Mohamed El Amine Seddik [1,2]    Cosme Louart [1,3]    Mohamed Tamaazousti [1]    Romain Couillet [2,3]

[1]CEA List    [2]CentraleSupélec    [3]GIPSA-Lab University of GrenobleAlpes

## Abstract



$$x_1, \ldots, x_n \in \mathbb{R}^p$$
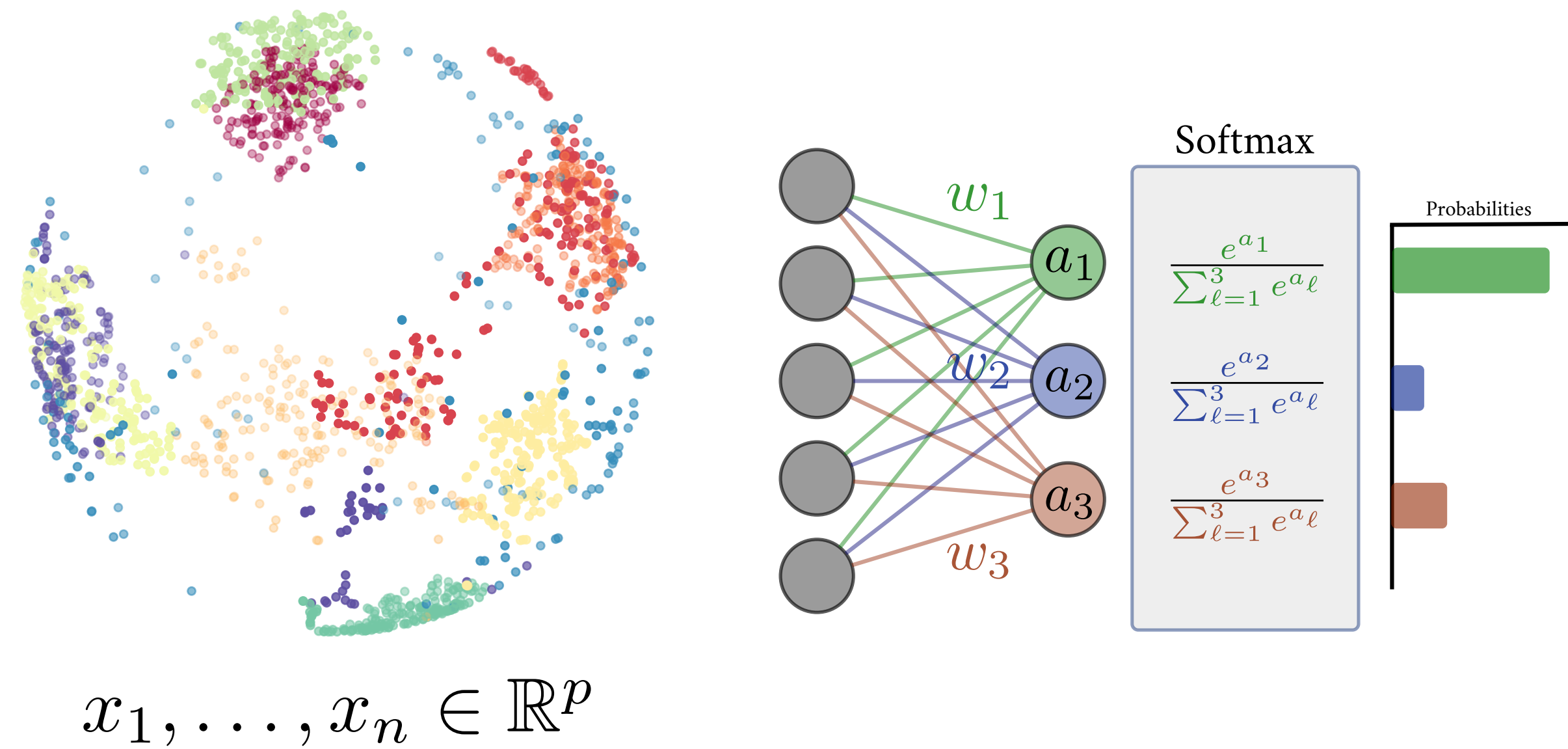
- RMT analysis of Softmax with high-dimensional concentrated inputs.
- Softmax weights depend only on data means and covariances.
- Asymptotic performance of Softmax derived based on first data moments.

## Notion of Concentrated Vectors

**Definition 1.** $\mathcal{X} \ni x$ is $q$-exponentially **concentrated** if for all $\varphi : \mathcal{X} \to \mathbb{R}$ 1-**Lipschitz**, there exists $C \geq 0$ independent of $\dim(\mathcal{X})$ and $\sigma > 0$ such that,

$$\forall t \geq 0, \quad \mathbb{P}\left(|\varphi(x) - \mathbb{E}\varphi(x)| > t\right) \leq C\, e^{-(t/\sigma)^q},$$

denoted $x \propto \mathcal{E}_q(\sigma)$ or $x \propto \mathcal{E}_q$ if $\sigma$ independent of $\dim(\mathcal{X})$.

**Examples:**

- $x = \mu + \Sigma^{1/2}z \in \mathbb{R}^p$ with $z \sim \mathcal{N}(0, I_p)$ and $\|\Sigma\| < \infty$, then $x \propto \mathcal{E}_2$.
- If $\mathcal{Z} \ni z \propto \mathcal{E}_q$ and $\mathcal{G} : \mathcal{Z} \to \mathcal{X}$ $L$-Lipschitz, then $\mathcal{G}(z) \propto \mathcal{E}_q(L)$.



Figure 1. Images generated by the BigGAN model [3].

## Model

(M) **Data matrix** (distributed in $k$ classes $\mathcal{C}_1, \mathcal{C}_2, \ldots, \mathcal{C}_k$):

$$\mathbb{R}^{p \times n} \ni X = \Big[\underbrace{x_1, \ldots, x_{n_1}}_{\sim \mathcal{L}(\mu_1, \Sigma_1)}, \underbrace{x_{n_1+1}, \ldots, x_{n_2}}_{\sim \mathcal{L}(\mu_2, \Sigma_2)}, \ldots, \underbrace{x_{n-n_k+1}, \ldots, x_n}_{\sim \mathcal{L}(\mu_k, \Sigma_k)}\Big] \propto \mathcal{E}_2$$

**Model statistics:** $\mu_\ell = \mathbb{E}_{x_i \in \mathcal{C}_\ell}[x_i], \quad \Sigma_\ell = \mathbb{E}_{x_i \in \mathcal{C}_\ell}[x_i x_i^\mathsf{T}] - \mu_\ell \mu_\ell^\mathsf{T}.$

## Softmax Classifier

Minimize:

$$\mathcal{L}(w_1, \ldots, w_k) = -\frac{1}{n}\sum_{i=1}^{n}\sum_{\ell=1}^{k} y_{i\ell}\log p_{i\ell} + \frac{1}{2}\sum_{\ell=1}^{k}\lambda_\ell\|w_\ell\|^2$$

$$p_{i\ell} = \frac{\exp(w_\ell^\mathsf{T} x_i)}{\sum_{j=1}^{k}\exp(w_j^\mathsf{T} x_i)}, \quad W \equiv [w_1^\mathsf{T}, \ldots, w_k^\mathsf{T}]^\mathsf{T} \in \mathbb{R}^{pk}$$

Implicit Equation (scalar case for some $f : \mathbb{R} \to \mathbb{R}$)

$$\mathbb{R}^p \ni w = \frac{1}{n}\sum_{i=1}^{n} f(w^\mathsf{T} x_i) x_i \quad \Rightarrow \quad \boxed{w = \Psi(w) \equiv \frac{1}{n} X f(X^\mathsf{T} w)}$$

$\Psi$ is requested to be $(1-\varepsilon)$-Lipschitz for some $\varepsilon > 0$ or equivalently

$$\mathcal{A}_w = \left\{\frac{1}{n}\|f'\|_\infty\|XX^\mathsf{T}\| \geq 1 - \varepsilon\right\} \quad \text{has low probability.}$$

## Assumptions

(A) **Growth rate assumptions:** As $p \to \infty$,

1. $p/n \to c \in (0, \infty)$ and $|\mathcal{C}_\ell|/n \to \gamma_\ell \in (0, 1)$.
2. $k$ fixed.
3. $\|\mu_\ell\| = \mathcal{O}(1)$ for each $\ell \in [k]$.
4. $\exists \varepsilon > 0$ independent of $p, n$ s.t. $\frac{1}{n}\|f'\|_\infty\|XX^\mathsf{T}\| \leq 1 - 2\varepsilon$.

## Main Result

Evaluate $\mu_w = \mathbb{E}[w] = \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}[f(x_i^\mathsf{T} w)x_i]$ and $\Sigma_w = \mathbb{E}[ww^\mathsf{T}] - \mu_w\mu_w^\mathsf{T}$
Under (M-A), $\mathbb{P}(\mathcal{A}_w) \propto e^{-n}$ and $w \propto \mathcal{E}_2\left(n^{-\frac{1}{2}}\right) \mid \mathcal{A}_w$, and there exists $(\delta, m, \sigma) \in (\mathbb{R}^k)^3$ satisfying

$$z_\ell \sim \mathcal{N}(m_\ell, \sigma_\ell^2); \quad \delta_\ell = \frac{1}{n}\mathrm{Tr}\left(\Sigma_\ell\left(I_p - K\right)^{-1}\right);$$

$$\tilde{\mu} \equiv \sum_{\ell=1}^{k}\gamma_\ell\mathbb{E}[g_\ell(z_\ell)]\mu_\ell; \quad \tilde{\Sigma} \equiv \sum_{\ell=1}^{k}\gamma_\ell\mathbb{E}[g_\ell(z_\ell)^2]\Sigma_\ell; \quad K \equiv \sum_{\ell=1}^{k}\gamma_\ell\mathbb{E}[g_\ell'(z_\ell)]\Sigma_\ell;$$

$$R_1 \equiv (I_p - K)^{-1}; \quad R_2(M) \equiv M + KR_2(M)K;$$

$$m_\ell \equiv \mu_\ell^\mathsf{T} R_1\tilde{\mu}; \quad \sigma_\ell^2 \equiv \frac{1}{n}\mathrm{Tr}(\Sigma_\ell R_2(\tilde{\Sigma})) + \tilde{\mu}^\mathsf{T} R_1\Sigma_\ell R_1\tilde{\mu};$$

Furthermore,

$$\|\mu_w - R_1\tilde{\mu}\| \leq \mathcal{O}\left(n^{-\frac{1}{2}}\right), \quad \|\Sigma_w - \frac{1}{n}R_2(\tilde{\Sigma})\|_* \leq \mathcal{O}\left(n^{-\frac{1}{2}}\right)$$

**Key Observation:** Only **first** and **second** order statistics matter!

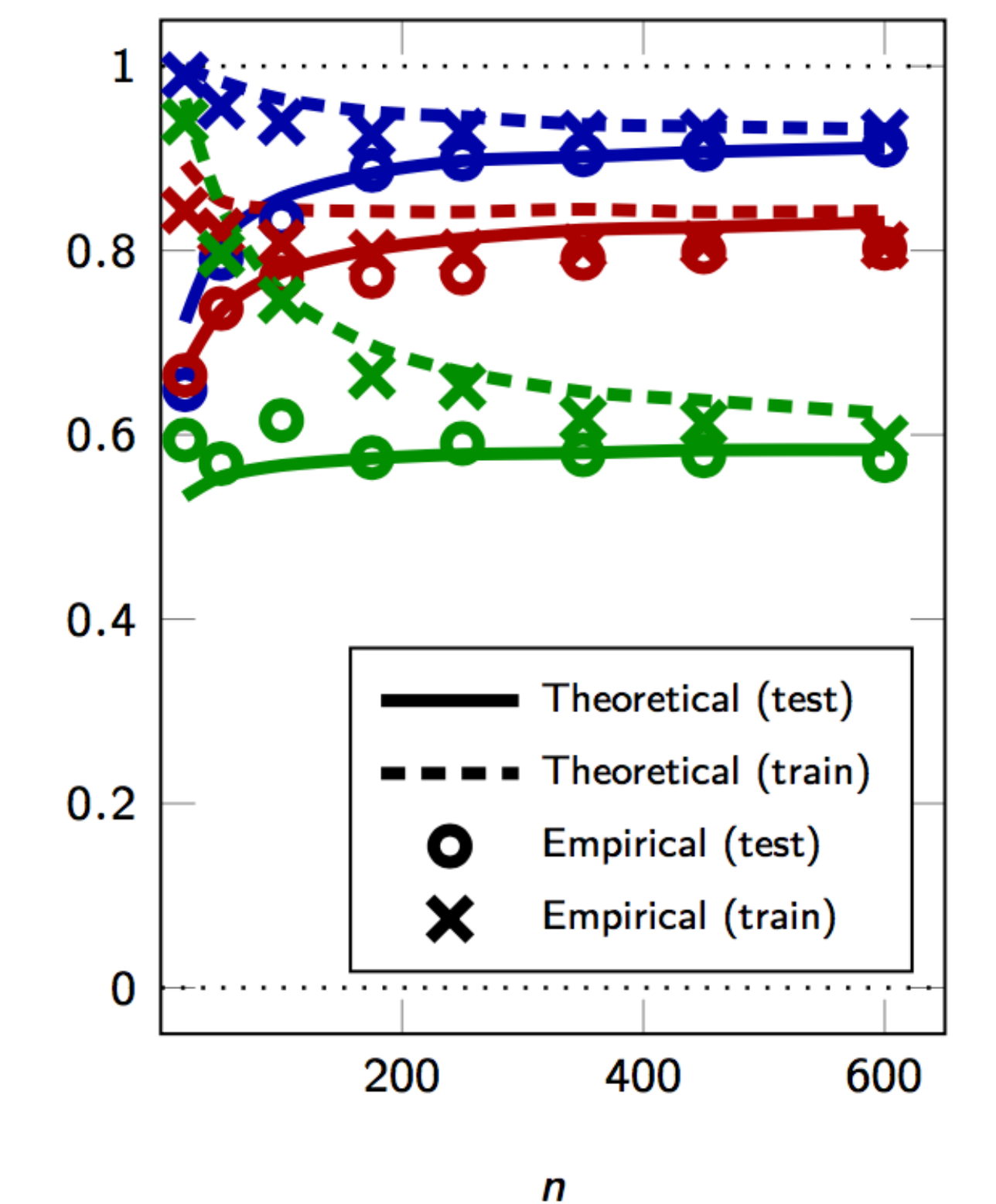## Simulations





Accuracy: $\lambda_1 = \lambda_2 = \lambda_3 = 30$

$\lambda_1 = 10, \lambda_2 = 20, \lambda_3 = 30$

- —— Theoretical (test)
- - - - Theoretical (train)
- ○ Empirical (test)
- ✕ Empirical (train)

## Conclusion

- (El-Karoui+'13, Mai+'19) analyzed logistic regression under **Gaussian** data.
- We generalized these ideas to a $k$-class mixture of **concentrated** data.
- **Universality:** *"Softmax treats input data as Gaussian random vectors"*.
- **Optimality:** Softmax is optimal for data with *strongly discriminative class-wise means* as suggested by distance-based image classification approaches (Mensink+'13).

## References

[1] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, "Generative adversarial nets", in NIPS 2014.

[2] Terence Tao, "Topics in random matrix theory, volume 132". American Mathematical Society Providence, RI, 2012.

[3] Andrew Brock, Jeff Donahue, and Karen Simonyan, "Large scale GAN training for high fidelity image synthesis", in ICLR 2019.