

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

Mohamed El Amine Seddik¹ Cosme Louart^{2,3} Romain Couillet³ Mohamed Tamaazousti²
¹LIX, Ecole Polytechnique ²CEA List ³GIPSA-Lab Grenoble-Alps University

February 17, 2021

Outline: This supplementary material provides the essential derivations to obtain the results of the main paper (Section 1), the extension of these results to the Softmax classifier (Section 2) and further experiments supporting our findings (Section 3). Specifically, Subsection 1.1 provides justifications about the CLT result for concentrated vectors as per Theorem 3.2. Subsection 1.2 next presents the main ingredients for obtaining Theorem 4.8. Notations about the Softmax classifier are presented In Subsection 2.1, while Subsection 2.2 provides an analog of Theorem 4.8 for the general case of the Softmax classifier. Finally, further experiments using multi-class classification with synthetic Gaussian data and MNIST data as well as CNN representations of real images are shown in Section 3.

Notation: The notation $z \in \tilde{z} \pm \mathcal{E}_q(\sigma)$ stands for $z \in \mathbb{R}$ satisfying $z \propto \mathcal{E}_q(\sigma)$ and $|\tilde{z} - \mathbb{E}[z]| \leq \mathcal{O}(\sigma)$, where $|\cdot|$ replaces Euclidean norm for vectors. \otimes stands for the Kronecker product.

Reproducibility: A github link for the project will be provided.

1 Proofs of the main paper results

1.1 CLT for concentrated random vectors

In this section, we consider an asymptotic quantity $p \in \mathbb{N}$ that will represent the dimension of the random vectors we consider, and will be the quantity implicitly tending to infinity when we will employ the notations $\mathcal{O}(\sigma)$ of $\mathcal{E}_q(\sigma)$. Let us first provide two fundamental result that were already evoked in the main paper: the concentration of the Gaussian distribution $\gamma = \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ and the uniform distribution on the Sphere $\sqrt{p}\mathbb{S}^{p-1} \equiv \{\mathbf{x} \in \mathbb{R}^p, \|\mathbf{x}\| = \sqrt{p}\}$.

Theorem 1.1 (Concentration of the Gaussian distribution and the uniform measure on the sphere). *If $\mathbf{x} \sim \gamma$ or $\mathbf{x} \sim \text{Unif}(\sqrt{p}\mathbb{S}^{p-1})$, then $\mathbf{x} \propto \mathcal{E}_2$.*

Those two distributions are very similar and almost equal when p is high¹. In general one can imagine concentrated vectors \mathbf{x} as distributed along a sphere of center $\mathbb{E}[\mathbf{x}]$ and radius $\sqrt{\text{Tr}(\mathbf{C}_\mathbf{x})}$ where

$$\mathbf{C}_\mathbf{x} \equiv \mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^\top,$$

thanks to the following lemma that is a mere consequence of the 1-Lipschitz character of the Euclidean norm and whose full proof can be found in [Led05].

Lemma 1.2 (Concentrated vectors are distributed close to the sphere). *If $\mathbf{x} \propto \mathcal{E}_2$, then:*

$$\mathbb{P}\left(\left|\|\mathbf{x} - \mathbb{E}[\mathbf{x}]\| - \sqrt{\text{Tr}(\mathbf{C}_\mathbf{x})}\right|\right) \leq Ce^{-ct^2}.$$

We have the following result that can also be found in [LC20a], which controls the moments of concentrated random variables.

¹A famous result of Poincaré states indeed that any orthogonal projection of a random vector uniformly distributed on the sphere $\sqrt{p}\mathbb{S}^{p-1}$ tends to a Gaussian vector when p tends to infinity.

Proposition 1.3 (Characterization of the concentration with a bound on the moments). *A random variable z follows the concentration $z \in \tilde{z} \pm \mathcal{E}_q(\sigma)$ if and only if there exists two constants $C, c > 0$ such that for all $p \in \mathbb{N}$, for all $r > q$:*

$$\mathbb{E}[|z - \tilde{z}|^r] \leq C \left(\frac{r}{q}\right)^{\frac{r}{q}} (c\sigma)^r.$$

We provide below a proposition (similar to the result of [Kla07] provided in Theorem 3.2 in the main paper but more simple) supporting the fact that concentrated vectors share a common *concentration* behavior with Gaussian vectors. Specifically, given a random concentrated vector $\mathbf{x} \propto \mathcal{E}_2$, if one samples a vector \mathbf{u} from the uniform distribution on the unit sphere sphere² \mathbb{S}^{p-1} , then most likely, the random variable $\mathbf{u}^\top \mathbf{x}$ behaves asymptotically (as p grows) as a Gaussian random variable. Technically, this observation allows us to compute expectations of functionals of $\mathbf{w}^\top \mathbf{x}$ when \mathbf{x} is independent from \mathbf{w} (for instance in $\mathbf{w}_i^\top \mathbf{x}_i$). Note however, that even though \mathbf{w} is *a priori* not guaranteed to be uniformly distributed on the unit sphere, our experiments (in the main paper and subsequently) tend to validate the Gaussianity of the random variable $\mathbf{w}^\top \mathbf{x}$.

Proposition 1.4. *Given an integer $p \in \mathbb{N}$ and a random vector $\mathbf{x} \in \mathbb{R}^p$, if $\mathbf{x} \propto \mathcal{E}_2$ and $\|\mathbb{E}[\mathbf{x}]\| \leq \mathcal{O}(1)$, then for any parameter $\kappa \geq 1$ there exists two constants $c, C > 0$ and a subset $\Theta \subset \mathbb{R}^p$ such that $\gamma(\Theta) \geq 1 - 2e^{-c\kappa^2}$, $\gamma = \mathcal{N}(\mathbf{0}, \frac{1}{p}\mathbf{I}_p)$ and for any λ -Lipschitz mapping $f : \mathbb{R} \rightarrow \mathbb{R}$:*

$$\forall \mathbf{u} \in \Theta : |\mathbb{E}[f(\mathbf{u}^\top \mathbf{x})] - \mathbb{E}[f(z)]| \leq \frac{C\kappa\lambda}{\sqrt{p}},$$

where $z \sim \mathcal{N}(0, \frac{1}{p}\mathbb{E}[\|\mathbf{x}\|])$ or $z \sim \mathcal{N}(0, \frac{1}{p}\mathbb{E}[\text{Tr}(\mathbf{C})])$, with $\mathbf{C} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^\top$.

Remark 1.5. *If $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \frac{1}{p}\mathbf{I}_p)$, then $\boldsymbol{\theta} \propto \mathcal{E}_2(1/\sqrt{p})$ and $\mathbb{P}(\|\boldsymbol{\theta}\| - 1 \geq t) \leq 2e^{-pt^2/2}$. Therefore, considering a vector $\mathbf{u} \in \mathbb{R}^p$, if \mathbf{u} is far from the sphere \mathbb{S}^{p-1} , although it could possibly be in the set Θ , it is highly unprobable. Our heuristic thus requires to apply Proposition 1.4 with the vector $\mathbf{u}/\|\mathbf{u}\|$ and the mapping $t \mapsto f(\|\mathbf{u}\|t)$.*

Proof. The set Θ appears from a concentration result where the deterministic vector \mathbf{u} can be seen as a drawing of the random vector $\boldsymbol{\theta} \sim \gamma$ that we chose to be independent of \mathbf{x} . We already know that $\boldsymbol{\theta} \propto \mathcal{E}_2(1/\sqrt{p})$ and we know that $\phi : \mathbf{u} \mapsto \mathbb{E}[f(\mathbf{u}^\top \mathbf{x})]$ is λ' -Lipschitz with:

$$\lambda' \leq \lambda \sup_{\|\mathbf{v}\| \leq 1} \mathbb{E}[|\mathbf{v}^\top \mathbf{x}|] \leq \lambda(\|\mathbb{E}[\mathbf{v}^\top \mathbf{x}]\| + \mathcal{O}(1)) \leq \mathcal{O}(\lambda)$$

thanks to Proposition 1.3, since $\mathbf{v}^\top \mathbf{x} \in 0 \pm \mathcal{E}_2$. Therefore $\phi(\boldsymbol{\theta}) \in \mathbb{E}[\phi(\boldsymbol{\theta})] \pm \mathcal{E}_2(\lambda/\sqrt{p})$ thanks to the concentration of $\boldsymbol{\theta}$ and noting $\Theta = \{\mathbf{u} \in \mathbb{R}^p \mid \phi(\mathbf{u}) - \mathbb{E}[\phi(\boldsymbol{\theta})] \leq \frac{\lambda\kappa}{\sqrt{p}}\}$, we know that there exist two constants $C, c > 0$ such that $\mathbb{P}(\boldsymbol{\theta} \notin \Theta) \leq 2e^{-c\kappa^2}$.

Besides, noting $\mathbb{E}_{\mathbf{x}}$ and $\mathbb{E}_{\boldsymbol{\theta}}$, respectively, the expectation on \mathbf{x} (with $\boldsymbol{\theta}$ fixed) and the expectation on $\boldsymbol{\theta}$, we can estimate:

$$\mathbb{E}[\phi(\boldsymbol{\theta})] = \mathbb{E}_{\mathbf{x}}[\mathbb{E}_{\boldsymbol{\theta}}[f(\boldsymbol{\theta}^\top \mathbf{x})]].$$

For a fixed \mathbf{x} , $\boldsymbol{\theta}^\top \mathbf{x}$ is a Gaussian random variable with mean equal to 0 and variance equal to $\frac{1}{p}\|\mathbf{x}\|^2$ (since $\mathbb{E}[\boldsymbol{\theta}\boldsymbol{\theta}^\top] = \frac{1}{p}\mathbf{I}_p$). Introducing a random variable $z \sim \mathcal{N}(0, 1)$, independent of \mathbf{x} , we have thus the identity $\mathbb{E}_{\boldsymbol{\theta}}[f(\boldsymbol{\theta}^\top \mathbf{x})] = \mathbb{E}_z[f(\|\mathbf{x}\|z/\sqrt{p})]$. Finally, we have the following bound thanks to Proposition 1.3:

$$\left| \mathbb{E} \left[f \left(\frac{\|\mathbf{x}\|z}{\sqrt{p}} \right) \right] - \mathbb{E}_z \left[\frac{f(\mathbb{E}[\|\mathbf{x}\|]z)}{\sqrt{p}} \right] \right| \leq \frac{\lambda}{\sqrt{p}} \mathbb{E}[\|\mathbf{x}\| - \mathbb{E}[\|\mathbf{x}\|]] \mathbb{E}[|z|] \leq \mathcal{O} \left(\frac{\lambda}{\sqrt{p}} \right).$$

Therefore, choosing well the constant $C > 0$ provides the result of the proposition. \square

However, the above proposition needs to be adapted for estimating quantities of the form $\mathbb{E}[f(\mathbf{u}^\top \mathbf{x})\mathbf{v}^\top \mathbf{x}]$ or $\mathbb{E}[f(\mathbf{u}^\top \mathbf{x})\mathbf{v}^\top \mathbf{x}\mathbf{x}^\top \mathbf{w}]$ for some deterministic vectors $\mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathbb{R}^p$, as we will need subsequently. In particular, we have the following result.

²Which is equivalent to the Gaussian distribution $\mathcal{N}(0, \frac{\mathbf{I}_p}{n})$.

Proposition 1.6. *In the setting of Proposition 1.4, let us note $\nu = \text{Unif}(\mathbb{S}^{p-1})$, for any parameter $\kappa > 0$, there exists two constants $C, c > 0$ and a subset $\Theta_2 \subset (\mathbb{R}^p)^2$ such that³ $\nu^{\times 2}(\Theta_2) \geq 1 - 2e^{-c\kappa^2}$ and for any λ -Lipschitz mapping $f : \mathbb{R} \rightarrow \mathbb{R}$ such that $|\mathbb{E}[f(\mathbf{u}^\top \mathbf{x})]| \leq \lambda$, we have*

$$\forall (\mathbf{u}, \mathbf{v}) \in \Theta_2 : |\mathbb{E}[f(\mathbf{u}^\top \mathbf{x})\mathbf{v}^\top \mathbf{x}]| \leq \frac{C\kappa\lambda}{\sqrt{p}}.$$

Remark 1.7. *Given two independent random vectors $\boldsymbol{\theta}, \boldsymbol{\sigma} \sim \mathcal{N}(0, \frac{\mathbf{I}_p}{p})$, the couple $(\boldsymbol{\theta}, \boldsymbol{\sigma}) \in (\mathbb{R}^p)^2$ has the distribution $\gamma^{\otimes 2}$ and one can show that*

$$\mathbb{P}(|\boldsymbol{\theta}^\top \boldsymbol{\sigma}| \geq t) \leq C'e^{-c'pt^2} + C'e^{-c'pt},$$

for some constants $C', c' > 0$ (independent of p). The same result particularly holds for $(\boldsymbol{\theta}, \boldsymbol{\sigma}) \sim \nu^{\times 2}$. Thus, we see that for sufficiently large κ , the couples of Θ_2 have high probability to be quasi-orthogonal. Hence, if one considers two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^p$, our heuristic requires not to employ directly Proposition 1.6 as if $(\mathbf{u}, \mathbf{v}) \in \Theta_2$, but rather employ it with the couple $(\mathbf{u}, \mathbf{v} - (\mathbf{u}^\top \mathbf{v})\mathbf{u})$, we end up therefore having the identity

$$\mathbb{E}[f(\mathbf{u}^\top \mathbf{x})\mathbf{v}^\top \mathbf{x}] = \mathbf{u}^\top \mathbf{v} \mathbb{E}[f(\mathbf{u}^\top \mathbf{x})\mathbf{u}^\top \mathbf{x}] + \mathcal{O}\left(\frac{\kappa\lambda}{\sqrt{p}}\right),$$

where the right hand side expectation can be estimated by Proposition 1.4, through the mapping $t \mapsto f(t)t$.

Proof. This proof follows the same steps as the proof of Proposition 1.4. Considering $(\boldsymbol{\theta}, \boldsymbol{\sigma}) \sim \nu^{\times 2}$, we know that $(\boldsymbol{\theta}, \boldsymbol{\sigma}) \propto \mathcal{E}_2(1/\sqrt{p})$ ($\nu^{\times 2}$ is similar to $\mathcal{N}(\mathbf{0}, \mathbf{I}_{2p}/2p)$), and the mapping $(\mathbf{u}, \mathbf{v}) \mapsto \mathbb{E}[f(\mathbf{u}^\top \mathbf{x})\mathbf{v}^\top \mathbf{x}]$ is λ' -Lipschitz for the Euclidean norm of $(\mathbb{R}^p)^2$, where we can bound thanks to Proposition 1.3 and Hölder inequality:

$$\begin{aligned} \lambda' &= \sqrt{\|\mathbb{E}[\lambda \mathbf{v}^\top \mathbf{x} \mathbf{x}]\|^2 + \|\mathbb{E}[f(\mathbf{u}^\top \mathbf{x})\mathbf{x}]\|^2} \\ &\leq \sqrt{\sup_{\|\mathbf{a}\| \leq 1} \lambda^2 |\mathbb{E}[(\mathbf{v}^\top \mathbf{x})^2]\mathbb{E}[(\mathbf{a}^\top \mathbf{x})^2]| + |\mathbb{E}[f(\mathbf{u}^\top \mathbf{x})^2]\mathbb{E}[(\mathbf{a}^\top \mathbf{x})^2]|} \leq \mathcal{O}(\lambda). \end{aligned}$$

Therefore, thanks again to Proposition 1.3, if we note:

$$\Theta_2 = \left\{ (\mathbf{u}, \mathbf{v}) \in \mathbb{R}^p \mid |\mathbb{E}[f(\mathbf{u}^\top \mathbf{x})\mathbf{v}^\top \mathbf{x}] - \mathbb{E}[f(\boldsymbol{\theta}^\top \mathbf{x})\boldsymbol{\sigma}^\top \mathbf{x}]| \leq \frac{\kappa\lambda}{\sqrt{p}} \right\},$$

then, there exist two constants $C, c > 0$ such that $\mathbb{P}((\boldsymbol{\theta}, \boldsymbol{\sigma}) \in \Theta_2) \leq Ce^{-c\kappa^2}$ and we retrieve the result of the proposition since (with the same notations as in the proof of Proposition 1.4):

$$\mathbb{E}[f(\boldsymbol{\theta}^\top \mathbf{x})\boldsymbol{\sigma}^\top \mathbf{x}] = \mathbb{E}_{\mathbf{x}} [\mathbb{E}_{\boldsymbol{\theta}}[f(\boldsymbol{\theta}^\top \mathbf{x})]\mathbb{E}_{\boldsymbol{\sigma}}[\boldsymbol{\sigma}^\top \mathbf{x}]] = 0.$$

□

1.2 Estimation of the weight statistics

The complete justifications leading to Theorem 4.8 are quite long and already provided in [LC20b] in the Gaussian setting. We thus only provide here the main ingredients to have some intuitions on the results. Typically, Propositions 1.4 and 1.6 will allow us to employ Stein-like identities as the ones given in the up coming proposition. They are specifically provided for Gaussian vectors \mathbf{x} , but, as we saw from the previous subsection, they can be extended to any concentrated vector, as long as $\mathbf{x} \propto \mathcal{E}_2$ and $\|\mathbb{E}[\mathbf{x}]\| \leq \mathcal{O}(1)$.

Proposition 1.8 (Stein identities). *Given $\mathbf{x} \in \mathbb{R}^p$, a Gaussian random vector satisfying $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ and, $f : \mathbb{R}^p \rightarrow \mathbb{R}$ some three times differentiable function such that f, f' and f'' are all λ -Lipschitz with $\lambda \leq \mathcal{O}(1)$, then for any $\mathbf{v}, \mathbf{w} \in \mathbb{R}^p$ and any $\mathbf{A} \in \mathcal{M}_p$, we have*

$$\begin{aligned} \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})\mathbf{v}^\top \mathbf{x}] &= \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})]\mathbf{v}^\top \boldsymbol{\mu} + \mathbb{E}[f'(\mathbf{w}^\top \mathbf{x})]\mathbf{v}^\top \mathbf{C}\mathbf{w}, \\ \mathbb{E}[f(\mathbf{w}^\top \mathbf{x}) \text{Tr}(\mathbf{x}\mathbf{x}^\top \mathbf{A})] &= \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})] \text{Tr}(\mathbf{A}(\boldsymbol{\mu}\boldsymbol{\mu}^\top + \mathbf{C})) + \mathbb{E}[f'(\mathbf{w}^\top \mathbf{x})]\boldsymbol{\mu}^\top (\mathbf{A} + \mathbf{A}^\top)\mathbf{C}\mathbf{w} + \mathbb{E}[f''(\mathbf{w}^\top \mathbf{x})]\mathbf{w}^\top \mathbf{C}\mathbf{A}\mathbf{C}\mathbf{w}. \end{aligned}$$

³ $\nu^{\times 2} = \nu \times \nu$ stands for the measure product.

Proof. The proposition is a straightforward result of applying the Stein's identity; for $z \sim \mathcal{N}(0, \sigma^2)$ and $f : \mathbb{R} \rightarrow \mathbb{R}$ differentiable, $\mathbb{E}[f(z)z] = \sigma^2[f'(z)]$. \square

Remark 1.9. In the case where \mathbf{x} is not Gaussian but rather only concentrated, such that $\mathbf{x} \propto \mathcal{E}_2$ and $\|\boldsymbol{\mu}\| \leq \mathcal{O}(1)$ ⁴, the formulas in Proposition 1.8 become estimations with a vanishing error of order $\mathcal{O}(n^{-\frac{1}{2}})$ ⁵ when $\|\mathbf{v}\| \leq \mathcal{O}(1)$ and $\|\mathbf{A}\| \leq \mathcal{O}(1/n)$. In particular, normalizing the second inequality with n and considering $\mathbf{A} \in \mathcal{M}_p$ such that $\|\mathbf{A}\| \leq 1$, we have

$$\frac{1}{n} \mathbb{E}[f(\mathbf{w}^\top \mathbf{x}) \mathbf{x}^\top \mathbf{A} \mathbf{x}] = \frac{1}{n} \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})] \text{Tr}(\mathbf{A} \mathbf{C}) + \mathcal{O}\left(n^{-\frac{1}{2}}\right),$$

since $\frac{1}{n} \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu} \leq \mathcal{O}(\frac{1}{n})$, $|\frac{1}{n} \boldsymbol{\mu}^\top (\mathbf{A} + \mathbf{A}^\top) \mathbf{C} \mathbf{w}| \leq \mathcal{O}(\frac{1}{n})$ and $|\frac{1}{n} \mathbb{E}[f''(\mathbf{w}^\top \mathbf{x})] \mathbf{w}^\top \mathbf{C} \mathbf{A} \mathbf{C} \mathbf{w}| \leq \mathcal{O}(\frac{1}{n})$.

We turn now to the estimation of the weight statistics, i.e., the quantities $\boldsymbol{\mu}_w$ and \mathbf{C}_w . Starting from the following estimation

$$\left\| \boldsymbol{\mu}_w - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathcal{A}_w} [\xi_{\delta_{k(i)}} (\mathbf{x}_i^\top \mathbf{w}_{-i}) \mathbf{x}_i] \right\| \leq \mathcal{O}(n^{-\frac{1}{2}}), \quad (1)$$

we apply Proposition 1.8 (more precisely Remark 1.9) which yields to; for most of vectors $\mathbf{v} \in \mathbb{R}^p$, such that $\|\mathbf{v}\| \leq \mathcal{O}(1)$,

$$\mathbb{E}_{\mathcal{A}_w} [\xi_{\delta_{k(i)}} (\mathbf{x}_i^\top \mathbf{w}_{-i}) \mathbf{v}^\top \mathbf{x}_i] = \mathbb{E}[\xi_{\delta_{k(i)}} (\tilde{z}_{k(i)})] \mathbf{v}^\top \boldsymbol{\mu}_{k(i)} + \mathbb{E}[\xi'_{\delta_{k(i)}} (\tilde{z}_{k(i)})] \mathbf{v}^\top \mathbf{C}_{k(i)} \mathbb{E}_{\mathcal{A}_w} [\mathbf{w}_{-i}] + \mathcal{O}\left(n^{-\frac{1}{2}}\right),$$

where, this time, $\tilde{z}_{k(i)}$ is a Gaussian random variable having the same mean and variance as $z_i \equiv \mathbf{x}_i^\top \mathbf{w}_{-i}$. Moreover, to show that \mathbf{w}_{-i} has the same statistics as \mathbf{w} , we need a supplementary result (Theorem 1.10) from [LC20b], which in turn leads to Theorem 4.4.

Theorem 1.10. Under Assumption 1-4, for any $i \in [n]$: $\|\mathbf{w} - \mathbf{w}_{-i} - \frac{1}{n} \mathbf{Q}_{-i} \mathbf{x}_i f(\mathbf{x}_i^\top \mathbf{w})\| \in 0 \pm \mathcal{E}_2(n^{-\frac{1}{2}})$.

Thus, since for any deterministic vector $\mathbf{u} \in \mathbb{R}^p$ such that $\|\mathbf{u}\| \leq \mathcal{O}(1)$, one can show that $\frac{1}{n} \mathbf{u}^\top \mathbf{Q}_{-i} \mathbf{x}_i \in 0 \pm \mathcal{E}_2(n^{-\frac{1}{2}})$, one can then deduce that $\mathbf{u}^\top \mathbf{w} - \mathbf{u}^\top \mathbf{w}_{-i} \in 0 \pm \mathcal{E}_2(n^{-\frac{1}{2}})$, which implies in particular that $\|\mathbb{E}_{\mathcal{A}_w} [\mathbf{w}_{-i}] - \boldsymbol{\mu}_w\| \leq \mathcal{O}(n^{-\frac{1}{2}})$. And studying now functionals of the form $\mathbf{w}^\top \mathbf{A} \mathbf{w}$, with $\mathbf{A} \in \mathcal{M}_p$ such that $\|\mathbf{A}\| \leq \mathcal{O}(1)$, one can also show that $\|\mathbb{E}_{\mathcal{A}_w} [\mathbf{w}_{-i} \mathbf{w}_{-i}^\top] - \mathbf{C}_w\| \leq \mathcal{O}(n^{-\frac{1}{2}})$.

Hence, setting $\mathbf{K} \equiv \sum_{\ell=1}^k \gamma_\ell \mathbb{E}[\xi'_{\delta_\ell}(\tilde{z}_\ell)] \mathbf{C}_\ell$ and $\tilde{\boldsymbol{\mu}} \equiv \sum_{\ell=1}^k \gamma_\ell \mathbb{E}[\xi_{\delta_\ell}(\tilde{z}_\ell)] \boldsymbol{\mu}_\ell$, we obtain from (1),

$$\|\boldsymbol{\mu}_w - \tilde{\boldsymbol{\mu}} - \mathbf{K} \boldsymbol{\mu}_w\| \leq \mathcal{O}\left(n^{-\frac{1}{2}}\right). \quad (2)$$

Now looking for an estimation of \mathbf{C}_w , we start from the identity:

$$\left\| \mathbf{C}_w + \boldsymbol{\mu}_w \boldsymbol{\mu}_w^\top - \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}_{\mathcal{A}_w} [\xi_{\delta_{k(i)}} (\mathbf{x}_i^\top \mathbf{w}_{-i})^2 \mathbf{x}_i \mathbf{x}_i^\top] - \frac{1}{n^2} \sum_{\substack{1 \leq i, j \leq n \\ i \neq j}} \mathbb{E}_{\mathcal{A}_w} [\xi_{\delta_{k(i)}} (\mathbf{x}_i^\top \mathbf{w}_{-i}) \xi_{\delta_{k(j)}} (\mathbf{x}_j^\top \mathbf{w}_{-j}) \mathbf{x}_i \mathbf{x}_j^\top] \right\|_* \leq \mathcal{O}(n^{-\frac{1}{2}}). \quad (3)$$

Then, we first deduce as per Remark 1.9 that for any $\ell \in [k]$:

$$\left\| \mathbb{E}_{\mathcal{A}_w} [\xi_{\delta_{k(i)}} (\mathbf{x}_i^\top \mathbf{w}_{-i})^2 \mathbf{x}_i \mathbf{x}_i^\top] - \mathbb{E}_{\mathcal{A}_w} [\xi_{\delta_{k(i)}} (\tilde{z}_{k(i)})^2] \mathbf{C}_{k(i)} \right\| \leq \mathcal{O}(n^{-\frac{1}{2}}).$$

For any $i, j \in [n]$, $i \neq j$, let us define $\mathbf{w}_{-i, j}$ as the unique solution $\mathbf{w}_{-i, j} = \frac{1}{n} \mathbf{X}_{-i, j} f(\mathbf{X}_{-i, j}^\top \mathbf{w}_{-i, j})$, where $\mathbf{X}_{-i, j}$ is the matrix \mathbf{X}_{-i} with a zero vector in the j^{th} column. Then we can show from Theorem 1.10 that

⁴The concentration $\mathbf{x} \propto \mathcal{E}_2$ automatically implies that $\|\mathbf{C}\| \leq \mathcal{O}(1)$. Indeed for all $\mathbf{v} \in \mathbb{R}^p$, $\mathbf{v}^\top \mathbf{x} \in \mathbf{v}^\top \boldsymbol{\mu} \pm \mathcal{E}_2$ and thus, Proposition 1.3 implies $\|\mathbf{C}\| = \sup_{\|\mathbf{v}\| \leq 1} \mathbf{v}^\top \mathbf{C} \mathbf{v} = \sup_{\|\mathbf{v}\| \leq 1} \mathbb{E}[(\mathbf{v}^\top \mathbf{x} - \mathbf{v}^\top \boldsymbol{\mu})^2] \leq \mathcal{O}(1)$.

⁵Recall that p (the dimension of the data) and n (the number of data) are of the same order i.e $p = \mathcal{O}(n)$ and $n = \mathcal{O}(p)$.

$\mathbf{x}_i^\top \mathbf{w}_{-i} - \mathbf{x}_i^\top \mathbf{w}_{-i,j} \in 0 \pm \mathcal{E}_2(1/\sqrt{n})$ (since \mathbf{x}_i is mutually independent with \mathbf{w}_{-i} and $\mathbf{w}_{-i,j}$), thus

$$\begin{aligned} & \mathbb{E}_{\mathcal{A}_w} [\xi_{\delta_{k(i)}}(\mathbf{x}_i^\top \mathbf{w}_{-i}) \xi_{\delta_{k(j)}}(\mathbf{x}_j^\top \mathbf{w}_{-j}) \mathbf{x}_i^\top \mathbf{A} \mathbf{x}_j] \\ &= \mathbb{E}_{\mathcal{A}_w} [\xi_{\delta_{k(i)}}(\mathbf{x}_i^\top \mathbf{w}_{-i,j}) \xi_{\delta_{k(j)}}(\mathbf{x}_j^\top \mathbf{w}_{-i,j}) \mathbf{x}_i^\top \mathbf{A} \mathbf{x}_j] + \mathcal{O}(n^{-\frac{1}{2}}) \\ &= \left(\mathbb{E}[\xi_{\delta_{k(i)}}(\tilde{z}_{k(i)})] \boldsymbol{\mu}_{k(i)} + \mathbb{E}[\xi'_{\delta_{k(i)}}(\tilde{z}_{k(i)})] \mathbf{C}_{k(i)} \boldsymbol{\mu}_w \right)^\top \mathbf{A} \left(\mathbb{E}[\xi_{\delta_{k(j)}}(\tilde{z}_{k(j)})] \boldsymbol{\mu}_{k(j)} + \mathbb{E}[\xi'_{\delta_{k(j)}}(\tilde{z}_{k(j)})] \mathbf{C}_{k(j)} \boldsymbol{\mu}_w \right) \\ & \quad + \mathbb{E}[\xi'_{\delta_{k(i)}}(\tilde{z}_{k(i)})] \mathbb{E}[\xi'_{\delta_{k(j)}}(\tilde{z}_{k(j)})] \text{Tr}(\mathbf{C}_{k(i)} \mathbf{A} \mathbf{C}_{k(j)} \mathbf{C}_w) + \mathcal{O}(n^{-\frac{1}{2}}). \end{aligned}$$

Therefore, summing up these estimations for $1 \leq i, j \leq n$ and $i \neq j$, we obtain:

$$\frac{1}{n^2} \sum_{\substack{1 \leq i, j \leq n \\ i \neq j}} \mathbb{E}_{\mathcal{A}_w} [\xi_{\delta_{k(i)}}(\mathbf{x}_i^\top \mathbf{w}_{-i}) \xi_{\delta_{k(j)}}(\mathbf{x}_j^\top \mathbf{w}_{-j}) \mathbf{x}_i^\top \mathbf{A} \mathbf{x}_j] = \boldsymbol{\mu}_w^\top \mathbf{A} \boldsymbol{\mu}_w + \frac{1}{n} \text{Tr}(\mathbf{C}_w \mathbf{K} \mathbf{A} \mathbf{K}) + \mathcal{O}(n^{-\frac{1}{2}}).$$

Moreover, if we note $\tilde{\mathbf{C}} \equiv \sum_{\ell=1}^k \gamma_\ell \mathbb{E}_{\mathcal{A}_w} [\xi_{\delta_\ell}(\tilde{z}_\ell)^2] \mathbf{C}_\ell$, equation (3) implies that,

$$\left\| \mathbf{C}_w - \tilde{\mathbf{C}} - \mathbf{K} \mathbf{C}_w \mathbf{K} \right\|_* \leq \mathcal{O}(n^{-\frac{1}{2}}). \quad (4)$$

Since $\|\mathbf{K}\| < 1 - \varepsilon'$ (see [LC20b] for details), for some constant $\varepsilon' > 0$, the matrix $\mathbf{R}_1 = (\mathbf{I}_p - \mathbf{K})^{-1}$ and the linear mapping \mathbf{R}_2 introduced in the main paper are well defined, and we have from (2) and (4), the estimations,

$$\left\| \boldsymbol{\mu}_w - \mathbf{R}_1 \tilde{\boldsymbol{\mu}} \right\| \leq \mathcal{O}(n^{-\frac{1}{2}}), \quad \left\| \mathbf{C}_w - \mathbf{R}_2(\tilde{\mathbf{C}}) \right\|_* \leq \mathcal{O}(n^{-\frac{1}{2}}).$$

We can then naturally approximate the quantities $\delta_\ell = \frac{1}{n} \text{Tr}(\mathbf{C}_\ell (\mathbf{I}_p - \mathbf{K})^{-1})$, $m_\ell = \mathbb{E}[\tilde{z}_\ell]$ and $\sigma_\ell = \mathbb{E}[\tilde{z}_\ell^2] - m_\ell^2$ as the solutions to the fixed point equation presented in Theorem 4.8 of the main paper. Note however that the proof of the existence and uniqueness of the solutions is quite an elaborate problem, but the simulations strongly confirm this conjecture.

2 Extension to the Softmax classifier

2.1 Position of the Problem and first properties

The main encountered difficulty with the Softmax classifier setting is that the mapping f goes from \mathbb{R}^k to \mathbb{R}^k instead of being scalar, thus its derivative is now a differential which brings some purely formal complexities. For any $\ell \in [k]$, let us denote

$$\begin{aligned} f_\ell : \mathbb{R}^k & \longrightarrow \mathbb{R}^k \\ \mathbf{v} & \longmapsto \left[\frac{\phi(v_a)/\lambda_a}{\sum_{b=1}^k \phi(v_b)} \sum_{b=1}^k \tilde{Y}_{\ell,b} \psi(v_b) - \tilde{Y}_{\ell,a} \psi(v_a) \right]_{1 \leq a \leq k}, \end{aligned} \quad (5)$$

where $\tilde{\mathbf{Y}} = \mathbf{I}_k$ contains in the ℓ^{th} column the label of the class ℓ . Then, the fixed point equation satisfied by the Softmax weights $\mathbf{W} = (\mathbf{w}_1^\top, \dots, \mathbf{w}_k^\top)^\top \in \mathbb{R}^{pk}$ (all stacked in a vector) writes in the simpler form,

$$\mathbf{W} = \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{x}}_i f_{k(i)}(\tilde{\mathbf{x}}_i^\top \mathbf{W}),$$

where $\tilde{\mathbf{x}}_i \equiv \mathbf{I}_k \otimes \mathbf{x}_i = \begin{pmatrix} \mathbf{x}_i & & \\ & \ddots & \\ & & \mathbf{x}_i \end{pmatrix} \in \mathcal{M}_{pk,k} \forall i \in [n]$. Let us introduce a constant $\varepsilon > 0$ and define the event \mathcal{A}_W as,

$$\mathcal{A}_W \equiv \left\{ \sup_{\substack{1 \leq \ell \leq k \\ \mathbf{z} \in \mathbb{R}^k}} \frac{\|df_\ell|_{\mathbf{z}}\|}{n} \|\tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top\| \leq 1 - \varepsilon \right\},$$

where $\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n) \in \mathcal{M}_{pk,nk}$. Then, replacing Assumptions 3 by its analog as follows,

Assumption 3 bis. $\sup_{\substack{1 \leq \ell \leq k \\ \mathbf{z} \in \mathbb{R}^k}} \frac{\|df_\ell|_{\mathbf{z}}\|}{n} \mathbb{E}[\|\tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top\|] \leq 1 - 2\varepsilon,$

we have the following analog to Lemma 4.2:

Lemma 2.1. *Under Assumption 1, 2 and 3 bis, there exists two constants $C, c > 0$ such that: $\mathbb{P}(\mathcal{A}_{\mathbf{W}}^c) \leq Ce^{-cn}.$*

Proof. It is a consequence of the result of concentration of product of concentrated variables provided in [LC20a] and that sets that since $\|\tilde{\mathbf{X}}\|/\sqrt{n} \propto \mathcal{E}_2(n^{-\frac{1}{2}})$ and $\mathbb{E}[\|\tilde{\mathbf{X}}\|/\sqrt{n}] \leq \mathcal{O}(1)$, then $\frac{1}{n}\|\tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top\| = \frac{1}{n}\|\tilde{\mathbf{X}}\|^2 \propto \mathcal{E}_2(1/\sqrt{n}) + \mathcal{E}_1(1/n)$. This signifies that there exist two constants $C, c > 0$ such that for any $p \in \mathbb{N}$:

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

Thus, in particular, Assumption 3 bis allows us to bound:

$$\mathbb{P}(\mathcal{A}_{\mathbf{W}}^c) \leq \mathbb{P}\left(\left|\frac{1}{n}\|\tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top\| - \mathbb{E}\left[\frac{1}{n}\|\tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top\|\right]\right| \geq \varepsilon\right) \leq Ce^{-cn\varepsilon^2} + Ce^{-cn\varepsilon},$$

which gives us the result of the Lemma, modifying slightly the constants $C, c > 0$. \square

Hence, in the case of the multiclass softmax classification, we obtain also the concentration of \mathbf{W} .

Theorem 2.2. *Under Assumption 1, 2 and 3 bis: $(\mathbf{W} \mid \mathcal{A}_{\mathbf{W}}) \propto \mathcal{E}_2(n^{-\frac{1}{2}}).$*

2.2 Analog to Theorem 4.8

Basically, the fixed point equation of Theorem 4.8 remains the same in the (multi-class) Softmax setting, we rather just need to adapt the objects $\mathbf{Q}, \delta, \xi_\ell, \tilde{\mathbf{C}}, \tilde{\boldsymbol{\mu}}, \mathbf{K}, \mathbf{R}_1, \mathbf{R}_2$ to the vectorial mapping introduced in equation (5). Hence, we are looking at first to an analog of Theorem 4.4 to link \mathbf{W} with the random matrix \mathbf{W}_{-i} defined as the only solution to $\mathbf{W}_{-i} = \frac{1}{n} \sum_{\substack{1 \leq j \leq n \\ j \neq i}} \tilde{\mathbf{x}}_j f_{k(j)}(\tilde{\mathbf{x}}_j^\top \mathbf{W}_{-i})$. For that an adaptation of Assumption 4 is needed.

Assumption 4 bis. $\forall \ell \in [k], \sup_{\mathbf{z} \in \mathbb{R}^k} \|df_\ell^2|_{\mathbf{z}}\| \leq \infty^6.$

Then introducing the notations $\tilde{\mathbf{X}}_{-i} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{i-1}, \mathbf{0}, \tilde{\mathbf{x}}_{i+1}, \dots, \tilde{\mathbf{x}}_n) \in \mathcal{M}_{pk, nk}$ and:

$$\mathbf{Q}_{-i} \equiv \left(\mathbf{I}_{kp} - \frac{1}{n} \tilde{\mathbf{X}}_{-i} \mathbf{D}^{(i)} \tilde{\mathbf{X}}_{-i}^\top \right)^{-1} \in \mathcal{M}_{kp}, \quad \mathbf{D}^{(i)} = \text{Diag}(\mathbf{D}_j^{(i)}) \in \mathcal{M}_{kn}, \quad \mathbf{D}_j^{(i)} \equiv df_{k(j)}|_{\tilde{\mathbf{x}}_j^\top \mathbf{W}_{-i}} \in \mathcal{M}_k.$$

We have the following proposition which is at the core of the Softmax classifier analysis.

Proposition 2.3 (Analog to Theorem 4.4). *Under Assumptions 1, 2, 3 bis and 4 bis, for any $i \in [n]$:*

$$\forall t > 0 : \mathbb{P}_{\mathcal{A}_X} \left(\left\| \tilde{\mathbf{x}}_i^\top \mathbf{W} - \tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i} + \frac{1}{n} \tilde{\mathbf{x}}_i^\top \mathbf{Q}_{-i} \tilde{\mathbf{x}}_i f_{k(i)}(\tilde{\mathbf{x}}_i^\top \mathbf{W}) \right\| \geq t \right) \leq Ce^{-cnt^2}.$$

Consequently, we retrieve an analog to Theorem 4.8 for the Softmax classifier.

Theorem 2.4 (Analog to Theorem 4.8). *Under Assumptions 1, 2, 3 bis and 4 bis, there exists three unique tuples $(\mathbf{m}_\ell)_{1 \leq \ell \leq k} \in (\mathbb{R}^k)^{[k]}$, $(\boldsymbol{\sigma}_\ell)_{1 \leq \ell \leq k} \in (\mathcal{M}_k)^{[k]}$ and $(\boldsymbol{\Delta}_\ell)_{1 \leq \ell \leq k} \in (\mathcal{M}_k)^{[k]}$ satisfying the conditions $(\forall \ell \in [k])$:*

- $\tilde{\mathbf{z}}_\ell \sim \mathcal{N}(\mathbf{m}_\ell, \boldsymbol{\sigma}_\ell^2);$
- $\xi_\ell : \mathbb{R}^k \rightarrow \mathbb{R}^k$ satisfying for any $\mathbf{z} \in \mathbb{R}^k$:

$$\xi_\ell(\mathbf{z}) = f_\ell(\mathbf{z} + \boldsymbol{\Delta}_\ell \xi_\ell(\mathbf{z}));$$
- $\mathbf{K} \equiv \sum_{a=1}^k \gamma_a \mathbb{E}[d\xi_a|_{\tilde{\mathbf{z}}_a}] \otimes \mathbf{C}_a \in \mathcal{M}_{kp};$
- $\tilde{\boldsymbol{\mu}} = \sum_{a=1}^k \gamma_a \mathbb{E}[\xi_a(\tilde{\mathbf{z}}_a)] \otimes \boldsymbol{\mu}_a \in \mathcal{M}_{kp};$
- $\tilde{\mathbf{C}} = \sum_{a=1}^k \gamma_a \mathbb{E}[\xi_a(\tilde{\mathbf{z}}_a) \xi_a(\tilde{\mathbf{z}}_a)^\top] \otimes \mathbf{C}_a \in \mathcal{M}_{kp};$
- $\bar{\mathbf{Q}} \equiv (\mathbf{I}_{kp} - \mathbf{K})^{-1} \in \mathcal{M}_{kp};$
- $\mathbf{m}_\ell = \boldsymbol{\mu}_\ell^\top \bar{\mathbf{Q}} \tilde{\boldsymbol{\mu}};$
- $\boldsymbol{\Delta}_\ell = \left[\frac{1}{n} \text{Tr} \left(\mathbf{C}_\ell \bar{\mathbf{Q}}_{\mathcal{I}_a, \mathcal{I}_b} \right) \right]_{1 \leq a, b \leq k} \in \mathcal{M}_k;$
- $\boldsymbol{\sigma}_\ell^2 = \left[\frac{1}{n} \text{Tr} \left(\mathbf{C}_\ell \mathbf{R}(\tilde{\mathbf{C}})_{\mathcal{I}_a, \mathcal{I}_b} \right) + \tilde{\boldsymbol{\mu}} \bar{\mathbf{Q}}_{\cdot, \mathcal{I}_a} \mathbf{C}_\ell \bar{\mathbf{Q}}_{\mathcal{I}_b, \cdot} \tilde{\boldsymbol{\mu}} \right]_{1 \leq a, b \leq k}.$

⁶Here, given $\mathbf{z} \in \mathbb{R}^k$, $df_\ell^2|_{\mathbf{z}}$ is a bilinear form on \mathbb{R}^k , thus $\|df_\ell^2|_{\mathbf{z}}\| = \sup_{\|\mathbf{a}\|, \|\mathbf{b}\| \leq 1} df_\ell^2|_{\mathbf{z}}(\mathbf{a}, \mathbf{b})$.

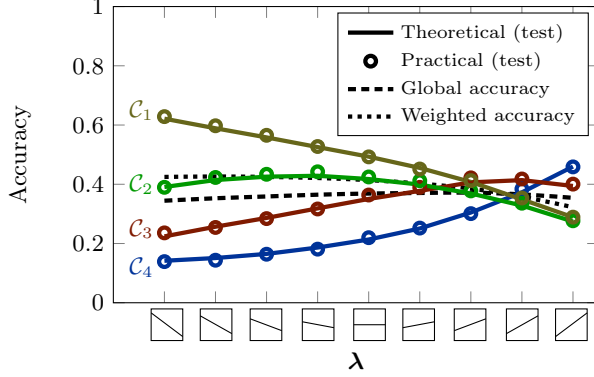


Figure 1: Theoretical and practical test accuracies on synthetic Gaussian data varying the regularizing vector λ . The considered parameters are $n = p = 200$, $k = 4$, $\gamma_1 = 2/5$, $\gamma_2 = 3/10$, $\gamma_3 = 1/5$ and $\gamma_4 = 1/10$. The statistics μ_ℓ and C_ℓ for $\ell \in [k]$ are sampled randomly as $\mu_\ell \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p/p)$ and $C_\ell \sim \text{Diag}(\mathbf{u}) + \sqrt{\mathbf{v}\mathbf{v}^\top}$ with $\mathbf{u} \sim \text{Unif}([0, 1]^p)$ and $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p/p)$.

Where $\forall a \in [k]$, $\mathcal{I}_a \equiv \{ap + 1, \dots, (a + 1)p\}$ and we define for any matrix $\mathbf{T} \in \mathcal{M}_{pk}$ and any $a, b \in [k]$:

$$\mathbf{T}_{\mathcal{I}_a, \mathcal{I}_b} \equiv [T_{i,j}]_{i \in \mathcal{I}_a, j \in \mathcal{I}_b}, \quad \mathbf{T}_{\mathcal{I}_a, \cdot} \equiv [T_{i,j}]_{i \in \mathcal{I}_a, j \in [pk]}, \quad \mathbf{T}_{\cdot, \mathcal{I}_b} \equiv [T_{i,j}]_{i \in [pk], j \in \mathcal{I}_b}.$$

Then one has the following estimations,

$$\|\mu_W - \bar{Q}\tilde{\mu}\| \leq \mathcal{O}(n^{-\frac{1}{2}}), \quad \|\mathbf{C}_W - \mathbf{R}\tilde{\mathbf{C}}\|_* \leq \mathcal{O}(n^{-\frac{1}{2}}),$$

and for a new datum \mathbf{x} in class \mathcal{C}_ℓ , independent of the training data $\mathbf{x}_1, \dots, \mathbf{x}_n$, one has

$$\|\mathbb{E}_{\mathcal{A}_W}[\tilde{\mathbf{x}}^\top \mathbf{W}] - \mathbf{m}_\ell\| \leq \mathcal{O}(n^{-\frac{1}{2}}), \quad \|\mathbb{E}_{\mathcal{A}_W}[\tilde{\mathbf{x}}^\top \mathbf{W} \mathbf{W}^\top \mathbf{x}] - \sigma_\ell^2 - \mathbf{m}_\ell \mathbf{m}_\ell^\top\| \leq \mathcal{O}(n^{-\frac{1}{2}}).$$

3 Further experiments

3.1 Synthetic Gaussian data

This subsection provides experiments using synthetic Gaussian data. Recalling the setting of Figure 1, which depicts the theoretical versus the practical accuracies (of a four-class classification problem with Softmax) for different choices of λ . One can first observe that the practical accuracies are perfectly estimated by their theoretical counterparts as per Theorem 2.4. We have also depicted in this figure the global accuracy of the classifier along with its weighted accuracy (computed as $\sum_\ell \gamma_\ell a_\ell$ with a_ℓ being the accuracy of class \mathcal{C}_ℓ). These quantities notably highlight the existence of an optimal choice of λ for maximising the classifier accuracy (for the less representative classes), which in turn might be anticipated through our present investigation.

3.2 GAN-generated MNIST images

We provide in this subsection experiments with a multi-class classification of data issued from a GAN⁷-generated version of the MNIST digits “1”, “2”, and “3” with the following settings:

1. We samples 50000 images of each class;
2. We multiplied each of the data by a randomly chosen matrix $\mathbf{A} \in \mathcal{M}_{p, p_{\text{MNIST}}}$, where $p = 200$ and $p_{\text{MNIST}} = 784$;
3. For each class \mathcal{C}_ℓ , $\ell \in [3]$, we divided each set of indexes $\mathcal{I}_\ell = [50000]$ into three subset: $\mathcal{I}_\ell^{\text{tr}}$, $\mathcal{I}_\ell^{\text{tst}}$ and $\mathcal{I}_\ell^{\text{pop}}$ such that $\#\mathcal{I}_\ell^{\text{tr}} = \#\mathcal{I}_\ell^{\text{tst}} = 5000$ and $\#\mathcal{I}_\ell^{\text{pop}} = 40000$;
4. For all $\ell \in [3]$, we computed the empirical mean μ_ℓ and covariance C_ℓ with the data indexed by $\mathcal{I}_\ell^{\text{pop}}$;

⁷We used a standard DC-GAN model [MO14] trained on the whole MNIST dataset.

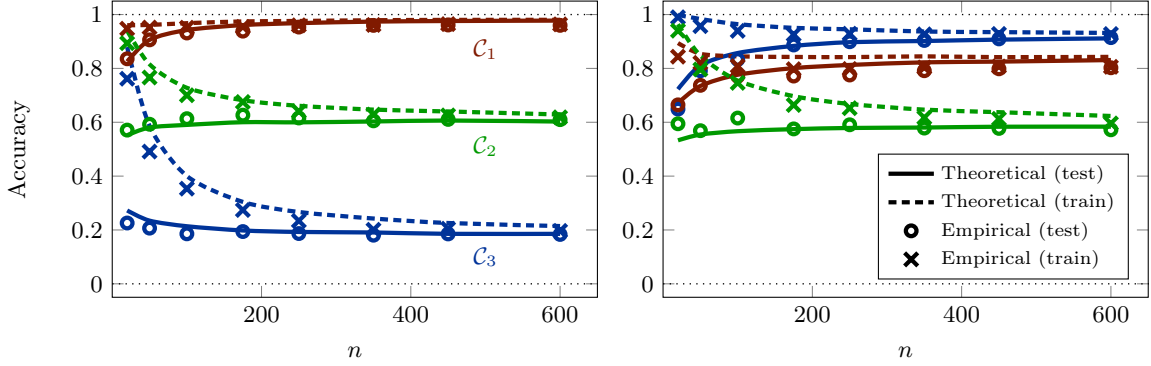


Figure 2: Classification accuracies of GAN-generated MNIST digits with unbalanced classes in the training set: $\gamma_1 = 1/2$, $\gamma_2 = 1/3$ and $\gamma_3 = 1/6$ for two choices of λ : $[30, 30, 30]$ (left) $[10, 20, 30]$ (right).

5. In order to make the classification problem harder, for all $\ell \in [3]$ we retrieved from each data indexed by I_ℓ^{tr} and I_ℓ^{tst} the vector $(1 - \alpha)\mu_\ell$, with $\alpha = 0.5$ and we corrected the mean ($\mu_\ell \leftarrow \alpha\mu_\ell$).

Figure 2 pictures with two different regularizing vector $\lambda \in \mathbb{R}^3$:

- The theoretical performance on the training set and on the test set based on the estimated class-wise means and covariances and relying on Theorem 2.4;
- The empirical performances computed on the data indexed by the sets $\mathcal{I}_\ell^{\text{tr}}$ and $\mathcal{I}_\ell^{\text{tst}}$, for $\ell \in [k]$.

As we can notice from Figure 2, the theoretical estimations closely match the empirical ones in both settings. Moreover, it is deductible that when the classes of the training set are unbalanced, the regularizing parameter vector λ should be fine-tuned accordingly. Besides, note that with an optimal choice of λ , the accuracy of classification of the less represented class is better than the accuracy of the two other classes. We therefore believe that our present investigation could be helpful for the automatic fine-tuning of λ .

3.3 CNN features of real Imagenet images

In this subsection, we provide further experiments using real images from the Imagenet dataset [DDS⁺09]. Figure 3-(left) depicts the learned Softmax weights against their expected large p, n asymptotics as predicted by Theorem 2.4. 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 3-(right) which strongly suggests that the conclusions of Theorem 2.4 generalize to real data.

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.
- [Kla07] B. Klartag. A central limit theorem for convex sets. *Inventiones mathematicae*, 168:91–131, 2007.
- [LC20a] 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.
- [LC20b] Cosme Louart and Romain Couillet. Concentration of solutions to random equations with concentration of measure hypotheses, 2020.
- [Led05] Michel Ledoux. *The concentration of measure phenomenon*. Number 89. American Mathematical Soc., 2005.

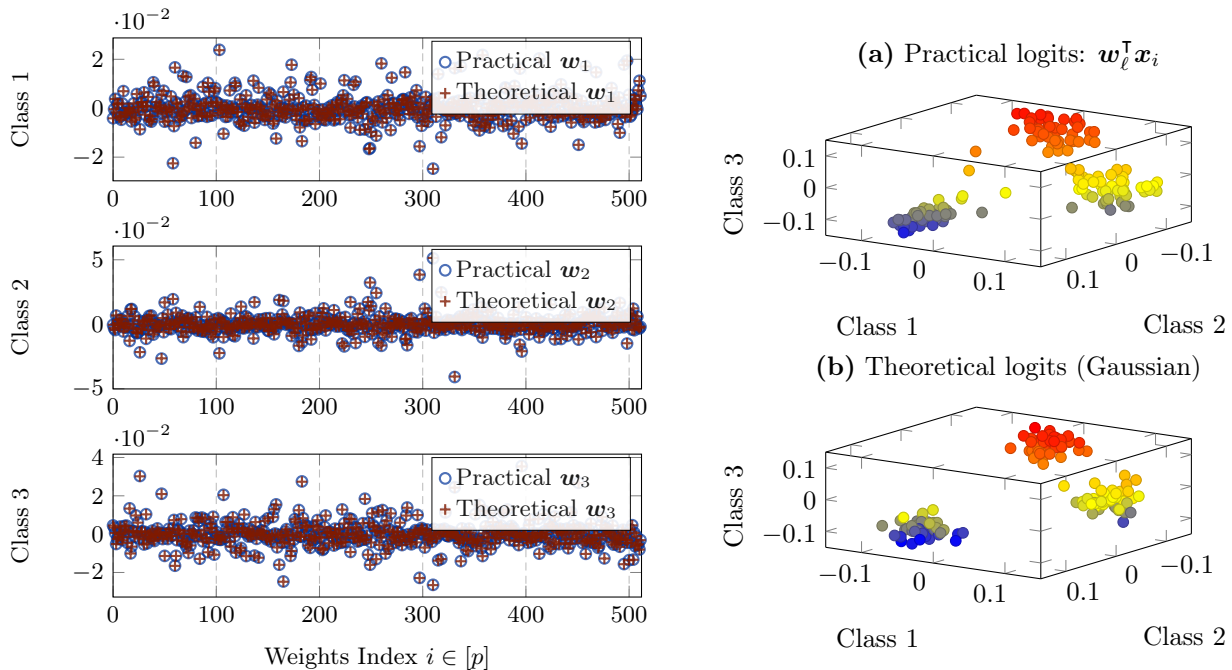


Figure 3: **(Left)** Learned weights (blue circles) versus theoretical estimates (red crosses) from Theorem 2.4. **(Right)** Practical (a) versus theory-predicted (b) logits, on a test set independent from the training set. The data are Resnet18 [SIVA17] representations ($p = 512$) of randomly selected images from the Imagenet dataset [DDS⁺09]; $k = 3$ classes: *hamburger*, *mushroom*, *pizza*; $n = 3000$; regularization constants $\lambda_1 = \lambda_2 = \lambda_3 = 1.5$; data normalized such that $\|\mathbf{x}_i\| = 0.1 \cdot \sqrt{p}$ to ensure \mathcal{A}_w .

- [MO14] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [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.