXR⁺19] and would only decrease the readability and accessibility of our results.

²When the GAN model has a controlled Lipschitz constant, which is practically ensured by Spectral Normalization as in the BigGAN model [BDS18].

neural network layers with controlled Lipschitz norms (as is again traditionally done) on concentrated random vectors also maintain the concentration and observable diameter.

As a consequence of the above remark, making the approximation that GAN-generated data are alike real data, we may assume that GAN data fed into the first layer of a deep neural network are output in the one-before-last layer as a concentrated random vector with observable diameter independent of its dimension. This is summarized into our present Softmax input data assumption:

Assumption 1 (Concentrated data). *Letting $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$, $\mathbf{X} \propto \mathcal{E}_2$.*

In the terms of Definition 1, Assumption 1 holds here for $s = (p, n)$ with $\mathbb{S} = \mathbb{N}^2$. In order to be able to transfer the concentration of \mathbf{X} to the Softmax weights $\mathbf{w}_1, \dots, \mathbf{w}_k$, a further condition is needed: the number of data n must scale with the data dimension p , i.e., $\mathbb{S} = \{(p, n) \in \mathbb{N}^2, \kappa p \leq n \leq Kp\}$ for some $K > \kappa > 0$.³ This is summarized by the request:

Assumption 2 (Growth rate). *$n = O(p)$ and $p = O(n)$.*

Concentrated vectors satisfy a host of interesting properties (the reader being referred to [Led05] for a detailed account and to [LLC⁺18] for their application to random matrix asymptotics, closer to the present work). We merely stress here one of these properties, of central importance to the present work, and which fundamentally justifies the appearance of Gaussian-like behaviors in large neural networks, even when the neural network input is far from Gaussian [KGC18, NBA⁺18].

Theorem 2.2 (CLT for concentrated vectors [Kla07, FGP07]). *Let $\mathbf{X} \in \mathbb{R}^p$ be a random vector with $\mathbb{E}[\mathbf{X}] = 0$ and $\mathbb{E}[\mathbf{X}\mathbf{X}^\top] = I_p$, and σ be the uniform measure on the sphere $\mathcal{S}^{p-1} \subset \mathbb{R}^p$ of radius 1. Then, if $\mathbf{X} \propto \mathcal{E}_2$, there exists two constants $C, c > 0$ and a set $\Theta \subset \mathcal{S}^{p-1}$ such that $\sigma(\Theta) \geq 1 - \sqrt{p}Ce^{-c\sqrt{p}}$ and $\forall \boldsymbol{\theta} \in \Theta$:*

$$\sup_{t \in \mathbb{R}} |\mathbb{P}(\boldsymbol{\theta}^\top \mathbf{X} \geq t) - G(t)| \leq p^{-1/4}$$

for G the cumulative distribution function of an $\mathcal{N}(0, 1)$ random variable.

3 Main results

3.1 Behavior of the Softmax classifier weights and performance estimation

This section characterizes the statistical behavior of the Softmax classifier weights $\mathbf{W} = [\mathbf{w}_1^\top, \dots, \mathbf{w}_k^\top]^\top \in \mathbb{R}^{pk}$ under Assumptions 1–2 and, as a result, accesses the asymptotic performances of the classifier.

To this end, our approach is to first write the implicit defining equation 1 of \mathbf{W} under the formal form $\mathbf{W} = \Psi(\mathbf{W})$, for $\Psi : \mathbb{R}^{pk} \rightarrow \mathbb{R}^{pk}$ to characterize, and then to *transfer* the concentration of \mathbf{X} (Assumption 1) to a concentration of \mathbf{W} .

For the concentration of \mathbf{X} to propagate into \mathbf{W} defined through the formal form $\mathbf{W} = \Psi(\mathbf{W})$, Ψ is required to have contracting properties, which in turn will enforce structural conditions on the operator ϕ and on the regularizers $(\lambda_\ell)_{\ell \in [k]}$. Specifically Ψ is requested to be $(1 - \varepsilon)$ -Lipschitz (for some $\varepsilon > 0$) so to ensure, thanks to the Banach fixed point theorem, the existence and uniqueness of $\mathbf{W} \in \mathbb{R}^{pk}$. However, being a *random map* depending on \mathbf{X} , Ψ is only contracting under the (asymptotically highly probable) event \mathcal{A}_X (see Section 2 of the Supplementary Material) that the norm of \mathbf{X} is not too large. With these informal steps in mind, we are in position to state our main results.

Theorem 3.1 (Concentration of \mathbf{W}). *Under Assumptions 1 and 2 and additional assumptions on ϕ and $(\lambda_\ell)_{\ell \in [k]}$ provided in Section 2 of the Supplementary Material (Assumptions 3, 4 and 5), there exist two constants $C, c > 0$ and an event \mathcal{A}_X with $\mathbb{P}(\mathcal{A}_X) > 1 - Ce^{-cn}$ such that⁴*

$$(\mathbf{W} \mid \mathcal{A}_X) \propto \mathcal{E}_2 \left(\sqrt{\log n/n} \right).$$

³Formally, in the present setting, it is sufficient that $p \leq \frac{1}{\kappa}n$. However, to obtain simpler expressions, it is convenient to assume, in addition, that $n \leq Kp$.

⁴Formally, the random vector \mathbf{W} is a measurable mapping $\Omega \rightarrow \mathbb{R}^{pk}$, where $(\Omega, \mathcal{F}, \mathbb{P})$ is the underlying probability space. If $\mathbb{P}(\mathcal{A}) > 0$, for $\mathcal{A} \in \mathcal{F}$, the random vector $(\mathbf{W} \mid \mathcal{A})$ is the measurable mapping $\mathcal{A} \rightarrow \mathbb{R}^{pk}$ such that, $\forall \omega \in \mathcal{A}$, $(\mathbf{W} \mid \mathcal{A})(\omega) = \mathbf{W}(\omega)$. The statistics of $(\mathbf{W} \mid \mathcal{A})$ are then computed in the probability space $(\mathcal{A}, \mathcal{F} \wedge \mathcal{A}, \mathbb{P}_{\mathcal{A}})$, where $\mathcal{F} \wedge \mathcal{A} = \{B \cap \mathcal{A}, B \in \mathcal{F}\}$ and $\forall B \in \mathcal{F}$, $\mathbb{P}_{\mathcal{A}}(B) = \mathbb{P}(B)/\mathbb{P}(\mathcal{A})$.

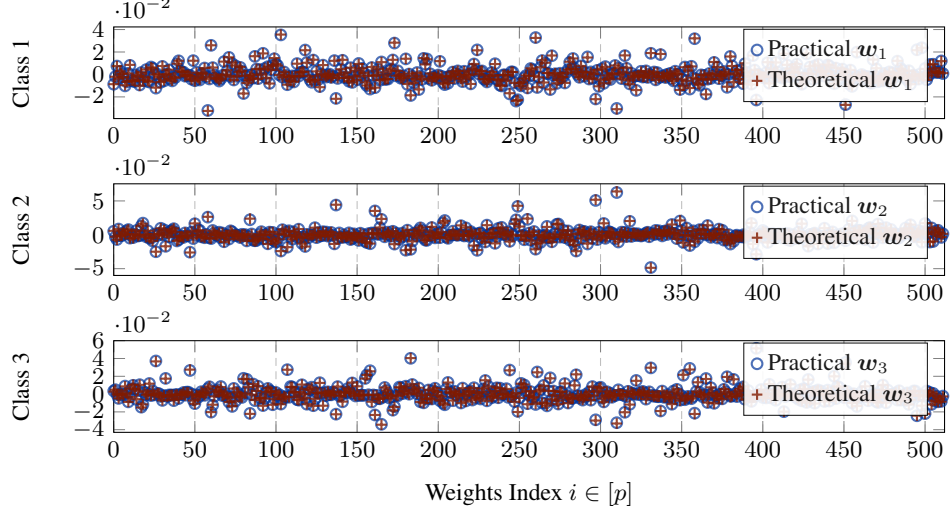


Figure 2: Learned weights (in blue circles) versus their theoretical estimations (in red crosses) as per Theorem 3.2. The used data are Resnet15 [SIVA17] representations ($p = 512$) of images generated by the BigGAN model [BDS18] which are concentrated vectors [SLTC20]. We considered $k = 3$ classes which are: *hamburger*, *mushroom*, *pizza*. $n = 3000$ and the regularization constants $\lambda_1, \lambda_2, \lambda_3 = 1.5$. The data are normalized such that their norm is $0.1 \cdot \sqrt{p}$ to ensure \mathcal{A}_X .

Since their observable diameter ($\sqrt{\log n/n}$) vanishes at large n , it therefore entails from Theorem 3.1 that the random weights vector \mathbf{W} tend to be deterministic as p, n grow large. The subsequent result further characterizes its first and second order statistics at large p, n .

Theorem 3.2 (Asymptotic statistics of \mathbf{W}). *Define the statistics*

$$\mathbf{m}_W \equiv \mathbb{E}[\mathbf{W}], \quad \mathbf{C}_W \equiv [\mathbf{W}\mathbf{W}^\top] - \mathbf{m}_W \mathbf{m}_W^\top.$$

Then, under Assumptions 1 and 2 and additional assumptions on ϕ and $(\lambda_\ell)_{\ell \in [k]}$ provided in Section 2 of the Supplementary Material (Assumptions 3, 4 and 5), there exists a deterministic mapping $\mathcal{F}_{\mu, \Sigma} = \mathcal{F}_{\mu, \Sigma}(\mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k) : \mathbb{R}^{pk} \times \mathcal{M}_{pk} \rightarrow \mathbb{R}^{pk} \times \mathcal{M}_{pk}$ depending only on the statistics μ_1, \dots, μ_k and $\Sigma_1, \dots, \Sigma_k$ of \mathbf{X} , such that the equation

$$(\mathbf{m}, \mathbf{C}) = \mathcal{F}_{\mu, \Sigma}(\mathbf{m}, \mathbf{C}) \quad \text{with} \quad \mathbf{m} \in \mathbb{R}^{pk}, \quad \mathbf{C} \in \mathcal{M}_{pk}$$

admits a unique solution $(\bar{\mathbf{m}}_W, \bar{\mathbf{C}}_W)$. Besides,

$$\|\bar{\mathbf{m}}_W - \mathbf{m}_W\| \leq O\left(\sqrt{\log n/n}\right) \quad \text{and} \quad \|\bar{\mathbf{C}}_W - \mathbf{C}_W\|_* \leq O\left(\sqrt{\log n/n}\right).$$

The exact expression of $\mathcal{F}_{\mu, \Sigma}$, explicitly given in (8), is rather elaborate and of limited interest at this point. Section 5 provides more extensive details.

The central outcome of Theorem 3.2 is that, under the data concentration Assumption 1, the behavior of the Softmax classifier only depends on the class-wise means and covariances of the input data. This arises as a direct consequence of the Lipschitz character of the Softmax classifier which preserves concentration (by the stability result of Remark 2.1), and of the presence of a projection of the parameter vectors \mathbf{w}_ℓ onto the concentrated data \mathbf{x}_i at the core of the optimization formulation: according to Theorem 2.2, these projections induce an asymptotic Gaussian behavior with mean and variance depending only on the first statistics of the data and the weights vector \mathbf{W} .

Once the Softmax classifier is trained, the probability for a new datum \mathbf{x} to belong to class $\ell \in [k]$ is explicitly given by $p_\ell(\mathbf{x}) = \phi(\mathbf{w}_\ell^\top \mathbf{x}) / \sum_{j \in [k]} \phi(\mathbf{w}_j^\top \mathbf{x})$. As a consequence of Theorem 2.2, $\mathbf{w}_\ell^\top \mathbf{x}$ has a high probability to be Gaussian (since \mathbf{x} is concentrated and \mathbf{w}_ℓ has a deterministic behavior). The performances of the Softmax classifier are therefore theoretically tractable.

Corollary 3.3 (Generalization performance of the Softmax classifier). *For $\ell \in [k]$, there exists $\bar{\kappa}^\ell \in \mathbb{R}^{k-1}$ and $\bar{\mathbf{K}}^\ell \in \mathcal{M}_{k-1}$ both depending only on μ_1, \dots, μ_k and $\Sigma_1, \dots, \Sigma_k$ such that the*

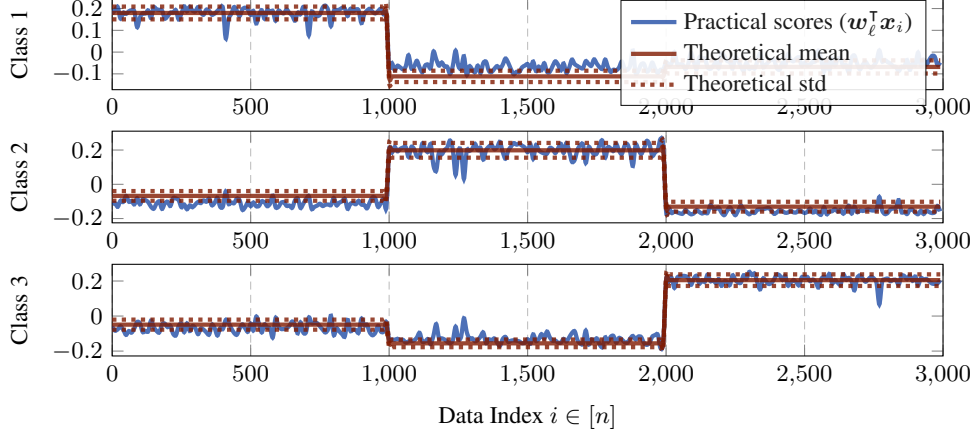


Figure 3: Scores (in blue) versus their theoretical estimations (in red) as per Corollary 3.3, with the theoretical means (through $\bar{\kappa}^\ell$) and standard deviations (through \bar{K}^ℓ), on a test set independent from the training set. The used data are Resnet15 [SIVA17] representations ($p = 512$) of images generated by the BigGAN model [BDS18] which are concentrated vectors [SLTC20]. We considered $k = 3$ classes which are: *hamburger*, *mushroom*, *pizza*. $n = 3000$ and the regularization constants $\lambda_1, \lambda_2, \lambda_3 = 1.5$. The data are normalized such that their norm is $0.1 \cdot \sqrt{p}$ to ensure \mathcal{A}_X .

asymptotic misclassification error $E_t(\mathbf{x} \in \mathcal{C}_\ell)$ of a new datum \mathbf{x} belonging to class $\ell \in [k]$ defined as $E_t(\mathbf{x} \in \mathcal{C}_\ell) \equiv 1 - \mathbb{P}(\forall j \in [k] \setminus \{\ell\} : p_\ell(\mathbf{x}) \geq p_j(\mathbf{x}))$ is

$$E_t(\mathbf{x} \in \mathcal{C}_\ell) = 1 - \mathbb{P}(\mathbf{Z}_\ell \in \mathbb{R}_+^{k-1}) \quad \text{with} \quad \mathbf{Z}_\ell \sim \mathcal{N}(\bar{\kappa}^\ell, \bar{K}^\ell). \quad (3)$$

In essence, Corollary 3.3 states that the generalization performance of the Softmax classifier reduces to the cumulative distribution of a low-dimensional Gaussian vector, the mean and covariance of which only depend on the class-wise means and covariances of the input data. This demonstrates the remarkable *universality* property of the Softmax classifier with respect to the data distribution, which we recall is only requested to satisfy a very loose concentration behavior (Assumption 1). The exact expressions of $\bar{\kappa}^\ell$ and \bar{K}^ℓ , along with a justification of the corollary, are provided in Sections 3–4 of the Supplementary Material.

3.2 Experimental validation

This section provides an experimental setup to support our theoretical findings. The input data $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ are independent Resnet15⁵ representations of size $p = 512$ [SIVA17] of images generated by the BigGAN generative adversarial network model [BDS18]: as such, being the composition of two neural networks (BigGAN and Resnet15) applied to random standard Gaussian noise (as per the BigGAN model), \mathbf{X} is concentrated by construction and satisfies Assumption 1, as requested (see [SLTC20] for a detailed analysis of the Lipschitz properties of these networks). Under this setting, Figure 2 depicts the learned Softmax weights against their expected large n, p asymptotics as per Theorem 3.2. Despite the finite p, n setting of the simulation, a perfect match is observed between the learned weights and the theoretical predictions. In Section 5 of the Supplementary Material, further experiments are performed on *real images* from the ImageNet dataset [DDS⁺09], which again show a perfect match between theory and practice, thereby strongly suggesting that the conclusions of Theorem 3.2 extend to real data.

Figure 3 next displays the class-wise scores of a practical Softmax classifier on an independent test set against their estimated statistics according to Corollary 3.3. An almost perfect match is again observed between empirical values and theoretical statistics. Section 5 of the Supplementary Material reports similar outputs for real ImageNet data. We importantly stress that, as per Corollary 3.3, the theoretical estimates were obtained using *only the empirical class-wise means and covariances* of the input data. Figure 3 thus confirms the theoretically predicted universality of the Softmax classifier. A Python implementation is attached for reproducibility of these experiments.

⁵We used its Pytorch implementation [PGM⁺19] pre-trained on the Imagenet dataset [DDS⁺09].

4 Concluding Remarks

As a consequence of Corollary 3.3, we have demonstrated that, even though the Softmax classifier has a non-linear nature, a property supposedly useful to extract “deep” non-linear features, for n, p rather large, the input data are in fact treated as if they were distributed as a mere Gaussian mixture model. This large dimensional universality phenomenon fundamentally revisits the conventional insights acquired along the years on non-linear classification methods.

As an aftermath, the Softmax classifier being optimal for Gaussian mixture inputs with common covariance, our study is strongly suggestive of the optimality of Softmax as the last layer of a deep neural network classifier.

Yet, the present study assumes a clear-cut separation between a back-end network training isolated from the front-end Softmax layer (as depicted in Figure 1). A thorough validation of the equivalence between full network training and this divided approach is a necessary final step to confirm the claimed optimality and anticipate the performances of Softmax classification.

5 Proof of the Main Results

We now provide the main ingredients to obtain the result of Theorem 3.2, which mainly unfolds from two essential steps: (i) the control of the statistical dependencies between \mathbf{W} and $\tilde{\mathbf{X}}$, presented in Subsection 5.1, and (ii) the estimation of the statistics \mathbf{m}_W and \mathbf{C}_W of the weight vector \mathbf{W} (in Subsection 5.2). We start by reformulating (1) in the compact and convenient form

$$\Lambda \mathbf{W} = \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{x}}_i f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}) \quad (4)$$

where $\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n] \in \mathcal{M}_{kp, kn}$, $f(\tilde{\mathbf{X}}^\top \mathbf{W}) \in \mathbb{R}^{kn}$ and $\Lambda = \text{Diag}(\lambda_1, \dots, \lambda_k) \otimes \mathbf{I}_p \in \mathcal{M}_{kp}$. For lack of space, the expressions of $\tilde{\mathbf{x}}_i \in \mathcal{M}_{pk, k}$ (which relates directly to \mathbf{x}_i), and of the functions $(f_i)_{i \in [n]}$ are deferred to Section 1 of the Supplementary Material.

5.1 Control of the dependencies

Applying the expectation operator both sides to (4), the main technical difficulty arises from the evaluation of $\mathbb{E}[\tilde{\mathbf{x}}_i f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W})]$ due to the elaborate dependencies between the weight vector \mathbf{W} and the data $\tilde{\mathbf{x}}_i$. Note that $\tilde{\mathbf{x}}_i^\top \mathbf{W}$ *a priori* has no reason of being Gaussian (even in the limit) and the performances of the Softmax classifier may depend on high order statistics of \mathbf{X} . These statistical dependencies are dealt with by introducing a mapping $\mathbf{W}_{-i} : [0, 1] \rightarrow \mathbb{R}^{pk}$, defined for $i \in [n]$, as the unique solution to:

$$\forall t \in [0, 1] : \Lambda \mathbf{W}_{-i}(t) = \frac{1}{n} \sum_{j \neq i} \tilde{\mathbf{x}}_j f_j(\tilde{\mathbf{x}}_j^\top \mathbf{W}_{-i}(t)) + \frac{1}{n} t \tilde{\mathbf{x}}_i f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i}(t)). \quad (5)$$

This mapping can be seen as a path between the weights vector $\mathbf{W} = \mathbf{W}_{-i}(1)$ of the Softmax classifier and $\mathbf{W}_{-i}(0)$ which is completely independent of $\tilde{\mathbf{x}}_i$ and which will be simply denoted \mathbf{W}_{-i} . Using the inverse function theorem, the mapping $t \mapsto \mathbf{W}_{-i}(t)$ is differentiable, we then deduce the following central close form formula:

$$\mathbf{W}_{-i}'(t) = \frac{1}{n} \mathbf{Q}_{-i}(t) \tilde{\mathbf{x}}_i \chi_i'(t) \quad \text{with} \quad \mathbf{Q}_{-i}(t) \equiv \left(\Lambda - \frac{1}{n} \tilde{\mathbf{X}}_{-i} \mathbf{D}^{(i)}(t) \tilde{\mathbf{X}}_{-i}^\top \right)^{-1} \in \mathcal{M}_{kp}, \quad (6)$$

where $\chi_i(t) \equiv t f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i}(t))$, $\mathbf{D}_j^{(i)}(t) \equiv df_j|_{\tilde{\mathbf{x}}_j^\top \mathbf{W}_{-i}(t)} \in \mathcal{M}_k$, $\mathbf{D}^{(i)}(t) \in \mathcal{M}_{kn}$ is a block-diagonal matrix with block-diagonal matrices $\mathbf{D}_j^{(i)}(t) \in \mathcal{M}_k$ for $j \in [n]$, and finally $\tilde{\mathbf{X}}_{-i} \equiv (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{i-1}, 0, \tilde{\mathbf{x}}_{i+1}, \dots, \tilde{\mathbf{x}}_n) \in \mathcal{M}_{pk, kn}$ ($\tilde{\mathbf{X}}_{-i}$ is independent of $\tilde{\mathbf{x}}_i$).

Relying on concentration of measure arguments [LC20], the random vector $\mathbf{Q}_{-i}(t) \tilde{\mathbf{x}}_i$ is almost constant in terms of t and thus almost equal to $\mathbf{Q}_{-i}(0) \tilde{\mathbf{x}}_i$. Moreover, the fact that $\mathbf{Q}_{-i}(0)$ (also simply denoted \mathbf{Q}_{-i}) is independent of $\tilde{\mathbf{x}}_i$ allows us to integrate the identity (6) to obtain the core

result of the article relating \mathbf{W} and \mathbf{W}_{-i} . Specifically, we have the following concentration inequality, for some constants $C, c > 0$:

$$\forall t > 0 : \mathbb{P} \left(\left\| \tilde{\mathbf{x}}_i^\top \mathbf{W} - \tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i} + \frac{1}{n} \tilde{\mathbf{x}}_i^\top \mathbf{Q}_{-i} \tilde{\mathbf{x}}_i f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}) \right\| \geq t \mid \mathcal{A}_X \right) \leq C e^{-cnt^2 / \log n}. \quad (7)$$

5.2 Estimating the mean and covariance of the Softmax weights

By breaking the statistical dependencies of the problem through \mathbf{W}_{-i} , we may now access and estimate the statistics \mathbf{m}_W and \mathbf{C}_W . This precisely comes from a deterministic approximation of the quadratic form $\frac{1}{n} \tilde{\mathbf{x}}_i^\top \mathbf{Q}_{-i} \tilde{\mathbf{x}}_i$ in (7) in the large n, p limit, a result inspired by [LC20]:

Proposition 5.1. *For any $\ell \in [k]$, let n_ℓ be the number of columns of \mathbf{X} in class \mathcal{C}_ℓ and, for any block diagonal matrix $\Delta = \text{Diag}(\Delta_\ell)_{1 \leq \ell \leq k} \in \mathcal{M}_{k^2}$ ($\forall \ell \in [k], \Delta_\ell \in \mathcal{M}_k$), let*

$$\bar{\mathbf{Q}}(\Delta) \equiv \left(\Lambda - \sum_{a=1}^k \frac{n_a}{n} \Gamma_a(\Delta_a) \otimes \mathbf{S}_a \right)^{-1} = \begin{pmatrix} \bar{\mathbf{Q}}_{1,1}(\Delta) & \dots & \bar{\mathbf{Q}}_{1,k}(\Delta) \\ \vdots & \ddots & \vdots \\ \bar{\mathbf{Q}}_{k,1}(\Delta) & \dots & \bar{\mathbf{Q}}_{k,k}(\Delta) \end{pmatrix} \in \mathcal{M}_{kp},$$

where $\Gamma_\ell(\Delta_a) = \mathbb{E}[(\mathbf{I}_k - \mathbf{D}_j^{-i}(0)\Delta_a)^{-1} \mathbf{D}_j^{-i}(0)]$ for \mathbf{x}_j in class \mathcal{C}_ℓ and $\mathbf{S}_\ell = \mathbb{E}[\mathbf{x}_j \mathbf{x}_j^\top]$. Then the fixed point equation

$$\Delta_\ell = \left[\frac{1}{n} \text{Tr}(\mathbf{S}_\ell \bar{\mathbf{Q}}_{a,b}(\Delta)) \right]_{1 \leq a,b \leq k}$$

admits a unique solution $\Delta \in \mathcal{M}_{k^2}$ that satisfies, for any \mathbf{x}_i is in class $\ell \in [k]$,

$$\forall t > 0 : \mathbb{P} \left(\left\| \frac{1}{n} \tilde{\mathbf{x}}_i^\top \mathbf{Q}_{-i} \tilde{\mathbf{x}}_i - \Delta_\ell \right\| \geq t \mid \mathcal{A}_X \right) \leq C e^{-cnt^2 / \log n} \quad \text{for some constants } C, c > 0.$$

From this result, using the identity (7), we then obtain an estimation for $\tilde{\mathbf{x}}_i^\top \mathbf{W}$:

Proposition 5.2. *For any $\mathbf{v} \in \mathbb{R}^k$, there exists a unique point $g_i(\mathbf{v}) \in \mathbb{R}^k$ satisfying:*

$$g_i(\mathbf{v}) = \mathbf{v} - \Delta_i f_i(g_i(\mathbf{v})),$$

and, for some constants $C, c > 0$,

$$\mathbb{P}(\|\tilde{\mathbf{x}}_i^\top \mathbf{W} - g_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i})\| \geq t \mid \mathcal{A}_X) \leq C e^{-cnt^2 / \log n}.$$

Therefor, letting $h_i = f_i \circ g_i$, by Hölder's inequality [Fin92], since $\tilde{\mathbf{x}}_i$ is concentrated,

$$\begin{aligned} \left\| \mathbf{m}_W - \frac{1}{n} \sum_{i=1}^n \Lambda^{-1} \mathbb{E}[\tilde{\mathbf{x}}_i h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i})] \right\| &= O \left(\sqrt{\frac{\log n}{n}} \right) \\ \left\| \mathbf{C}_W - \frac{1}{n^2} \sum_{i=1}^n \Lambda^{-1} \mathbb{E}[\tilde{\mathbf{x}}_i h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i}) h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i})^\top \tilde{\mathbf{x}}_i^\top] \Lambda^{-1} \right\|_* &= O \left(\sqrt{\frac{\log n}{n}} \right) \end{aligned}$$

Knowing from Theorem 2.2 that $\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i}$ is asymptotically Gaussian, $\mathbb{E}[\tilde{\mathbf{x}}_i h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i})]$ and $\mathbb{E}[\tilde{\mathbf{x}}_i h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i}) h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i})^\top \tilde{\mathbf{x}}_i^\top]$ can be explicitly evaluated (for instance using Stein's Lemma [LN08]), and only depend on the statistical means and covariances of $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$ and of \mathbf{W}_{-i} (which has the same statistics as \mathbf{W}). Their exact expressions are provided in Section 3 of the Supplementary Material. Finally, let us introduce the $2k$ functions $m_1, \dots, m_k : \mathbb{R}^{kp} \times \mathcal{M}_{kp} \rightarrow \mathbb{R}^{kp}$ and $c_1, \dots, c_k : \mathbb{R}^{kp} \times \mathcal{M}_{kp} \rightarrow \mathcal{M}_{kp}$ defined, $\forall i \in [n]$, by

$$\begin{aligned} m_{k(i)}(\mathbf{m}_W, \mathbf{C}_W) &= \Lambda^{-1} \mathbb{E}[\tilde{\mathbf{x}}_i h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i})] \\ c_{k(i)}(\mathbf{m}_W, \mathbf{C}_W) &= \frac{1}{n} \Lambda^{-1} \mathbb{E}[\tilde{\mathbf{x}}_i h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i}) h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i})^\top \tilde{\mathbf{x}}_i^\top] \Lambda^{-1}, \end{aligned}$$

where $k(i)$ denotes the class of $\tilde{\mathbf{x}}_i$ and \mathbf{m}_W and \mathbf{C}_W are respectively the mean and covariance of \mathbf{W} . The mappings $(m_\ell)_{1 \leq \ell \leq k}$ and $(c_\ell)_{1 \leq \ell \leq k}$ are uniquely determined by the means μ_1, \dots, μ_n and the covariances $\Sigma_1, \dots, \Sigma_n$. In particular, the deterministic pair $(\bar{\mathbf{m}}_W, \bar{\mathbf{C}}_W)$, defined as the unique solution of

$$\bar{\mathbf{m}}_W = \sum_{\ell=1}^k \frac{n_\ell}{n} m_\ell(\bar{\mathbf{m}}_W, \bar{\mathbf{C}}_W) \quad \text{and} \quad \bar{\mathbf{C}}_W = \sum_{\ell=1}^k \frac{n_\ell}{n} c_\ell(\bar{\mathbf{m}}_W, \bar{\mathbf{C}}_W), \quad (8)$$

is a good approximation for $(\mathbf{m}_W, \mathbf{C}_W)$ as stated in Theorem 3.2.

Broader Impact

In this work, we provided a theoretical analysis of the widely used Softmax component in neural network classifiers, by deriving the exact high dimensional asymptotic performances of this classifier in terms of its inputs statistics. An important positive impact of this work is the in-depth analysis and understanding of the interplay between the data representations and the last trained layer of neural networks. The developed theoretical approach in this study presents fair and non-offensive societal consequences.

References

- [BDS18] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- [CHM⁺15] Anna Choromanska, Mikael Henaff, Michael Mathieu, Gérard Ben Arous, and Yann LeCun. The loss surfaces of multilayer networks. In *Artificial intelligence and statistics*, pages 192–204, 2015.
- [dCPS⁺18] Remi Tachet des Combes, Mohammad Pezeshki, Samira Shabanian, Aaron Courville, and Yoshua Bengio. On the learning dynamics of deep neural networks. *arXiv preprint arXiv:1809.06848*, 2018.
- [DDS⁺09] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [FGP07] B. Fleury, O. Guédon, and G. Paouris. A stability result for mean width of l_p -centroid bodies. *Advances in Mathematics*, 214:865–877, 2007.
- [Fin92] Helmut Finner. A generalization of holder’s inequality and some probability inequalities. *The Annals of probability*, pages 1893–1901, 1992.
- [GCM18] Samantha Guerriero, Barbara Caputo, and Thomas Mensink. Deep nearest class mean classifiers. In *International Conference on Learning Representations, Worskhop Track*, 2018.
- [GP17] Bolin Gao and Lacra Pavel. On the properties of the softmax function with application in game theory and reinforcement learning. *arXiv preprint arXiv:1704.00805*, 2017.
- [GPAM⁺14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [HRU⁺17] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in neural information processing systems*, pages 6626–6637, 2017.
- [KGC18] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7482–7491, 2018.
- [Kla07] B. Klartag. A central limit theorem for convex sets. *Inventiones mathematicae volume*, 168:pages91–131, 2007.
- [KXR⁺19] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis. Decoupling representation and classifier for long-tailed recognition, 2019.
- [LC20] Cosme Louart and Romain Couillet. Concentration of measure and large random matrices with an application to sample covariance matrices. *submitted to Random Matrices: Theory and Applications*, 2020.

- [Led05] Michel Ledoux. *The concentration of measure phenomenon*. Number 89. American Mathematical Soc., 2005.
- [LLC⁺18] Cosme Louart, Zhenyu Liao, Romain Couillet, et al. A random matrix approach to neural networks. *The Annals of Applied Probability*, 28(2):1190–1248, 2018.
- [LN08] Zinoviy Landsman and Johanna Nešlehová. Stein’s lemma for elliptical random vectors. *Journal of Multivariate Analysis*, 99(5):912–927, 2008.
- [LWL⁺17] Xuezhi Liang, Xiaobo Wang, Zhen Lei, Shengcai Liao, and Stan Z Li. Soft-margin softmax for deep classification. In *International Conference on Neural Information Processing*, pages 413–421. Springer, 2017.
- [MVPC13] Thomas Mensink, Jakob Verbeek, Florent Perronnin, and Gabriela Csurka. Distance-based image classification: Generalizing to new classes at near-zero cost. *IEEE transactions on pattern analysis and machine intelligence*, 35(11):2624–2637, 2013.
- [NBA⁺18] Roman Novak, Yasaman Bahri, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Sensitivity and generalization in neural networks: an empirical study. *arXiv preprint arXiv:1802.08760*, 2018.
- [PB17] Jeffrey Pennington and Yasaman Bahri. Geometry of neural network loss surfaces via random matrix theory. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2798–2806. JMLR. org, 2017.
- [PGM⁺19] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019.
- [PRU⁺18] Kristina Preuer, Philipp Renz, Thomas Unterthiner, Sepp Hochreiter, and Günter Klambauer. Fréchet chemblnet distance: A metric for generative models for molecules. *arXiv preprint arXiv:1803.09518*, 2018.
- [PW17] Jeffrey Pennington and Pratik Worah. Nonlinear random matrix theory for deep learning. In *Advances in Neural Information Processing Systems*, pages 2637–2646, 2017.
- [SIVA17] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Thirty-first AAAI conference on artificial intelligence*, 2017.
- [SLTC20] Mohamed El Amine Seddik, Cosme Louart, Mohamed Tamaazousti, and Romain Couillet. Random matrix theory proves that deep learning representations of gan-data behave as gaussian mixtures. *arXiv preprint arXiv:2001.08370*, 2020.
- [SMG13] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- [YKYR18] Chih-Kuan Yeh, Joon Kim, Ian En-Hsu Yen, and Pradeep K Ravikumar. Representer point selection for explaining deep neural networks. In *Advances in Neural Information Processing Systems*, pages 9291–9301, 2018.
- [YW19] Yaqiong Yao and HaiYing Wang. Optimal subsampling for softmax regression. *Statistical Papers*, 60(2):235–249, 2019.

Supplementary Material: The Unexpected Deterministic and Universal Behavior of Large Softmax Classifiers

Mohamed El Amine Seddik^{1,2}, Cosme Louart^{1,3}, Romain Couillet^{2,3}, Mohamed Tamaazousti¹

¹CEA List, ²Centralesupélec, ³GIPSA Lab Grenoble-Alpes University

¹Palaiseau, France, ²Gif-sur-Yvette, France, ³Grenoble, France

firstname.lastname@{cea,centralesupelec}.fr

1 Notations

We consider the following notations: $\tilde{\mathbf{x}}_i = \begin{pmatrix} \mathbf{x}_i & & \\ & \ddots & \\ & & \mathbf{x}_i \end{pmatrix} \in \mathcal{M}_{pk,k} \forall i \in [n]$, and introduce the functions:

$$\begin{aligned} f_i : \mathbb{R}^k &\longrightarrow \mathbb{R}^k \\ \mathbf{v} &\longmapsto \left[\frac{\phi(\mathbf{v}_\ell)}{\sum_{j=1}^k \phi(\mathbf{v}_j)} \sum_{j=1}^k y_j^{(i)} \psi(\mathbf{v}_j) - y_\ell^{(i)} \psi(\mathbf{v}_\ell) \right]_{1 \leq \ell \leq k}, \end{aligned}$$

Introduce the matrix $\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n) \in \mathcal{M}_{kp,kn}$, equation (1) of the Main Paper becomes:

$$\mathbf{\Lambda} \mathbf{W} = \frac{1}{n} \tilde{\mathbf{X}} f(\tilde{\mathbf{X}}^\top \mathbf{W}) \quad (1)$$

where $\mathbf{\Lambda} = \text{Diag}(\mathbf{\Lambda}_1, \dots, \mathbf{\Lambda}_k) \otimes \mathbf{I}_p$ and $f(\tilde{\mathbf{X}}^\top \mathbf{W}) \in \mathbb{R}^{kn}$ concatenates the elements $f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}) \in \mathbb{R}^k$.

2 Concentration of the classifier weights

In this section, we provide the conditions under which the weights vector \mathbf{W} is concentrated as per Theorem 3.1 of the Main Paper. Rewriting (1) as $\mathbf{W} = \Psi(\mathbf{W})$, one sees that Ψ is contracting if $\|d\Psi|_{\mathbf{w}}\| \leq 1 - \varepsilon$ for some $\varepsilon > 0$. Now, since $\|d\Psi|_{\mathbf{w}}\| \leq \frac{1}{n} \|\mathbf{\Lambda}^{-1}\| \|\tilde{\mathbf{X}}\|^2 \|df|_{\tilde{\mathbf{X}}^\top \mathbf{w}}\|$ we first need $\|df\|_\infty = \sup_{\mathbf{v} \in \mathbb{R}^k} \|df|_{\mathbf{v}}\|$ to be bounded, that is ensured by:

Assumption 3 (Regularity). $\|\frac{\phi'}{\phi}\|_\infty \leq \infty$ and $\|\frac{\phi''}{\phi}\|_\infty \leq \infty$.

Besides, we also need to be able to bound $\frac{1}{n} \|\tilde{\mathbf{X}}\|^2 = \frac{1}{n} \|\mathbf{X}\|^2$. The spectral norm being lower than the Frobenius norm (involved in Assumption 1 giving the concentration of $\tilde{\mathbf{X}}$), it is a 1-Lipschitz observation of $\tilde{\mathbf{X}}$ thus there exists two constants $C, c > 0$ (independent of p, n), such that:

$$\mathbb{P} \left(\left| \frac{1}{\sqrt{n}} \|\tilde{\mathbf{X}}\| - \frac{1}{\sqrt{n}} \mathbb{E}[\|\tilde{\mathbf{X}}\|] \right| \geq t \right) \leq C e^{-cnt^2}. \quad (2)$$

This concentration inequality proves that the random variable $\frac{1}{\sqrt{n}} \|\tilde{\mathbf{X}}\|$ is almost deterministic and equal to $\frac{1}{\sqrt{n}} \mathbb{E}[\|\tilde{\mathbf{X}}\|]$. Bounding this last quantity necessities the following result:

$$\left| \mathbb{E}[\|\tilde{\mathbf{X}}\|] - \|\mathbb{E}[\tilde{\mathbf{X}}]\| \right| = O(\sqrt{p+n}), \quad (3)$$

It suffices to have a bound on $\|\mathbb{E}[\tilde{\mathbf{X}}]\|$ which comes from the following assumption:

Assumption 4. $\sup_{1 \leq \ell \leq k} \|\boldsymbol{\mu}_\ell\| \leq O(1)$ ¹

This assumption implies that $\|\mathbb{E}[\tilde{\mathbf{X}}]\| = O(\sqrt{n})$, and we can deduce from equation 3 that $\mathbb{E}[\|\tilde{\mathbf{X}}\|] = O(\sqrt{n})$ (recall from Assumption 2 that $p = O(n)$). Now, assuming

Assumption 5. $\frac{1}{n} \mathbb{E}[\|\mathbf{X}\|^2] \|df\|_\infty \|\Lambda^{-1}\| < 1$.

and noting $\varepsilon = \frac{1}{2} - \frac{1}{2} \|\Lambda^{-1}\| \|f'\|_\infty \mathbb{E}[\|\tilde{\mathbf{X}}\|/\sqrt{n}]^2$, it can be deduced from equation 2 that the event:

$$\mathcal{A}_X = \left\{ \frac{1}{n} \left| \|\tilde{\mathbf{X}}\|^2 - \mathbb{E}[\|\tilde{\mathbf{X}}\|^2] \right| \leq \frac{\varepsilon}{2 \|\Lambda^{-1}\| \|df\|_\infty} \right\}$$

has a very high probability to happen (bigger than $1 - Ce^{-cn}$ for two constant $C, c > 0$) and it satisfies $\mathcal{A}_X \subset \{\|d\Psi\|_\infty \leq 1 - \varepsilon\}$. As a consequence, our fixed point \mathbf{W} is uniquely determined under the event \mathcal{A}_X that appears in the concentration result of Theorem 3.1. The concentration of the random vector \mathbf{W} is far from being trivial because \mathbf{W} is not explicitly written as a Lipschitz transformation of \mathbf{X} and additional tools are necessary to prove the concentration of \mathbf{W} . The complete proof will be provided in an extended version of this paper.

3 Computation of mean and covariance of the parameter vector.

This section provides the exact computation of the fixed point equation (8) of the Main Paper which is the main result of our study as per Theorem 3.2. We start by introducing the main tools to perform the calculations.

Theorem 2.2 giving the central limit theorem for concentrated vectors was originally proven for uniform distributions on convex subspaces of \mathbb{R}^p , but it was quickly understood that the result is true for a larger class of random vectors satisfying a so-called “thin shell property” (see [Fre19] for a simple and complete proof of this inference). The thin shell property expresses the fact that a random vector \mathbf{X} lies principally on a thin shell around a sphere with the following inequality satisfied for some $\varepsilon > 0$

$$\mathbb{P} \left(\left| \frac{\|\mathbf{X}\|}{\sqrt{p}} - 1 \right| \geq \varepsilon \right) \leq \varepsilon.$$

For a concentrated random vector $\mathbf{X} \propto \mathcal{E}_2$, such that $\mathbb{E}[\mathbf{X}\mathbf{X}^T] = I_p$ and $\sqrt{p} = O(\mathbb{E}[\|\mathbf{X}\|])$ ² the norm being a 1-Lipschitz observation, we know (as in equation 2) that there exist two constants $C, c > 0$ such that:

$$\mathbb{P}(\|\mathbf{X}\| - \mathbb{E}[\|\mathbf{X}\|] \geq \varepsilon) \leq Ce^{-c\varepsilon^2}.$$

Integrating this concentration inequality for $\varepsilon \in [0, \infty)$ with Fubini Theorem, we have $\forall r \geq 2$:

$$\mathbb{E} \left[\left| \frac{\|\mathbf{X}\|}{\mathbb{E}[\|\mathbf{X}\|]} - 1 \right|^r \right] \leq C(r/2c)^{r/2} \mathbb{E}[\|\mathbf{X}\|]^r.$$

Therefore, since $\mathbb{E}[\|\mathbf{X}\|] \leq \sqrt{\mathbb{E}[\|\mathbf{X}\|^2]} = \sqrt{p} = O(\mathbb{E}[\|\mathbf{X}\|])$, we can deduce from Hölder’s inequality that:

$$\mathbb{E} \left[\left| \frac{\|\mathbf{X}\|^2}{\mathbb{E}[\|\mathbf{X}\|]^2} - 1 \right| \right] \leq \sqrt{\mathbb{E} \left[\left| \frac{\|\mathbf{X}\|}{\mathbb{E}[\|\mathbf{X}\|]} - 1 \right|^2 \right] \mathbb{E} \left[\left| \frac{\|\mathbf{X}\|}{\mathbb{E}[\|\mathbf{X}\|]} + 1 \right|^2 \right]} = O(1/\sqrt{p}),$$

¹For classification problems, if $\forall a, b \in [k], a \neq b, \|\boldsymbol{\mu}_a - \boldsymbol{\mu}_b\| \gg 1$ then the classification becomes trivial in the large dimensional regime, a reasonable assumption then is $\|\boldsymbol{\mu}_a - \boldsymbol{\mu}_b\| \leq O(1)$. However, in cases where $\sup_{1 \leq \ell \leq k} \|\boldsymbol{\mu}_\ell\|$ is of order $O(\sqrt{p})$, it is still possible to work with the data matrix $\tilde{\mathbf{X}} - \frac{1}{n} \tilde{\mathbf{X}} \mathbf{1}_n \mathbf{1}_n^T$ that satisfies Assumption 4 but has the drawback that the columns are then dependent. However, this is not a big issue since this dependence is very small and it can be managed thanks to some technical consideration that we want to avoid here.

²This is a very loose hypothesis needed to set that $\frac{\|\mathbf{X}\|}{\mathbb{E}[\|\mathbf{X}\|]}$ is sufficiently concentrated. Generally if $\sqrt{p} \ll \mathbb{E}[\|\mathbf{X}\|]$, that means that one can obtain a better concentration than $\mathbf{X} \propto \mathcal{E}_2$

from which we conclude that: $\mathbb{E}[\|X\|]^2 = \mathbb{E}[\|X\|^2] + O(\sqrt{p})$ and therefore $\mathbb{E}[\|X\|] = \sqrt{p} + O(p^{1/4})$. Choosing $\varepsilon = p^{-1/4}$ yields from the concentration of $\|X\|$ to the existence of some constant $K > 0$ such that:

$$\mathbb{P}\left(\left|\frac{\|X\|}{\sqrt{p}} - 1\right| \geq Kp^{-1/4}\right) \leq Ce^{-cp^{1/2}} \leq Kp^{-1/4},$$

for p large enough. We can then infer (see [Fre19]), that the projections on small dimensional vector spaces of a concentrated vector are Gaussian vectors with high probability.

Theorem A (CLT for concentrated vectors [Kla07, FGP07]). *Given a random vector $\mathbf{X} \in \mathbb{R}^p$, and noting G , the cumulative distribution function of a Gaussian variable of zero mean and unit variance. If $\mathbf{X} \propto \mathcal{E}_2$, $\mathbb{E}[\mathbf{X}] = 0$ and $\mathbb{E}[\mathbf{X}\mathbf{X}^T] = I_p$, then for any integer $k \in \mathbb{N}$, small compared to p , for any $\eta \in (0, 1)$, there exists two constants $C, c > 0$ and a set $\Theta \subset \mathcal{S}^{p-1}$ such that $\sigma(\Theta) \geq 1 - \sqrt{p}Ce^{-c\sqrt{p}}$ and $\forall \boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k) \in \Theta^k$, there exists a Gaussian vector $\mathbf{Z} \sim \mathcal{N}(0, \boldsymbol{\theta}^T \boldsymbol{\theta})$ such that:*

$$\forall \mathbf{a} \in \mathbb{R}^k : \sup_{t \in \mathbb{R}} |\mathbb{P}(\mathbf{a}^T \boldsymbol{\theta}^T \mathbf{X} \geq t) - \mathbb{P}(\mathbf{a}^T \mathbf{Z} \geq t)| \leq Cp^{-1/4}.$$

We need a simple preliminary Lemma to state our Stein-like formula for concentrated vectors from Theorem A.

Lemma B. *Given two random variables $X, Y \in \mathbb{R}^k$, if:*

$$\sup_{t \in \mathbb{R}} |\mathbb{P}(\mathbf{a}^T X \geq t) - \mathbb{P}(\mathbf{a}^T Y \geq t)| \leq \varepsilon,$$

then for any differentiable mapping $f : \mathbb{R}^k \rightarrow \mathbb{R}$ integrable and bounded around ∞ by f_∞ :

$$|\mathbb{E}[f(X)] - \mathbb{E}[f(Y)]| \leq \frac{f_\infty + \int |f|}{\varepsilon}$$

Proof. Let us prove it in the case $k = 1$:

$$\begin{aligned} |\mathbb{E}[f(X)] - \mathbb{E}[f(Y)]| &\leq \left| \int_{t=-\infty}^{\infty} f(t)(d\mathbb{P}_X(t) - d\mathbb{P}_Y(t)) \right| \\ &\leq \left| [f(t)(\mathbb{P}(X \geq t) - \mathbb{P}(Y \geq t))]_{-\infty}^{\infty} \right| \\ &\quad + \left| \int_{t=-\infty}^{\infty} f'(t)(\mathbb{P}(X \geq t) - \mathbb{P}(Y \geq t))dt \right| \\ &\leq \frac{f_\infty + \int |f|}{\varepsilon} \end{aligned}$$

□

We have the following Stein-like [LN08] theorem for concentrated vectors which is the central tool to express the fixed point mappings in equation (8) of the Main Paper.

Proposition C. *Let $\mathbf{x} \in \mathbb{R}^p$ be a random vector satisfying the concentration $\mathbf{x} \propto \mathcal{E}_2$ (denote $\mathbf{m} \equiv \mathbb{E}[\mathbf{x}]$ and $\mathbf{C} \equiv \mathbb{E}[\mathbf{x}\mathbf{x}^T]$) and let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be some three times differentiable function such that f, f' and f'' satisfy the hypotheses of Lemma B with $f_\infty, \int |f|, f'_\infty, \int |f'|, f''_\infty, \int |f''| = O(1)$. Then, there exists a subset $\Theta \subset \mathbb{S}^{p-1}$ such that: $\sigma(\Theta) \geq 1 - Ce^{-cp/\log p}$ and $\forall \mathbf{w}, \mathbf{v}, \mathbf{u} \in \Theta$:*

$$\begin{aligned} \mathbb{E}[f(\mathbf{w}^T \mathbf{x})\mathbf{v}^T \mathbf{x}] &= \mathbb{E}[f(\mathbf{w}^T \mathbf{x})\mathbf{v}^T \mathbf{m}] + \mathbb{E}[f'(\mathbf{w}^T \mathbf{x})\mathbf{v}^T \mathbf{C} \mathbf{w}] + O\left(\frac{1}{n^{1/4}}\right) \\ \mathbb{E}[f(\mathbf{w}^T \mathbf{x})\mathbf{v}^T \mathbf{x} \mathbf{x}^T \mathbf{u}] &= \mathbb{E}[f(\mathbf{w}^T \mathbf{x})\mathbf{v}^T (\mathbf{m} \mathbf{m}^T + \mathbf{C}) \mathbf{u}] + \mathbb{E}[f'(\mathbf{w}^T \mathbf{x})\mathbf{v}^T (\mathbf{C} \mathbf{w} \mathbf{m}^T + \mathbf{m} \mathbf{w}^T \mathbf{C}) \mathbf{u}] \\ &\quad + \mathbb{E}[f''(\mathbf{w}^T \mathbf{x})\mathbf{v}^T \mathbf{C} \mathbf{w} \mathbf{w}^T \mathbf{C} \mathbf{u}] + O\left(\frac{1}{n^{1/4}}\right) \end{aligned}$$

Proof. Let us first consider a random vector \mathbf{z} with zero mean and identity covariance satisfying $\mathbf{z} \propto \mathcal{E}_2$. Considering the subset $\Theta \subset \mathbb{S}^{p-1}$ mentioned in Theorem A (for $k = 3$), we know that

$\sigma(\Theta) \geq 1 - Ce^{-cp^{1-\varepsilon}}$ for two constants $C, c > 0$ and furthermore for any $\theta = (\mathbf{w}, \mathbf{v}, \mathbf{u}) \in \Theta^3$, and some Gaussian random vector $\mathbf{Z} \sim \mathcal{N}(0, \theta^\top \theta)$:

$$\forall a \in \mathbb{R}^2 : \sup_{t \in \mathbb{R}} |\mathbb{P}(a^\top \theta_i^\top \mathbf{z} \geq t) - \mathbb{P}(a^\top \mathbf{Z} \geq t)| = O\left(\frac{1}{n^{1/4}}\right).$$

Given a mapping $g : \mathbb{R} \rightarrow \mathbb{R}$, we know thanks to Lemma B and Stein's identity:

$$\mathbf{w}^\top \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})] = \mathbb{E}[g(\mathbf{e}_1^\top \mathbf{Z})\mathbf{e}_1^\top \mathbf{Z}] + O\left(\frac{1}{n^{1/4}}\right) = \mathbb{E}[g'(\mathbf{w}^\top \mathbf{z})]\|\mathbf{w}\|^2 + O\left(\frac{1}{n^{1/4}}\right),$$

where $\mathbf{e}_1 = (1, 0)$. Second, if we note $\tilde{\theta} = (\mathbf{w}, \tilde{\mathbf{v}}, \tilde{\mathbf{u}}) = \mathbf{Q}\mathbf{R}$, the QR-decomposition of θ , and $\mathbf{Q} = (\mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3)$ (of course, $\mathbf{q}_1 = \mathbf{e}_1$), we know that $\mathbf{e}_1^\top \mathbf{Z}$ and $\mathbf{q}_2^\top \mathbf{Z}$ are independent (\mathbf{Z} is Gaussian and $\mathbb{E}[\mathbf{e}_1^\top \mathbf{Z} \mathbf{q}_2^\top \mathbf{Z}] = \mathbf{e}_1^\top \theta^\top \mathbf{q}_2 = \mathbf{w}^\top \tilde{\mathbf{v}} = 0$). We can therefore estimate:

$$\tilde{\mathbf{v}}^\top \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})] = \mathbb{E}[g(\mathbf{e}_1^\top \mathbf{Z})]\mathbb{E}[\mathbf{q}_2^\top \mathbf{Z}] + O\left(\frac{1}{n^{1/4}}\right) = O\left(\frac{1}{n^{1/4}}\right).$$

Combing those 2 estimations, we see that for any differentiable function g $\mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{v}^\top \mathbf{z}] = \mathbb{E}[g'(\mathbf{w}^\top \mathbf{z})]\mathbf{v}^\top \mathbf{w}$. Therefore if we take for g the mapping $t \mapsto f(\mathbf{w}^\top \mathbf{m} + t)$ (satisfying $f(\mathbf{w}^\top \mathbf{x}) = g(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})$), we get the identity:

$$\begin{aligned} \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})\mathbf{v}^\top \mathbf{x}] &= \mathbb{E}[g(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})\mathbf{v}^\top \mathbf{m}] + \mathbb{E}[g(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})]\mathbf{v}^\top \mathbf{C}^{1/2} \mathbf{w} + O\left(\frac{1}{n^{1/4}}\right) \\ &= \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})]\mathbf{v}^\top \mathbf{m} + \mathbb{E}[g'(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})]\mathbf{v}^\top \mathbf{C} \mathbf{w} + O\left(\frac{1}{n^{1/4}}\right) \end{aligned}$$

With the same method, let us first compute:

$$\tilde{\mathbf{v}}^\top \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{z}\mathbf{z}^\top] \mathbf{w} = \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{w}\mathbf{z}]\mathbb{E}[\mathbf{z}^\top \tilde{\mathbf{v}}] = O\left(\frac{1}{n^{1/4}}\right) = \mathbf{w}^\top \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{z}\mathbf{z}^\top] \tilde{\mathbf{u}}.$$

Second:

$$\tilde{\mathbf{v}}^\top \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{z}\mathbf{z}^\top] \tilde{\mathbf{u}} = \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})]\mathbb{E}[\tilde{\mathbf{v}}^\top \mathbf{z}\mathbf{z}^\top \tilde{\mathbf{u}}] + O\left(\frac{1}{n^{1/4}}\right) = \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})]\tilde{\mathbf{v}}^\top \tilde{\mathbf{u}} + O\left(\frac{1}{n^{1/4}}\right).$$

Third:

$$\begin{aligned} \mathbf{w}^\top \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{z}\mathbf{z}^\top] \mathbf{w} &= \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{w}^\top \mathbf{z}\mathbf{w}^\top \mathbf{z}] = \mathbb{E}[g(\mathbf{w}^\top \mathbf{z}) + g'(\mathbf{w}^\top \mathbf{z})\mathbf{w}^\top \mathbf{z}]\|\mathbf{w}\|^2 + \\ &= \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})]\|\mathbf{w}\|^2 + \mathbb{E}[g''(\mathbf{w}^\top \mathbf{z})]\|\mathbf{w}\|^4 + O\left(\frac{1}{n^{1/4}}\right). \end{aligned}$$

Therefore, $\mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\tilde{\mathbf{v}}^\top \mathbf{z}\mathbf{z}^\top \tilde{\mathbf{v}}] = \mathbb{E}[(g(\mathbf{w}^\top \mathbf{z}))\tilde{\mathbf{v}}^\top \tilde{\mathbf{u}} + \mathbb{E}[g''(\mathbf{w}^\top \mathbf{z})]\tilde{\mathbf{v}}^\top \mathbf{w}\mathbf{w}^\top \tilde{\mathbf{u}}]$, and we can conclude as before that:

$$\begin{aligned} \tilde{\mathbf{v}}^\top \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})\mathbf{x}\mathbf{x}^\top] \mathbf{u} &= \tilde{\mathbf{v}}^\top \mathbf{m} \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})] \mathbf{m}^\top \mathbf{u} + \tilde{\mathbf{v}}^\top \mathbf{m} \mathbb{E}[g(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})\mathbf{z}^\top] \mathbf{C}^{1/2} \mathbf{u} \\ &\quad + \tilde{\mathbf{v}}^\top \mathbf{C}^{1/2} \mathbb{E}[g(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})\mathbf{z}] \mathbf{m} \mathbf{u} + \tilde{\mathbf{v}}^\top \mathbf{C}^{1/2} \mathbb{E}[g(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})\mathbf{z}\mathbf{z}^\top] \mathbf{C}^{1/2} \mathbf{u} + O\left(\frac{1}{n^{1/4}}\right) \\ &= \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})]\tilde{\mathbf{v}}^\top (\mathbf{m}\mathbf{m}^\top + \mathbf{C}) \mathbf{u} \\ &\quad + \mathbb{E}[f'(\mathbf{w}^\top \mathbf{x})]\tilde{\mathbf{v}}^\top (\mathbf{C}\mathbf{w}\mathbf{m}^\top + \mathbf{m}\mathbf{C}\mathbf{w}^\top) \mathbf{u} \\ &\quad + \mathbb{E}[f''(\mathbf{w}^\top \mathbf{x})]\tilde{\mathbf{v}}^\top \mathbf{C}\mathbf{w}\mathbf{w}^\top \mathbf{C} \mathbf{u} + O\left(\frac{1}{n^{1/4}}\right) \end{aligned}$$

□

Let us now employ Proposition C to express the mappings (defined here for a random vector \mathbf{x}_i in the class \mathcal{C}_ℓ)

$$\begin{aligned} m_\ell(\mathbf{m}_W, \mathbf{C}_W) &= \mathbf{\Lambda}^{-1} \mathbb{E}[\tilde{\mathbf{x}}_i h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i})] \\ c_\ell(\mathbf{m}_W, \mathbf{C}_W) &= \frac{1}{n} \mathbf{\Lambda}^{-1} \mathbb{E}[\tilde{\mathbf{x}}_i h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i}) h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i})^\top \tilde{\mathbf{x}}_i^\top] \mathbf{\Lambda}^{-1}, \end{aligned}$$

from the class-wise means and covariances of x_1, \dots, x_n and \mathbf{m}_W and \mathbf{C}_W (that are respectively the mean and covariance of \mathbf{W} but also of \mathbf{W}_{-i}). We are going to fix successively \mathbf{x}_i and \mathbf{W}_{-i} to be able to compute the expectations appearing in the formulations of m_ℓ and c_ℓ (it is made possible since \mathbf{x}_i and \mathbf{W}_{-i} are independent). Although it is not fully rigorous, we employ Proposition C as if the estimations of $\mathbb{E}[f(\mathbf{w}^\top \mathbf{x}) \mathbf{v}^\top \mathbf{x}]$ and $\mathbb{E}[f(\mathbf{w}^\top \mathbf{x}) \mathbf{v}^\top \mathbf{x} \mathbf{x}^\top \mathbf{u}]$ for \mathbf{w}, \mathbf{v} and \mathbf{u} belonging to a big subset of \mathbb{S}^{p-1} of measure bigger than $1 - Ce^{-cp-1/4}$ implied that the result would be true for all vectors $\mathbf{w}, \mathbf{v}, \mathbf{u} \in \mathbb{S}^{p-1}$. It is of course not rigorously correct, however, in practice, for the vectors \mathbf{w}, \mathbf{v} and \mathbf{u} we are considering, it appears to be valid.

To simplify the expression of the derivative in $\tilde{\mathbf{x}}_i$ of $h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W})$, let us replace our couple of variables $(\tilde{\mathbf{x}}_i, \mathbf{W}) \in \mathbb{R}^{p \times k} \times \mathbb{R}^{p \times k}$ by the variables $(\mathbf{x}_i, \tilde{\mathbf{W}}) \in \mathbb{R}^p \times \mathbb{R}^{p \times k}$ where $\tilde{\mathbf{W}} = (\mathbf{w}_1, \dots, \mathbf{w}_k)$. We have then $h_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}) = h_i(\tilde{\mathbf{W}}^\top \mathbf{x}_i)$. Given a twice differentiable mapping $\phi : \mathbb{R}^k \rightarrow \mathbb{R}$ and $\psi : \mathbf{v} \mapsto \phi(\tilde{\mathbf{W}}^\top \mathbf{v})$ we have the identities:

$$\nabla \psi(\mathbf{v}) = \tilde{\mathbf{W}} \nabla \phi(\tilde{\mathbf{W}}^\top \mathbf{v}) \quad \text{and} \quad d^2|_{\mathbf{v}} = \tilde{\mathbf{W}} d^2 \phi|_{\tilde{\mathbf{W}}^\top \mathbf{v}} \tilde{\mathbf{W}}^\top$$

Now, let us follow our new notations and try to compute $m_\ell(\mathbf{m}_W, \mathbf{C}_W) \equiv \mathbb{E}_{-i} \left[\sqrt{n} \mathbb{E}_i[\mathbf{x}_i h^\ell(\tilde{\mathbf{W}}_{-i}^\top \mathbf{x}_i)] \right]$, where we noted $h^\ell : \mathbb{R}^k \rightarrow \mathbb{R}^k$, the mapping h_i for $k(i) = \ell$ (recall that $k(i)$ provides the class of \mathbf{x}_i). Let us decompose the matrix $\mathbf{S}_W = \mathbb{E}[\mathbf{W} \mathbf{W}^\top]$ followingly:

$$\mathbf{S}_W = \left(\begin{array}{c|c|c} \mathbf{S}_W^{1,1} & \dots & \mathbf{S}_W^{1,k} \\ \hline \vdots & & \vdots \\ \hline \mathbf{S}_W^{k,1} & \dots & \mathbf{S}_W^{k,k} \end{array} \right) \in \mathcal{M}_{kp},$$

where for all $a, b \in [k]$, $\mathbf{S}_W^{a,b} \in \mathcal{M}_p$, so that we can introduce the low-dimensional random vector $z \sim \mathcal{N}(\boldsymbol{\kappa}^\ell, \mathbf{K}^\ell)$ with:

$$\boldsymbol{\kappa}^\ell \equiv \boldsymbol{\mu}_\ell^\top \mathbf{m}_W \quad \text{and} \quad \mathbf{K}^\ell \equiv (\text{Tr}((\boldsymbol{\Sigma}_\ell + \boldsymbol{\mu}_\ell \boldsymbol{\mu}_\ell^\top) \mathbf{S}_W^{a,b}))_{1 \leq a, b \leq k} - \boldsymbol{\mu}_\ell^\top \mathbf{m}_W \mathbf{m}_W^\top \boldsymbol{\mu}_\ell. \quad (4)$$

With such a choice, z has the same distribution as $\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i}$ for $k(i) = \ell$. Then we can compute thanks to Proposition C:

$$\begin{aligned} m_\ell(\mathbf{m}_W, \mathbf{C}_W) &= \boldsymbol{\mu}_\ell \mathbb{E} \left[h^\ell(\tilde{\mathbf{W}}_{-i}^\top \mathbf{x}_i)^\top \right] + \boldsymbol{\Sigma}_\ell \mathbb{E}_{-i} \left[\tilde{\mathbf{W}}_{-i} \mathbb{E}_i[dh_i^\top | \tilde{\mathbf{W}}_{-i}^\top \mathbf{x}_i] \right] \\ &= \boldsymbol{\mu}_\ell \mathbb{E}[h^\ell(z)^\top] + \boldsymbol{\Sigma}_\ell \mathbf{m}_W \mathbb{E}[dh_i^\top |_z] + O_{\|\cdot\|} \left(\frac{1}{n^{1/4}} \right), \end{aligned}$$

since $\|\mathbf{C}_W\| \leq O(\sqrt{\log n/n})$ (because $\|\mathbf{W}\| \propto \mathcal{E}_2(\sqrt{\log n/n})$ on \mathcal{A}_X).

To estimate the mapping σ_ℓ , let us note for simplicity $H : \mathbf{v} \mapsto h^\ell(\tilde{\mathbf{W}}_{-i}^\top \mathbf{v}) h^\ell(\tilde{\mathbf{W}}_{-i}^\top \mathbf{v})^\top \in \mathcal{M}_{k,k}$, we know that: $\nabla H_{a,b}(\mathbf{v}) = \tilde{\mathbf{W}}_{-i} J(\tilde{\mathbf{W}}_{-i}^\top \mathbf{v})_{a,b}$ and $d^2 H_{a,b}|_{\mathbf{v}} = \tilde{\mathbf{W}}_{-i} K(\tilde{\mathbf{W}}_{-i}^\top \mathbf{v})_{a,b} \tilde{\mathbf{W}}_{-i}^\top$, where we introduce for any $\mathbf{u} \in \mathbb{R}^k$ the objects:

$$\begin{aligned} J(\mathbf{u})_{a,b} &= h_a^\ell(\mathbf{u}) \nabla h_b^\ell(\mathbf{u}) + h_b^\ell(\mathbf{u}) \nabla h_a^\ell(\mathbf{u}) \in \mathbb{R}^k \\ K(\mathbf{u})_{a,b} &= \nabla h_a^\ell(\mathbf{u}) \nabla h_b^\ell(\mathbf{u})^\top + h_b^\ell(\mathbf{u}) d^2 h_a^\ell|_{\mathbf{u}} + h_a^\ell(\mathbf{u}) d^2 h_b^\ell|_{\mathbf{u}} \in \mathcal{M}_k \end{aligned}$$

Following the same strategy as previously, we can show thanks to Proposition C that with the decomposition:

$$c_\ell(\mathbf{m}_W, \mathbf{C}_W) = \left(\begin{array}{c|c|c} c_\ell(\mathbf{m}_W, \mathbf{C}_W)_{1,1} & \dots & c_\ell(\mathbf{m}_W, \mathbf{C}_W)_{1,k} \\ \hline \vdots & & \vdots \\ \hline c_\ell(\mathbf{m}_W, \mathbf{C}_W)_{k,1} & \dots & c_\ell(\mathbf{m}_W, \mathbf{C}_W)_{k,k} \end{array} \right) \in \mathcal{M}_{kp},$$

for any $a, b \in [k]$ and $\ell \in [k]$ such that $k(i) = \ell$:

$$\begin{aligned}\Lambda c_\ell(\mathbf{m}_W, \mathbf{C}_W)_{a,b} \Lambda &= \frac{1}{n} \mathbb{E}[H_{a,b}(\mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top] \\ &= \frac{1}{n} \Sigma_\ell \tilde{\mathbf{m}}_W \mathbb{E}[J(z)_{a,b}] \boldsymbol{\mu}_\ell^\top + \frac{1}{n} \boldsymbol{\mu}_\ell \mathbb{E}[J(z)_{a,b}]^\top \tilde{\mathbf{m}}_W^\top \Sigma_\ell \\ &\quad + \frac{1}{n} \Sigma_\ell \left(\tilde{\mathbf{m}}_W \mathbb{E}[K(z)_{a,b}] \tilde{\mathbf{m}}_W^\top + \sum_{1 \leq c, d \leq k} S_W^{c,d} (\mathbb{E}[K(z)_{a,b}])_{c,d} \right) \Sigma_\ell \\ &\quad + \frac{1}{n} \mathbb{E}[h_a^\ell(z) h_b^\ell(z)] (\boldsymbol{\mu}_\ell \boldsymbol{\mu}_\ell^\top + \Sigma_\ell) + O_{\|\cdot\|} \left(\frac{1}{n^{1/4}} \right)\end{aligned}$$

We are then left to estimating the derivatives of h to be able to compute m_ℓ and c_ℓ . From the implicit expression of g_i given by Proposition 4.3 one can deduce the formulas, for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^k$:

$$dg_i|_{\mathbf{v}} = -(\Delta_i df_i|_{g(\mathbf{v})} - \mathbf{I}_k)^{-1} \quad \text{and} \quad d^2 g_i|_{\mathbf{v}} \cdot \mathbf{u} = dg_i|_{\mathbf{v}} \Delta_i (d^2 f_i|_{g(\mathbf{v})} \cdot \mathbf{u}) dg_i|_{\mathbf{v}}$$

Then, recalling the identity $h_i = f_i \circ g_i$, one can derive from the upper formulas the expression of the differentiates of h_i thanks to the identities for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^k$:

$$dh_i|_{\mathbf{v}} = \Delta_i^{-1} (dg_i|_{\mathbf{v}} - \mathbf{I}_k) \quad \text{and} \quad d^2 h_i|_{\mathbf{v}} \cdot \mathbf{u} = \Delta_i^{-1} d^2 g_i|_{\mathbf{v}} \cdot \mathbf{u}$$

4 Performance Evaluation

The expected misclassification error $E_t(\mathbf{x} \in \mathcal{C}_\ell)$ on a test data \mathbf{x} belonging to class $\ell \in [k]$ expresses followingly:

$$E_t(\mathbf{x} \in \mathcal{C}_\ell) = 1 - \mathbb{P}(\forall j \in [k] \setminus \{\ell\} : p_\ell(\mathbf{x}) \geq p_j(\mathbf{x})) \quad \text{where} \quad p_j(\mathbf{x}) = \frac{\phi(\mathbf{w}_j^\top \mathbf{x})}{\sum_{h=1}^k \phi(\mathbf{w}_h^\top \mathbf{x})}.$$

Since \mathbf{x} and \mathbf{w} are both concentrated and independent Theorem A allows us to assume that for all $j \in [k]$, $\tilde{\mathbf{x}}^\top \mathbf{W} = (\mathbf{w}_j^\top \mathbf{x})_{1 \leq j \leq k} \sim \mathcal{N}(\boldsymbol{\kappa}^\ell, \mathbf{K}^\ell)$ (the objects $\boldsymbol{\kappa}^\ell$ and \mathbf{K}^ℓ where introduced in equation 4). To simplify the problem, we are going to make the additional hypothesis that ϕ is increasing since in that case

$$\forall j \in [k] \setminus \{\ell\} : p_\ell(\mathbf{x}) \geq p_j(\mathbf{x}) \iff \mathbf{w}_\ell^\top \mathbf{x} \mathbf{1}_k - \tilde{\mathbf{x}}^\top \mathbf{W} \in \mathbb{R}_+^k$$

Since the ℓ^{th} coordinate of $\tilde{\mathbf{x}}^\top \mathbf{W}$ is, by definition, equal to $\mathbf{w}_\ell^\top \mathbf{x}$, only the $k-1$ other are interesting. Let us then introduce the Gaussian vector $\mathbf{Z}_\ell \in \mathbb{R}^{k-1}$ defined for all $j \in [k] \setminus \ell$ as: $[\mathbf{Z}_\ell]_j = (\mathbf{w}_\ell - \mathbf{w}_j)^\top \mathbf{x}$. Such a vector \mathbf{Z}_ℓ has then the mean $\bar{\boldsymbol{\kappa}} \equiv^\ell \mathbf{P} \boldsymbol{\kappa}^\ell$ and the covariance $\bar{\mathbf{K}} \equiv \mathbf{P} \mathbf{K}^\ell \mathbf{P}^\top$ with:

$$\mathbf{P} = \mathbf{1}_{k-1} \mathbf{e}_\ell^\top - \mathbf{I}_k^{-\ell} \in \mathcal{M}_{k-1,k}$$

where $\mathbf{1}_{k-1} \in \mathbb{R}^{k-1}$ is a vector full of 1, \mathbf{e}_ℓ is the ℓ^{th} vector of the canonical basis of \mathbb{R}^k (full of zeros with a one in the ℓ^{th} coordinate) and $\mathbf{I}_k^{-\ell}$ is the identity matrix of \mathcal{M}_k deprived of the ℓ^{th} row. We can then express

$$E_t(\mathbf{x} \in \mathcal{C}_\ell) = 1 - \mathbb{P}(\mathbf{Z}_\ell \in \mathbb{R}_+^{k-1})$$

5 Further experiments

In this section we provide further experiments using real images from the Imagenet dataset [DDS⁺09]. Figure 1 depicts the learned Softmax weights against their expected large p, n asymptotics as predicted by Theorem 3.2. As for GAN generated images, we observe a perfect match between the learned weights and the theoretical predictions. An almost perfect match is also observed for the scores (between the practical scores and their theoretical counterparts) as depicted in Figure 2 which strongly suggests that the conclusions of Theorem 3.2 generalize to real data. A Python implementation is attached to this Supplementary Material which implements Theorem 3.2 and Corollary 3.3 using the previous derivations.

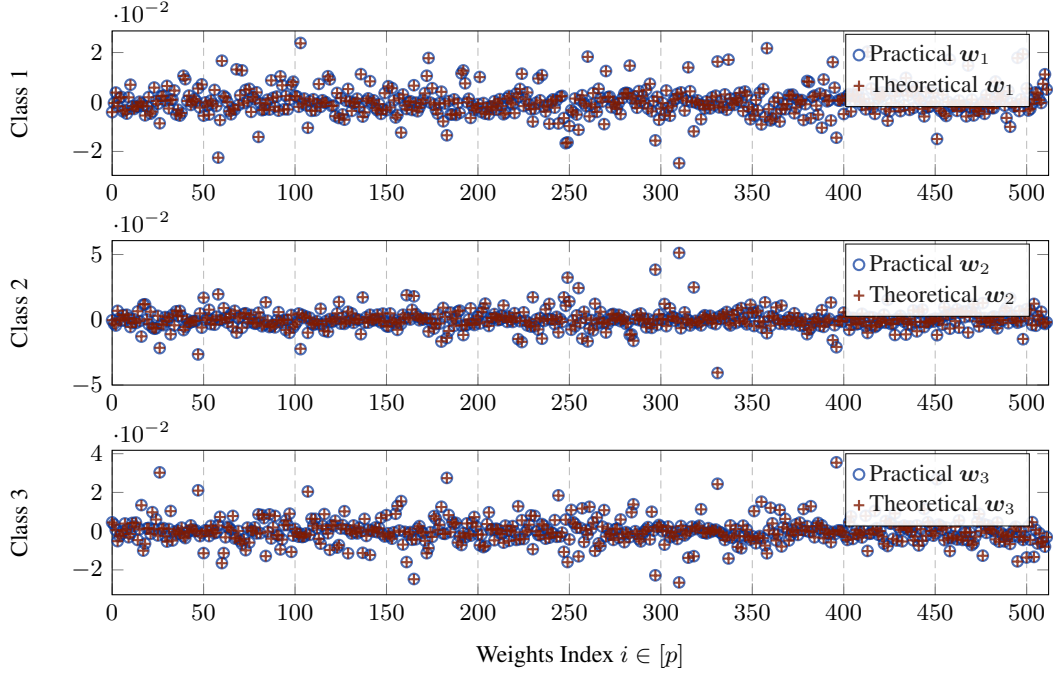


Figure 1: Learned weights (in blue circles) versus their theoretical estimations (in red crosses) as per Theorem 3.2 of the Main Paper. The used data are Resnet18 [SIVA17] representations ($p = 512$) of real images from the Imagenet dataset [DDS⁺09]. We considered $k = 3$ classes which are: *hamburger*, *mushroom*, *pizza*. $n = 3000$ and the regularization constants $\lambda_1, \lambda_2, \lambda_3 = 1.5$. The data are normalized such that their norm is $0.1 \cdot \sqrt{p}$ to ensure \mathcal{A}_X .

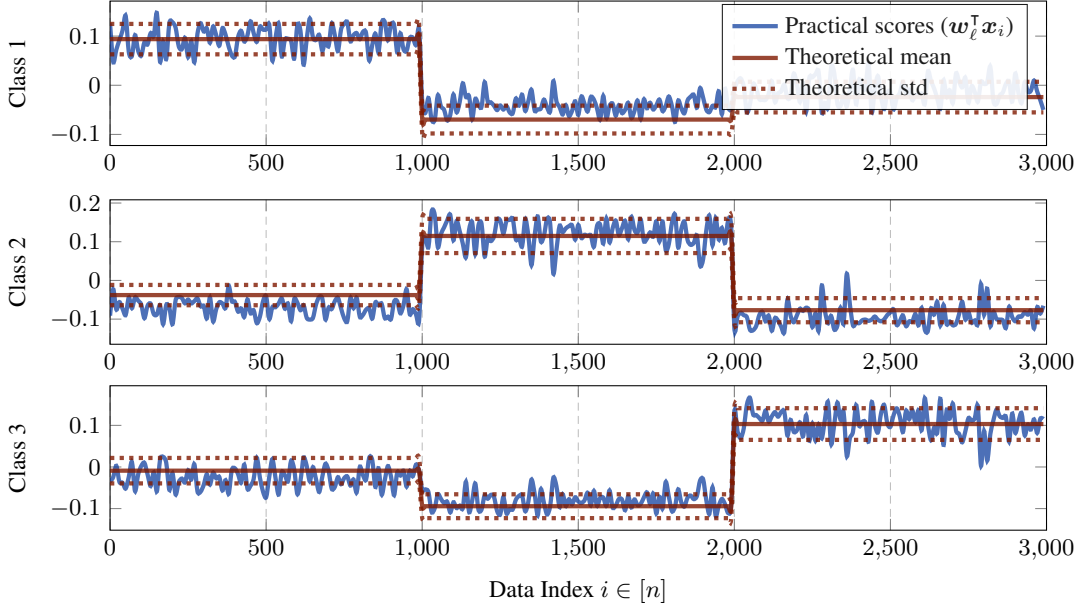


Figure 2: Scores (in blue) versus their theoretical estimations (in red) as per Corollary 3.3 of the Main Paper, with the theoretical means (through $\bar{\kappa}^\ell$) and standard deviations (through \bar{K}^ℓ), on a test set independent from the training set. The used data are Resnet18 [SIVA17] representations ($p = 512$) of real images from the Imagenet dataset [DDS⁺09]. We considered $k = 3$ classes which are: *hamburger*, *mushroom*, *pizza*. $n = 3000$ and the regularization constants $\lambda_1, \lambda_2, \lambda_3 = 1.5$. The data are normalized such that their norm is $0.1 \cdot \sqrt{p}$ to ensure \mathcal{A}_X .

References

- [DDS⁺09] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [FGP07] B. Fleury, O. Guédon, and G. Paouris. A stability result for mean width of l_p -centroid bodies. *Advances in Mathematics*, 214:865–877, 2007.
- [Fre19] Daniel J. Fresen. A simplified proof of clt for convex bodies. *arXiv preprint arXiv:1907.06785*, 2019.
- [Kla07] B. Klartag. A central limit theorem for convex sets. *Inventiones mathematicae volume*, 168:pages91–131, 2007.
- [LN08] Zinoviy Landsman and Johanna Nešlehová. Stein’s lemma for elliptical random vectors. *Journal of Multivariate Analysis*, 99(5):912–927, 2008.
- [SIVA17] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Thirty-first AAAI conference on artificial intelligence*, 2017.