

Softmax is $1/2$ -Lipschitz: A tight bound across all ℓ_p norms

Pravin Nair

Department of Electrical Engineering
Indian Institute of Technology, Madras

pravin@ee.iitm.ac.in

Reviewed on OpenReview: <https://openreview.net/forum?id=6dowaHsa6D>

Abstract

The softmax function is a basic operator in machine learning and optimization, used in classification, attention mechanisms, reinforcement learning, game theory, and problems involving log-sum-exp terms. Existing robustness guarantees of learning models and convergence analysis of optimization algorithms typically consider the softmax operator to have a Lipschitz constant of 1 with respect to the ℓ_2 norm. In this work, we prove that the softmax function is contractive with the Lipschitz constant $1/2$, uniformly across all ℓ_p norms with $p \geq 1$. We also show that the local Lipschitz constant of softmax attains $1/2$ for $p = 1$ and $p = \infty$, and for $p \in (1, \infty)$, the constant remains strictly below $1/2$ and the supremum $1/2$ is achieved only in the limit. To our knowledge, this is the first comprehensive norm-uniform analysis of softmax Lipschitz continuity. We demonstrate how the sharper constant directly improves a range of existing theoretical results on robustness and convergence. We further validate the sharpness of the $1/2$ Lipschitz constant of the softmax operator through empirical studies on attention-based architectures (ViT, GPT-2, Qwen3-8B) and on stochastic policies in reinforcement learning.

1 Introduction

The softmax function has applications across diverse areas of machine learning and optimization. In classification models, model output logits are normalized using the softmax function to form a probability distribution (Goodfellow et al., 2016). In Transformer architectures, softmax normalizes attention scores to compute weighted combinations of input token features, enabling high-quality feature refinement (Vaswani et al., 2017). These transformer architectures have shown state-of-the-art performance in various applications in natural language processing (Tunstall et al., 2022), computer vision (Khan et al., 2022), reinforcement learning (Li et al., 2023), etc. Since softmax is the gradient of the log-sum-exp function, its properties are crucial for analyzing optimization algorithms (Gao & Pavel, 2017; Nachum et al., 2017; Nesterov, 2005). In reinforcement learning, softmax is used to convert action-value estimates (Q-values) to stochastic policies. This results in a probability distribution over actions that favours actions with higher Q-values and still assigns probabilities to suboptimal actions. In entropy-regularized reinforcement learning, the softmax operator arises naturally as the solution of the policy optimization problem (Sutton & Barto, 2018; Nachum et al., 2017).

The softmax operator maps any real vector to a probability distribution. Formally, the softmax function $\sigma_\lambda : \mathbb{R}^n \rightarrow \Delta_n^\circ$ with inverse-temperature parameter $\lambda > 0$ is defined as

$$\sigma_\lambda(\mathbf{x})_i = \frac{\exp(\lambda x_i)}{\sum_{j=1}^n \exp(\lambda x_j)}, \quad i = 1, \dots, n, \quad (1)$$

where $\Delta_n = \{\mathbf{u} \in \mathbb{R}^n : u_i \geq 0, \sum_i u_i = 1\}$ is the probability simplex and $\Delta_n^\circ = \{\mathbf{u} \in \mathbb{R}^n : u_i > 0, \sum_i u_i = 1\}$ denotes its interior. Since $\exp(x_i) > 0$ for all $x_i \in \mathbb{R}$, $\sigma_\lambda(\mathbf{x})$ never attains the boundary of the simplex $\partial\Delta_n = \{\mathbf{u} \in \mathbb{R}^n : u_i = 0 \text{ for some } i\}$. The coefficient λ measures the smoothness of the output distribution. A larger λ makes the distribution more peaked on the largest entries, and a smaller λ makes it smoother with more evenly spread probability values. The standard softmax function is when $\lambda = 1$.

Accurate analysis of softmax’s Lipschitz constant has important applications across machine learning and optimization. It enables accurate robustness analysis of attention mechanisms, where Lipschitz bounds can be used to set hyperparameters that ensure stable training (Qi et al., 2023). In game-theoretic learning, modeling action selection via softmax allows Lipschitz continuity to be directly leveraged for establishing convergence rates toward Nash equilibria (Gao & Pavel, 2017). Many existing generalization results for learning models require computing the Lipschitz constant of the loss function with respect to either network input or network parameters. When softmax appears in the final layer or in attention, an accurate Lipschitz analysis becomes essential for deriving sharp generalization bounds (Asadi & Abbe, 2020). Moreover, when a model is expressed as a composition of layers, standard network analyses bound the global Lipschitz constant by an appropriate composition (often a product) of layer-wise constants Bartlett et al. (2017); Miyato et al. (2018); Virmaux & Scaman (2018); Fazlyab et al. (2019). In such bounds, the softmax Lipschitz constant acts as a multiplicative factor, so tightening it leads directly to sharper network-level guarantees. In classification networks, robustness and domain adaptation depend on how sensitive output probabilities are to input perturbations, which is again controlled by the Lipschitz constant of softmax (Chen et al., 2022). The log-sum-exp (LSE) function is a smooth approximation of the maximum operator, and the gradient of LSE is the softmax function. Hence, in optimization problems with LSE-based objective functions, step-size rules for gradient descent and smoothness-based convergence rates can be sharpened using a tight softmax Lipschitz bound. Lipschitz constant of softmax is also used in analysing several other frameworks, including meta-learning (Jeon et al., 2024), structural causal models (Le Priol et al., 2021), prior-data fitted networks (Nagler, 2023), and sparse training (Lei et al., 2024). The above-mentioned diverse set of applications highlights that deriving a tight Lipschitz constant of softmax has a substantial practical impact.

Despite the theoretical and empirical significance of softmax’s Lipschitz constant, an accurate analysis is missing in the literature. A number of works (Gao & Pavel, 2017; Laha et al., 2018; Asadi & Abbe, 2020; Gouk et al., 2021; Le Priol et al., 2021; Chen et al., 2022; Nagler, 2023; Jeon et al., 2024; Lei et al., 2024) assume the softmax function to be 1-Lipschitz with respect to the ℓ_2 norm. This assumption is often attributed to Proposition 4 in Gao & Pavel (2017). In some works, this is mentioned as a straightforward fact, perhaps because showing the Lipschitz constant to be bounded by 1 is trivial. We also note that Xu et al. (2022) reports a bound of $1/4$ for the softmax’s Lipschitz constant. The main motivation for this work is the discrepancy between the commonly assumed Lipschitz constant of the softmax function and the smaller values observed in our empirical studies on transformer architectures. This prompted us to derive a tighter bound for the Lipschitz constant of the softmax function. In this regard, our contributions are as follows:

- We prove that softmax is $1/2$ -Lipschitz across all ℓ_p norms ($p \geq 1$), thereby improving the usually assumed bound of 1.
- Furthermore, we show the tightness of this bound by proving that the local Lipschitz constant is attained for $p = 1$ and $p = \infty$ as $1/2$, while for $p \in (1, \infty)$ the local Lipschitz constants remain strictly below $1/2$ and the supremum is only approached in the limit.
- We demonstrate that our tight Lipschitz bound improves the existing theoretical results in robustness and convergence. We also validate our derived Lipschitz constant through experiments on large-scale transformer models and reinforcement learning policies.

2 Preliminaries

In this section, we provide the necessary background for our main results. We start with the definition of the ℓ_p norm for vectors and the corresponding induced operator norm for matrices.

Definition 2.1 (ℓ_p norm). For a vector $\mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$ and for $1 \leq p < \infty$, the ℓ_p norm of \mathbf{x} is defined as

$$\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}.$$

For $p = \infty$, the ℓ_∞ norm is defined as

$$\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|.$$

We can also define the operator norm of a matrix $\mathbf{A} = (A_{ij}) \in \mathbb{R}^{m \times n}$ induced by the ℓ_p norm as

$$\|\mathbf{A}\|_p := \sup_{\mathbf{v} \neq 0} \frac{\|\mathbf{A}\mathbf{v}\|_p}{\|\mathbf{v}\|_p} = \sup_{\|\mathbf{v}\|_p=1} \|\mathbf{A}\mathbf{v}\|_p.$$

Throughout the paper, we overload $\|\cdot\|_p$ to represent both the vector norm and the corresponding induced norm on matrices. Next, we define the Lipschitz property of functions.

Definition 2.2 (Lipschitz continuity in ℓ_p norm). Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and let $\|\cdot\|_p$ denote the ℓ_p norm with $1 \leq p \leq \infty$. The function f is said to be Lipschitz continuous with respect to $\|\cdot\|_p$ if there exists a constant $L_p \geq 0$ such that

$$\|f(\mathbf{x}) - f(\mathbf{y})\|_p \leq L_p \|\mathbf{x} - \mathbf{y}\|_p, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n. \quad (2)$$

The smallest such constant L_p is called the Lipschitz constant of f with respect to the ℓ_p norm.

A mapping f is contractive if $L_p < 1$, and non-expansive if $L_p \leq 1$. It is firmly non-expansive if $\|f(\mathbf{x}) - f(\mathbf{y})\|_p^2 \leq \langle f(\mathbf{x}) - f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$ for all \mathbf{x}, \mathbf{y} , and co-coercive with constant $\beta > 0$ if $\langle f(\mathbf{x}) - f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \beta \|f(\mathbf{x}) - f(\mathbf{y})\|_p^2$.

Note that the Lipschitz constant in Eq. 2 is a global quantity, since the inequality holds for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. By definition, L_p can be characterized as,

$$L_p = \sup_{\mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \mathbf{x} \neq \mathbf{y}} \frac{\|f(\mathbf{x}) - f(\mathbf{y})\|_p}{\|\mathbf{x} - \mathbf{y}\|_p}. \quad (3)$$

We next introduce the notion of a local Lipschitz constant (Rockafellar & Wets, 1998) via the Jacobian of f , which in turn provides another characterization of the global Lipschitz constant.

Definition 2.3 (Jacobian matrix (Boyd & Vandenberghe, 2004)). Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a differentiable function, where $f(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x}))$ and each $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is scalar-valued function. The Jacobian matrix of f at \mathbf{x} , denoted $\mathbf{J}_f(\mathbf{x}) \in \mathbb{R}^{m \times n}$, is the matrix of all partial derivatives,

$$\mathbf{J}_f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}) & \cdots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{x}) & \cdots & \frac{\partial f_m}{\partial x_n}(\mathbf{x}) \end{bmatrix}.$$

If the function is continuously differentiable, the local Lipschitz constant of a function at any $\mathbf{x} \in \mathbb{R}^n$ is equal to the ℓ_p norm of the Jacobian at that point.

Definition 2.4 (Local Lipschitz constant (Rockafellar & Wets, 1998)). Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$. The local Lipschitz constant of f at a point $\mathbf{x} \in \mathbb{R}^n$ with respect to the ℓ_p norm is defined as

$$L_p(\mathbf{x}) := \limsup_{\mathbf{y} \rightarrow \mathbf{x}, \mathbf{y} \neq \mathbf{x}} \frac{\|f(\mathbf{y}) - f(\mathbf{x})\|_p}{\|\mathbf{y} - \mathbf{x}\|_p}.$$

Intuitively, $L_p(\mathbf{x})$ characterizes the tightest Lipschitz bound in an arbitrarily small neighborhood of \mathbf{x} . If f is differentiable at \mathbf{x} , then

$$L_p(\mathbf{x}) = \|\mathbf{J}_f(\mathbf{x})\|_p,$$

where $\|\mathbf{J}_f(\mathbf{x})\|_p$ denotes ℓ_p -induced operator norm of the Jacobian matrix.

We now relate the global Lipschitz constant of a function to its local Lipschitz constant, following a result from Hytönen et al. (2016).

Lemma 1 (Lipschitz constant via the Jacobian). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be continuously differentiable. Then, for $1 \leq p \leq \infty$, the global Lipschitz constant of f with respect to ℓ_p norm is given as*

$$L_p = \sup_{\mathbf{x} \in \mathbb{R}^n} L_p(\mathbf{x}) = \sup_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{J}_f(\mathbf{x})\|_p.$$

Thus, the global Lipschitz constant of a function can be expressed in a variational form involving the ℓ_p -induced operator norm of its Jacobian matrix. To apply this principle to the softmax operator, we recall its Jacobian formulation from Gao & Pavel (2017), as stated below.

Lemma 2 (Jacobian of the softmax). *Let $\sigma_\lambda : \mathbb{R}^n \rightarrow \Delta_n^\circ$ be the softmax function as in Definition 1 and let $\mathbf{s} = \sigma_\lambda(\mathbf{x})$. Then, the Jacobian of σ_λ at \mathbf{x} is*

$$\mathbf{J}_{\sigma_\lambda}(\mathbf{x}) = \lambda (\text{Diag}(\mathbf{s}) - \mathbf{s}\mathbf{s}^\top),$$

so, in coordinates, if $\mathbf{s} = \{s_1, s_2, \dots, s_n\}$,

$$[\mathbf{J}_{\sigma_\lambda}(\mathbf{x})]_{ij} = \begin{cases} \lambda s_i(1 - s_i), & i = j, \\ -\lambda s_i s_j, & i \neq j. \end{cases}$$

Refer to Proposition 2 in Gao & Pavel (2017) for the proof of Lemma 2. In Section 3, we will make use of Lemma 1 and the result in Lemma 2 for the softmax’s Jacobian to establish our main results.

3 Main Results

We begin by deriving an inequality result for ℓ_p -induced operator norms for matrices.

Proposition 1 (Norm Interpolation). *Let $\mathbf{A} \in \mathbb{R}^{n \times n}$. Then:*

$$(a) \|\mathbf{A}\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |A_{ij}|.$$

$$(b) \|\mathbf{A}\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |A_{ij}|.$$

$$(c) \text{ For } 1 < p < \infty,$$

$$\|\mathbf{A}\|_p \leq \|\mathbf{A}\|_1^{1/p} \|\mathbf{A}\|_\infty^{1-1/p}.$$

Proposition 1(a) and (b) state the well-known results that the ℓ_1 -induced operator norm of a matrix is equal to its maximum absolute column sum, and the ℓ_∞ -induced operator norm is equal to the maximum absolute row sum (Golub & Van Loan, 2013). Proposition 1(c) establishes an interpolation inequality that upper bounds the ℓ_p -induced operator norm of a matrix in terms of $\|\mathbf{A}\|_1$ and $\|\mathbf{A}\|_\infty$. While this inequality is a special case of the Riesz–Thorin interpolation theorem (Riesz, 1927; Stein, 1956; Thorin, 1939), we provide a self-contained proof in the Appendix A.1 for clarity. Importantly, Proposition 1(c) forms the key technical tool for this section, as it enables us to derive norm-uniform bounds on the Lipschitz constant of the softmax operator. In particular, we reformulate the optimization problem in Lemma 1 using the Jacobian form of the softmax provided in Lemma 2.

Lemma 3 (Lipschitz constant of the softmax operator). *Let $\sigma_\lambda : \mathbb{R}^n \rightarrow \Delta_n^\circ$ be the softmax function as in Definition 1, and let $\mathbf{s} = \sigma_\lambda(\mathbf{x})$. Then, for any $1 \leq p \leq \infty$, the Lipschitz constant of σ_λ with respect to ℓ_p norm is given by*

$$L_p = \lambda \sup_{\mathbf{s} \in \Delta_n^\circ} \|\text{Diag}(\mathbf{s}) - \mathbf{s}\mathbf{s}^\top\|_p. \quad (4)$$

The key implication of Lemma 3 is that the global Lipschitz constant reduces to the supremum of the induced matrix norm in Eq. 4 over the interior of the probability simplex. The boundary points of the simplex, $\partial\Delta_n$, do not influence the formulation of the Lipschitz constant, since all components of the softmax output are strictly positive.

Next, we combine the interpolation inequality from Proposition 1 with the Jacobian formulation of the softmax operator from Lemma 3 to establish its Lipschitz constant across all ℓ_p norms. This result is summarized as the main contribution of our work in Theorem 1.

Theorem 1 (Lipschitz constant of softmax function). *Irrespective of the ℓ_p norms ($p \geq 1$),*

$$(a) \|\mathbf{J}_{\sigma_1}(\mathbf{x})\|_p \leq \frac{1}{2}$$

$$(b) \|\sigma_\lambda(\mathbf{x}) - \sigma_\lambda(\mathbf{y})\|_p \leq \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|_p$$

The upper bound on the Jacobian of the softmax operator across ℓ_p norms, established in Theorem 1, directly provides the global Lipschitz constant of the softmax function. Refer to Appendix A.2 for a detailed proof. A natural question arises that while many works assume the constant to be 1, we show it is in fact $1/2$; but could it be even smaller, for instance $1/4$ as claimed by Xu et al. (2022). Proposition 2 resolves this question by proving that $\lambda/2$ is indeed the tight global Lipschitz constant of the softmax operator for any $\lambda > 0$.

Proposition 2 (Tightness of the Lipschitz constant for softmax). *Consider the optimization problem,*

$$\sup_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{J}_{\sigma_1}(\mathbf{x})\|_p = \sup_{\mathbf{s} \in \Delta_n^\circ} \|\text{Diag}(\mathbf{s}) - \mathbf{s}\mathbf{s}^\top\|_p,$$

where Δ_n° is the interior of the probability simplex.

- (a) For $p = 1$ or $p = \infty$, the supremum is $1/2$ and is attained at an interior point in Δ_n° . In particular, for $\mathbf{x} = (\ln(n-1), 0, 0, \dots, 0) \in \mathbb{R}^n$, the softmax output $\mathbf{s} = \sigma_1(\mathbf{x})$ satisfies $\mathbf{s} \in \Delta_n^\circ$ and

$$\|\text{Diag}(\mathbf{s}) - \mathbf{s}\mathbf{s}^\top\|_p = \frac{1}{2}.$$

- (b) For $p \in (1, \infty)$, the supremum value is again $1/2$ but is not attained in Δ_n° for $n > 2$. Instead, there exists a sequence $\{\mathbf{s}_k\}_{k \geq 1} \subset \Delta_n^\circ$ converging to the boundary of the probability simplex $\partial\Delta_n$ such that

$$\lim_{k \rightarrow \infty} \|\text{Diag}(\mathbf{s}_k) - \mathbf{s}_k\mathbf{s}_k^\top\|_p = \frac{1}{2},$$

with $\mathbf{s}_k = \sigma_1(\mathbf{x}_k)$ for some $\mathbf{x}_k \in \mathbb{R}^n$. For $n = 2$, the supremum $\frac{1}{2}$ is attained at a point in Δ_2° .

The proof of Proposition 2 is deferred to Appendix A.3. In particular, part (b) relies on the fact that for $1 < p < \infty$, the interpolation inequality in Proposition 1 is indeed strict for $\mathbf{J}_{\sigma_1}(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}^n$. This is because the supremum in the optimization problem of Proposition 2 can be attained only on the boundary points of the probability simplex. Specifically, the extremal points are the permutations of $(1/2, 1/2, 0, \dots, 0)$, which lie in $\partial\Delta_n$.

Remark. There are multiple works showing that the Lipschitz constant of the softmax operator with respect to the ℓ_2 norm is upper bounded by $1/2$. Alghamdi et al. (2022) upper bounds $\|\mathbf{J}_{\sigma_1}(\mathbf{x})\|_2$ by first controlling the Frobenius norm of the softmax Jacobian (using $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F$) and then maximising this bound over the simplex to obtain the same constant $1/2$. Yudin et al. (2025) analyse the Jacobian matrix $\mathbf{M}(\mathbf{s}) = \text{diag}(\mathbf{s}) - \mathbf{s}\mathbf{s}^\top$ on the simplex and upper bound its spectral norm via an eigenvalue argument (reducing the worst case to a distribution supported on two classes), yielding $\|\mathbf{M}(\mathbf{s})\|_2 \leq 1/2$. Similarly, Newhouse (2025) gives an independent proof by applying Gershgorin's circle theorem to $\text{diag}(\mathbf{s}) - \mathbf{s}\mathbf{s}^\top$, and also shows tightness in the ℓ_2 setting. In contrast, we prove that the bound $1/2$ is tight uniformly across all ℓ_p norms for $p \geq 1$, and give a detailed analysis of attainability on the simplex for different p .

We next provide an example of a pair of points for which the empirically computed Lipschitz ratio approaches $1/2$. This gives empirical evidence for the tightness of our theoretical bound.

Example 1 (Tightness of the Lipschitz bound). *Let $K = 20$ and $\varepsilon = 10^{-4}$. Consider the vectors*

$$\mathbf{x} = (0, 0, -K, -K, \dots, -K) \in \mathbb{R}^{10}, \quad \mathbf{y} = \mathbf{x} + \varepsilon \mathbf{v},$$

where \mathbf{v} is eigenvector corresponding to the maximum eigenvalue of $\mathbf{J}_{\sigma_1}(\mathbf{x})$. Then the ratio

$$\frac{\|\sigma_1(\mathbf{y}) - \sigma_1(\mathbf{x})\|_p}{\|\mathbf{y} - \mathbf{x}\|_p}$$

evaluates approximately to 0.49999999504472 for all $p \geq 1$. This demonstrates a concrete pair (\mathbf{x}, \mathbf{y}) where the Lipschitz constant of the softmax function is nearly attained. We can extend the example to any dimension by adding $-K$ as a value to the other added dimensions.

4 Implications of the Improved Softmax Lipschitz Constant

Our result in Theorem 1 establishes that the softmax function is $\lambda/2$ -Lipschitz with respect to all ℓ_p norms for $p \geq 1$. This tighter characterization enables us to revisit existing results where softmax Lipschitz continuity is relevant, leading to sharper constants or simplified identities. Next, we illustrate this in a few representative settings.

4.1 Refinement of (Gao & Pavel, 2017, Cor. 3).

Leveraging our sharper Lipschitz estimate, Corollary 3 in Gao & Pavel (2017) can be strengthened as follows.

Corollary 1 (Softmax regularity with improved constants). *For any $\lambda > 0$, the softmax map $\sigma_\lambda : \mathbb{R}^n \rightarrow \Delta_n^\circ$ satisfies*

$$\begin{aligned} (\text{Lipschitz}) \quad & \|\sigma_\lambda(\mathbf{x}) - \sigma_\lambda(\mathbf{y})\|_p \leq \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|_p, \quad \forall p \geq 1, \\ (\text{co-coercive}) \quad & \langle \sigma_\lambda(\mathbf{x}) - \sigma_\lambda(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \frac{2}{\lambda} \|\sigma_\lambda(\mathbf{x}) - \sigma_\lambda(\mathbf{y})\|_2^2. \end{aligned}$$

In particular,

- σ_λ is nonexpansive and firmly nonexpansive for $\lambda \in (0, 2]$;
- σ_λ is contractive for $\lambda \in (0, 2)$.

The proof of this corollary follows directly from the analysis in Gao & Pavel (2017), with the Lipschitz constant of the softmax operator replaced by the sharper value $\lambda/2$.

4.2 Lipschitz Analysis of Attention variant

The self-attention module (Vaswani et al., 2017) refines an input feature matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ by projecting it into queries, keys, and values via learnable matrices $\mathbf{W}^Q \in \mathbb{R}^{d \times d_k}$, $\mathbf{W}^K \in \mathbb{R}^{d \times d_k}$, and $\mathbf{W}^V \in \mathbb{R}^{d \times d_v}$. The attention operation is then given by

$$\text{Att}(\mathbf{X}; \mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V) = \text{softmax}\left(\frac{\mathbf{X}\mathbf{W}^Q(\mathbf{X}\mathbf{W}^K)^\top}{\sqrt{d_k}}\right)\mathbf{X}\mathbf{W}^V, \quad (5)$$

where $\text{softmax}(\cdot)$ denotes the row-wise application of the function σ_1 . Although we obtain a tight global Lipschitz constant for the softmax operator σ_1 acting on the rows of the attention score matrix, the self-attention map $\mathbf{X} \mapsto \text{Att}(\mathbf{X}; \mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V)$ defined on an unbounded input domain is not globally Lipschitz. Kim et al. (Kim et al., 2021) show that, for this map, the Lipschitz constant with respect to any ℓ_p norm is infinite. Motivated by this negative result, Qi et al. (2023) propose *scaled cosine similarity attention* (SCSA) as a Lipschitz-controlled variant of self-attention, which normalizes the projections $\mathbf{X}\mathbf{W}^Q$ and $\mathbf{X}\mathbf{W}^K$ to unit norm and applies an inverse temperature coefficient $\tau > 0$:

$$\text{SCSA}(\mathbf{X}; \mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V, \nu, \tau) = \nu \text{softmax}\left(\tau \frac{\mathbf{X}\mathbf{W}^Q(\mathbf{X}\mathbf{W}^K)^\top}{\|\mathbf{X}\mathbf{W}^Q\| \|\mathbf{X}\mathbf{W}^K\|}\right)\mathbf{V}. \quad (6)$$

SCSA can be viewed as a Lipschitz variant of the standard attention mechanism. Theorem 1 in Qi et al. (2023) establishes an Lipschitz bound for scaled cosine similarity attention (SCSA), but their bound depends on the choice $(n-1)/n$ as a Lipschitz constant for the softmax operator. Using our sharper analysis of the softmax Lipschitz constant, we can replace $(n-1)/n$ by $1/2$ in their argument. Thus Lipschitz bound for SCSA can be tightened as shown below.

Theorem 2 (Refined ℓ_2 -Lipschitz bound for SCSA). *Let $L_2(\text{SCSA})$ denote the global Lipschitz constant of $\mathbf{X} \mapsto \text{SCSA}(\mathbf{X}, \mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V, \nu, \tau)$ with respect to the ℓ_2 norm. Then, for fixed $(\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V, \nu, \tau)$,*

$$L_2(\text{SCSA}) \leq n^2 \nu \tau \varepsilon^{-1/2} \|\mathbf{W}^K\|_2 + n \nu \tau \varepsilon^{-1/2} \|\mathbf{W}^Q\|_2 + 2n \nu \varepsilon^{-1/2} \|\mathbf{W}^V\|_2.$$

In particular, compared with the Lipschitz bound stated in Theorem 1 of Qi et al. (2023), Theorem 2 improves the estimate by a factor of 2, removing the extra factor 2 that appears in their original bound. Such refinements can directly impact theoretical robustness guarantees for attention-based architectures.

4.3 Entropy-regularized zero-sum games via double-softmax fixed point

Consider a two-player zero-sum matrix game with a payoff matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$. The standard minimax formulation is

$$\min_{\mathbf{x} \in \Delta_n} \max_{\mathbf{y} \in \Delta_m} \mathbf{x}^\top \mathbf{A} \mathbf{y},$$

where Δ_n, Δ_m denote the probability simplices. This formulation characterizes the Nash equilibrium of the game, capturing the row player’s strategy that minimizes the worst-case loss against an adversarial opponent, and thus serves as a foundation for robust decision making, reinforcement learning, and adversarial training. To ensure unique mixed strategies and to improve algorithmic stability, it is common to add entropic regularization terms (Cen et al., 2021). The entropy-regularized problem becomes

$$\min_{\mathbf{x} \in \Delta_n} \max_{\mathbf{y} \in \Delta_m} \mathbf{x}^\top \mathbf{A} \mathbf{y} + \tau (H(\mathbf{x}) - H(\mathbf{y})),$$

where, $\tau > 0$ is the regularization parameter and $H(\cdot)$ is the Shannon entropy, defined for any probability vector $\mathbf{x} \in \Delta_n$ as $H(\mathbf{x}) = -\sum_{i=1}^n x_i \log x_i$.

A well-known algorithm to solve the above entropy regularization problem is the double-softmax fixed-point (DSFP) iteration (McKelvey & Palfrey, 1995), which updates

$$\mathbf{y}_{k+1} \leftarrow (1 - \alpha) \mathbf{y}_k + \alpha \sigma_{1/\tau} \left(\mathbf{A}^\top \sigma_{1/\tau} (-\mathbf{A} \mathbf{y}_k) \right).$$

Intuitively, each player responds via a softmax update, yielding an equilibrium. Our improved Lipschitz constant for softmax in Theorem 1 provides a tighter condition to choose the regularization parameter τ depending on the payoff matrix \mathbf{A} , such that the DSFP iteration is contractive and convergent.

Theorem 3 (Convergence of DSFP under entropy regularization). *Let $\mathbf{A} \in \mathbb{R}^{n \times m}$ be a payoff matrix and let $\tau > 0$ denote the entropy regularization parameter. For all $1 \leq p \leq \infty$, if $\tau > \|\mathbf{A}\|_p / 2$ and $\alpha \in (0, 1]$, then the DSFP iteration is a Banach contraction on $(\Delta_m, \|\cdot\|_p)$ with contraction factor $\|\mathbf{A}\|_p^2 / (4\tau^2)$. Hence, the iterates \mathbf{y}_k converge linearly in the ℓ_p norm to the unique fixed point $\mathbf{y}^* \in \Delta_m$. The corresponding strategy of the row player is recovered as*

$$\mathbf{x}^* = \sigma_{1/\tau} (-\mathbf{A} \mathbf{y}^*),$$

which is also convergent since $\mathbf{x}_k = \sigma_{1/\tau} (-\mathbf{A} \mathbf{y}_k)$ at each step.

In summary, for DSFP iterations to converge linearly under the ℓ_p norm, we can derive a condition for the regularization parameter τ to satisfy and thereby guarantee well-defined equilibrium strategies.

5 Empirical Validation of Lipschitz Constant

In this section, we empirically validate our theoretical result that the softmax operator has a global Lipschitz constant of $1/2$, which is norm-uniform and tight. Our experiments involve a broad range of architectures across vision, language, and reinforcement learning, under varying datasets, prompts, and inverse temperature coefficient λ . All experiments were conducted on NVIDIA A40 GPUs. We compute the empirical Lipschitz constant using the definition in equation 3, using pairs of input obtained from each experimental setting. In particular, we compute,

$$\text{Empirical } L_p = \max_{1 \leq i \leq M} \frac{\|\sigma_\lambda(\mathbf{x}_i) - \sigma_\lambda(\mathbf{x}_i + \delta \mathbf{x}_i)\|_p}{\|\delta \mathbf{x}_i\|_p}, \quad (7)$$

where M is the number of inputs. For each input sample \mathbf{x}_i , we apply small random perturbations $\delta \mathbf{x}_i$ normalized to have $\|\delta \mathbf{x}_i\|_p = \epsilon$, where ϵ denotes the perturbation magnitude. We report the computed Empirical L_p constant across multiple values of ϵ and p .

5.1 Vision Models

We empirically evaluate the Lipschitz constant of the softmax operator, which appears in the attention mechanism of vision models. In particular, we consider Vision Transformer (ViT) models (Dosovitskiy et al., 2021) in three variants, Base (86M parameters), Large (307M parameters), and Huge (632M parameters). For each model, we select images from both the CIFAR-100 and ImageNet datasets, pass them through the transformer, and extract the pre-softmax attention scores $\mathbf{QK}^\top / \sqrt{d_k}$, which is the input to the softmax in Eq. 5. We then apply perturbations to these scores and compute the empirical Lipschitz constant row-wise according to Definition 7. M is set to be 100 images.

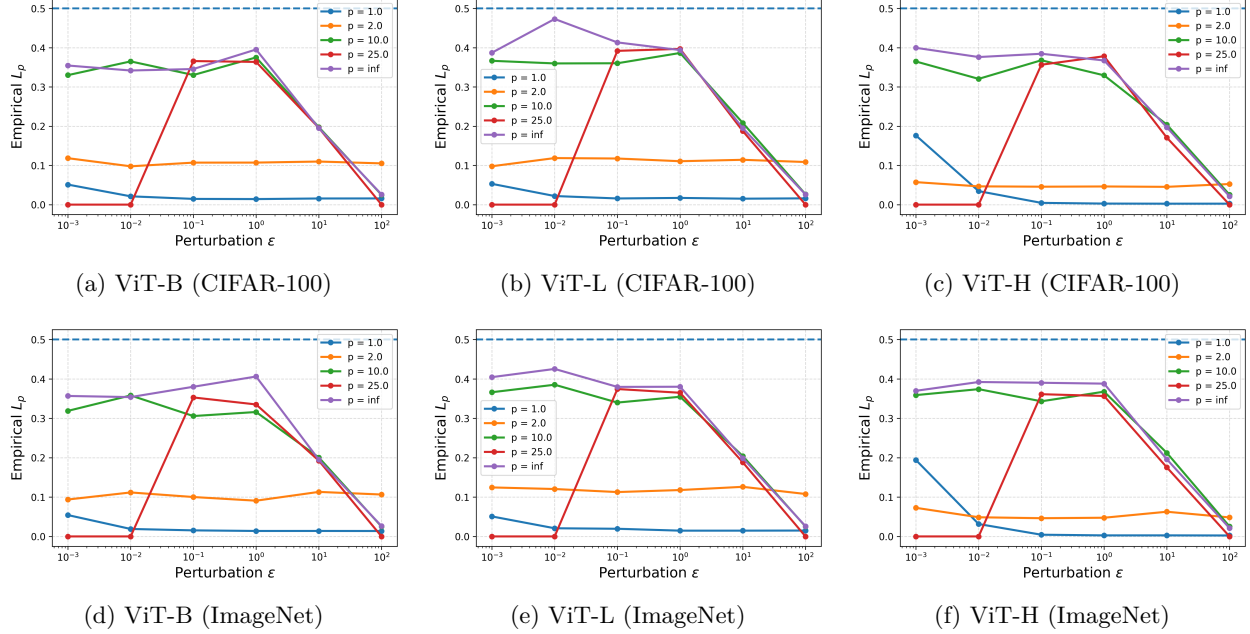


Figure 1: Empirical L_p of the softmax operator over attention scores from three Vision Transformer (ViT) variants on CIFAR-100 (top row) and ImageNet (bottom row), across varying perturbation magnitudes ϵ for multiple ℓ_p norms. In all cases, the empirical values remain below the derived bound of $1/2$.

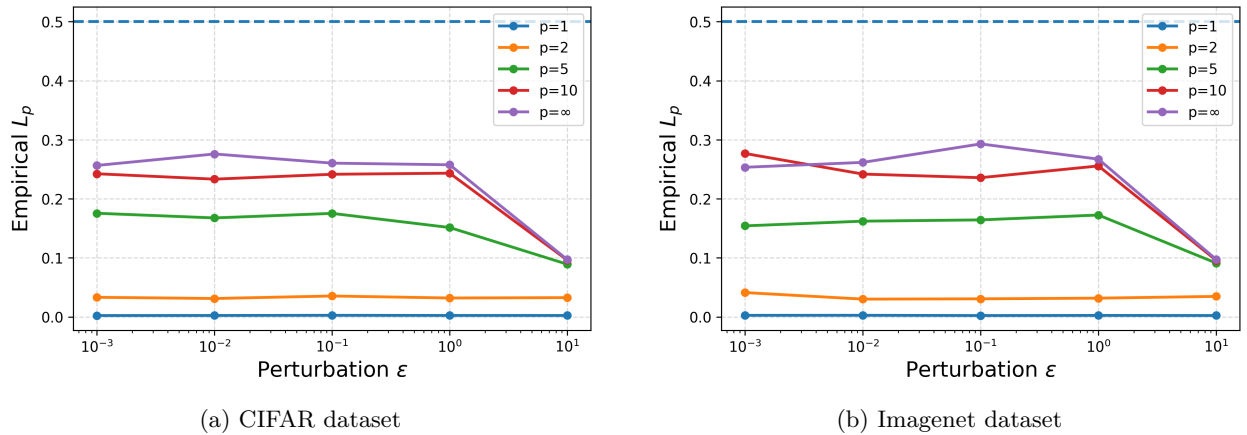


Figure 2: Empirical L_p of the softmax operator for classification logits of ResNET50 model on (a) CIFAR-100 and (b) ImageNet dataset, across varying perturbation magnitudes ϵ for multiple ℓ_p norms. In all cases, the empirical values remain well below the derived bound of $1/2$.

The results are presented as plots of empirical Lipschitz values versus the perturbation norm in Fig. 1 for CIFAR-100 (top row) and ImageNet (bottom row) datasets, with the theoretical baseline of 0.5 indicated. We report results across multiple ℓ_p norms. In each case, the plotted value corresponds to the maximum empirical Lipschitz constant over all attention heads, layers, tokens, and images. In particular, the total number of tested vectors is significant, which is equal to $N \times H \times M \times 12$, where N is the number of tokens and H the number of heads. The plots consistently show that the empirical Lipschitz constant remains below our theoretical value of $1/2$, regardless of the choice of norm or ViT variant.

We further evaluate the Lipschitz constant of the softmax operator on the classification logits of ResNET50 model (He et al., 2016), using images from the CIFAR-100 dataset in Fig. 2a and ImageNet dataset in Fig. 2b. In this experiment, we use $M = 25,000$ images and introduce input perturbations to the logits and compute the empirical L_p constant, for different p and perturbation magnitude ϵ . Consistent with our findings for attention scores, the empirical L_p remains strictly below the theoretical bound of $1/2$.

5.2 Language Models

We next evaluate the Lipschitz constant of the softmax operator on attention scores of large language models. Specifically, we sample 100 prompts each from the HellaSwag (Zellers et al., 2019) and the PIQA (Bisk et al., 2020) datasets, process them through GPT-2 (Radford et al., 2019) (518M parameters) and Qwen3 (Bai et al., 2023) (8B parameters), extract the pre-softmax attention scores, apply perturbations, and compute the empirical Lipschitz constant as defined in Definition 7.

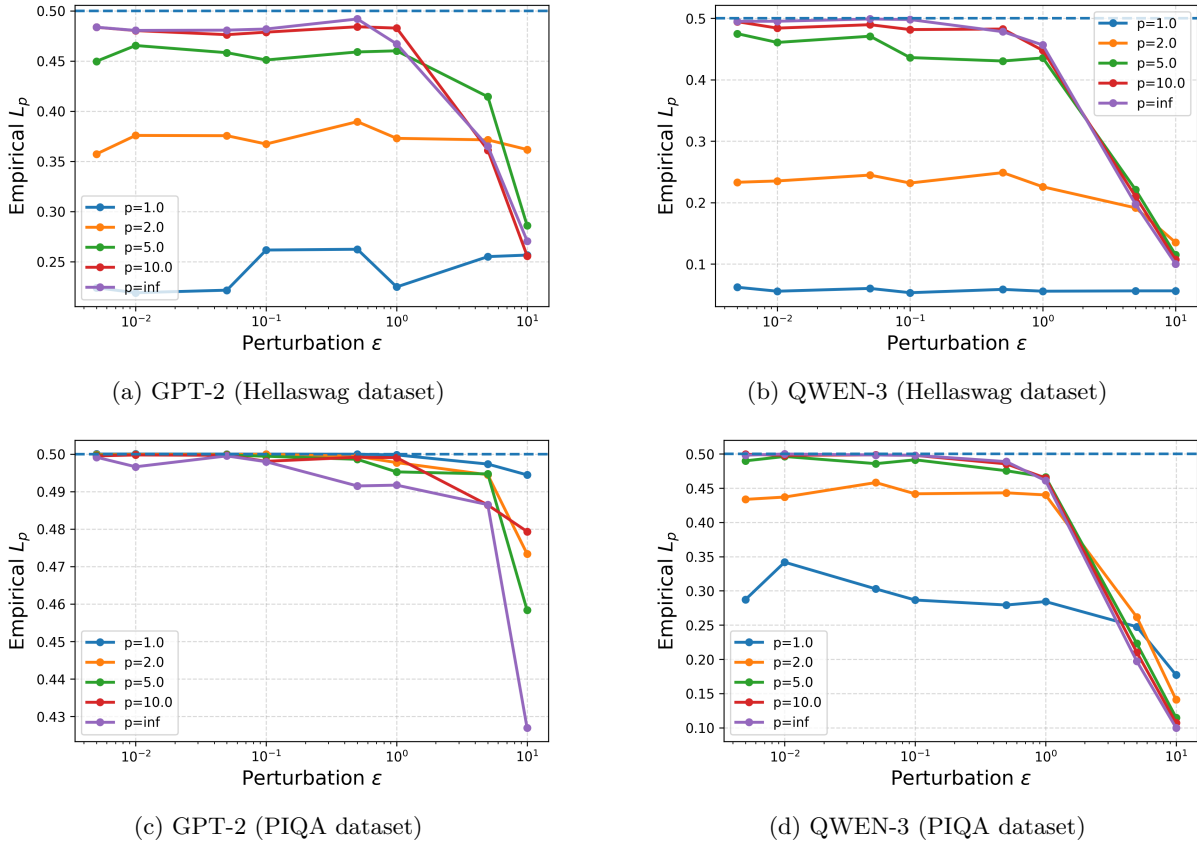


Figure 3: Empirical L_p of the softmax operator over attention scores from two Large Language models, GPT-2 and QWEN-3, on HellaSwag dataset (top row) and PIQA dataset (bottom row), across varying perturbation magnitudes ϵ for multiple ℓ_p norms. Across all configurations, the empirical values remain below the theoretical bound of $1/2$, with several instances approaching this limit, thereby confirming the tightness of the derived bound.

The scale of this evaluation is substantial here as well. The total number of tested vectors equals $H \times 100 \times 12$ per token per prompt, where the number of tokens varies with prompt length. Consistent with our experiments on vision models, the results in Fig. 3 for the Hellaswag dataset (top row) and the PIQA dataset (bottom row) demonstrate that the empirical Lipschitz constant remains strictly below the theoretical value of $1/2$ across all norms. Notably, for the PIQA dataset, we observe empirical constants of 0.4995 for the Qwen3 model and 0.4999 for GPT-2, providing empirical evidence for the tightness of our theoretically derived Lipschitz bound.

5.3 Reinforcement Learning Policies

We next investigate the Lipschitz behavior of softmax within stochastic policies in reinforcement learning. In policy-gradient methods (Sutton & Barto, 2018), a stochastic policy is typically parameterized as

$$\pi_\lambda(a | s) = \frac{\exp(\lambda Q(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\lambda Q(s, a'))},$$

where $Q(s, a)$ are the action logits for state $s \in \mathcal{S}$, \mathcal{A} is the action space, and $\lambda > 0$ is the inverse temperature coefficient. We evaluate this mapping $Q(s, \cdot) \mapsto \pi_\lambda(\cdot | s)$ in the **LunarLander** environment (discrete action space of size $|\mathcal{A}| = 4$) and the **CartPole** environment ($|\mathcal{A}| = 2$) from the OpenAI Gym benchmark (Brockman et al., 2016). We implement a PPO (Proximal Policy Optimization) agent (Schulman et al., 2017) using the **stable_baselines3** library, initialized randomly, and executed to collect a large set of visited states. For each state, we perturb the corresponding action logits $Q(s, \cdot)$ with random perturbations, and compute the empirical Lipschitz constant as in Definition 7. The reported value is obtained by averaging over multiple states ($M = 25000$) and perturbation trials for different coefficients λ . The results in Fig. 4 for Cartpole and Lunarlander environments confirm that the empirical constants scale proportionally to $\lambda/2$, in agreement with our theoretical prediction. Note that for the **CartPole** environment, where the softmax input dimension is $n = 2$, we obtain an empirical constant of 1.9996 for $p = 8$ and $\lambda = 4.0$, which is very close to the derived bound.

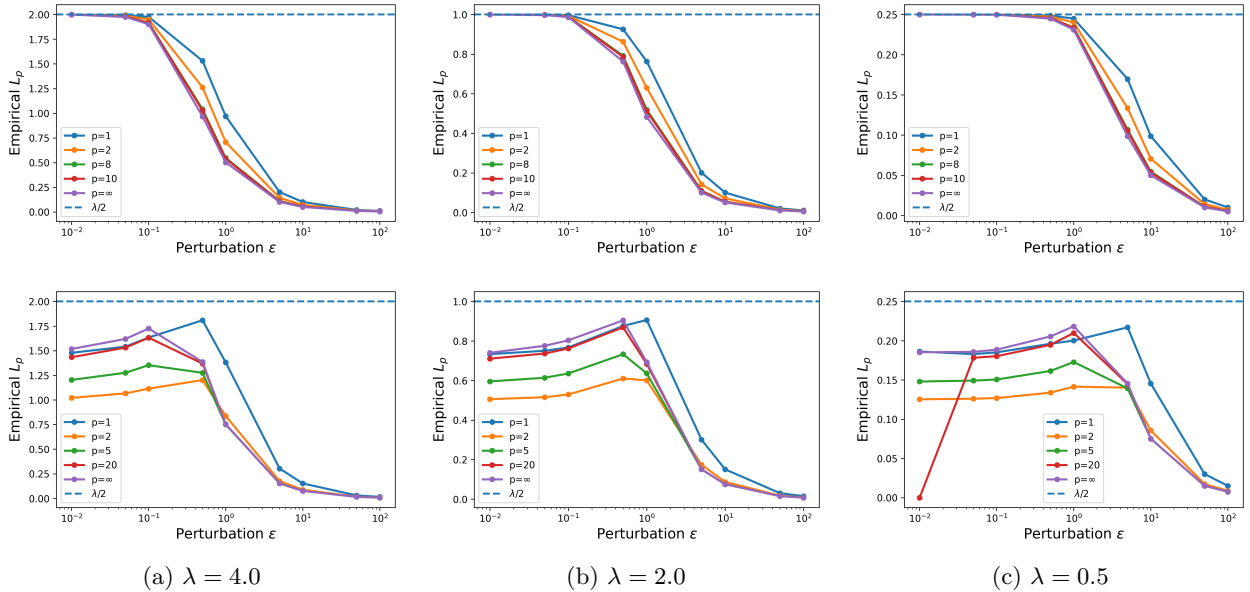


Figure 4: Empirical Lipschitz sensitivity of the softmax policy in RL environments, Cartpole (top row) and Lunarlander (bottom row), across varying perturbation magnitudes ϵ , under varying coefficients λ for different p -norms. The empirical values scale with λ and remain below the theoretical bound of $\lambda/2$, thereby empirically confirming the derived bound.

In summary, across all models, datasets, and settings considered, the empirical Lipschitz constant of softmax never exceeded $1/2$, thereby providing strong empirical validation of our theoretical bound. Furthermore, in

several cases, the empirical estimates approached $1/2$, illustrating that the derived bound is indeed tight in practice.

Conclusion

In this work, we derived that the softmax function admits a Lipschitz constant of $1/2$ uniformly across all ℓ_p norms. This result clarifies a commonly adopted assumption in the machine learning literature, where the Lipschitz constant of the softmax operator has often been considered as 1 with respect to the ℓ_2 norm. We further prove the tightness of our derived bound and demonstrate its utility in strengthening existing theoretical analyses. We also estimate the empirical Lipschitz constants for softmax transformations applied to attention score matrices from transformer-based architectures such as ViT, GPT-2, and Qwen3, classification logits from ResNet-50, and stochastic policies in reinforcement learning agents. In all cases, the empirical estimates remain strictly below the theoretical bound of $1/2$, thereby guaranteeing the tightness and generality of our results.

Broader Impact Statement

This work is primarily theoretical. By itself, this result does not introduce new datasets, architectures, or applications, but it provides sharper analytical tools for reasoning about the stability and robustness of existing models. On the positive side, such tools can help practitioners design systems with better controlled sensitivity to perturbations, complementing existing Lipschitz-based approaches used to improve robustness and training stability in deep networks (Cisse et al., 2017)(Miyato et al., 2018). At the same time, any advance that makes robust and stable deployment easier can be dual-use - it may also strengthen models used in societally sensitive settings, including surveillance and biometric monitoring (Almeida et al., 2022). A further risk is overstatement or misinterpretation - guarantees for a single component can be misconstrued as end-to-end assurances of safety, fairness, or reliability, which could encourage premature use in high-stakes contexts (AI, 2023).

References

- NIST AI. Artificial intelligence risk management framework (ai rmf 1.0). URL: <https://nvlpubs.nist.gov/nistpubs/ai/nist.ai>, pp. 100–1, 2023.
- Wael Alghamdi, Hsiang Hsu, Haewon Jeong, Hao Wang, Peter Michalak, Shahab Asodeh, and Flavio Calmon. Beyond adult and compas: Fair multi-class prediction via information projection. *Proc. Advances in Neural Information Processing Systems*, 35:38747–38760, 2022.
- Denise Almeida, Konstantin Shmarko, and Elizabeth Lomas. The ethics of facial recognition technologies, surveillance, and accountability in an age of artificial intelligence: a comparative analysis of us, eu, and uk regulatory frameworks. *AI and Ethics*, 2(3):377–387, 2022.
- Amir R Asadi and Emmanuel Abbe. Chaining meets chain rule: Multilevel entropic regularization and training of neural networks. *Journal of Machine Learning Research*, 21(139):1–32, 2020.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Peter L. Bartlett, Dylan J. Foster, and Matus Telgarsky. Spectrally-normalized margin bounds for neural networks. *Proc. Advances in Neural Information Processing Systems*, 30, 2017.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language. *Proc. AAAI Conference on Artificial Intelligence (AAAI)*, 34(05):7432–7439, 2020.
- Stephen P Boyd and Lieven Vandenbergh. *Convex optimization*. Cambridge University Press, 2004.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.

- Shicong Cen, Yuting Wei, and Yuejie Chi. Fast policy extragradient methods for competitive games with entropy regularization. *Proc. Advances in Neural Information Processing Systems*, 34:27952–27964, 2021.
- Lin Chen, Huaian Chen, Zhixiang Wei, Xin Jin, Xiao Tan, Yi Jin, and Enhong Chen. Reusing the task-specific classifier as a discriminator: Discriminator-free adversarial domain adaptation. *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7181–7190, 2022.
- Moustapha Cisse, Piotr Bojanowski, Edouard Grave, Yann Dauphin, and Nicolas Usunier. Parseval networks: Improving robustness to adversarial examples. *Proc. International Conference on Machine Learning*, 70: 854–863, 2017.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *Proc. International Conference on Learning Representations*, 2021.
- Mahyar Fazlyab, Manfred Morari, and George J. Pappas. Efficient and accurate estimation of lipschitz constants for deep neural networks. *Proc. Advances in Neural Information Processing Systems*, 32, 2019.
- 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.
- Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- Henry Gouk, Eibe Frank, Bernhard Pfahringer, and Michael J Cree. Regularisation of neural networks by enforcing lipschitz continuity. *Machine Learning*, 110(2):393–416, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Tuomas Hytönen, Jan Van Neerven, Mark Veraar, and Lutz Weis. *Analysis in Banach spaces*, volume 12. Springer, 2016.
- Hong Jun Jeon, Jason D. Lee, Qi Lei, and Benjamin Van Roy. An information-theoretic analysis of in-context learning. *Proc. International Conference on Machine Learning*, 2024.
- Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *ACM computing surveys (CSUR)*, 54(10s):1–41, 2022.
- Hyunjik Kim, George Papamakarios, and Andriy Mnih. The lipschitz constant of self-attention. *Proc. International Conference on Machine Learning*, 139:5562–5571, 2021.
- Anirban Laha, Saneem Ahmed Chemmengath, Priyanka Agrawal, Mitesh Khapra, Karthik Sankaranarayanan, and Harish G Ramaswamy. On controllable sparse alternatives to softmax. *Proc. Advances in Neural Information Processing Systems*, 31, 2018.
- Rémi Le Priol, Reza Babanezhad, Yoshua Bengio, and Simon Lacoste-Julien. An analysis of the adaptation speed of causal models. *Proc. International Conference on Artificial Intelligence and Statistics*, pp. 775–783, 2021.
- Bowen Lei, Dongkuan Xu, Ruqi Zhang, and Bani Mallick. Embracing unknown step by step: Towards reliable sparse training in real world. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.
- Wenzhe Li, Hao Luo, Zichuan Lin, Chongjie Zhang, Zongqing Lu, and Deheng Ye. A survey on transformers in reinforcement learning. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856.
- Richard D. McKelvey and Thomas R. Palfrey. Quantal response equilibria for normal form games. *Games and Economic Behavior*, 10(1):6–38, 1995.

- Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *Proc. International Conference on Learning Representations*, 2018.
- Ofir Nachum, Mohammad Norouzi, Kelvin Xu, and Dale Schuurmans. Bridging the gap between value and policy based reinforcement learning. *Proc. Advances in neural information processing systems*, 30, 2017.
- Thomas Nagler. Statistical foundations of prior-data fitted networks. *Proc. International Conference on Machine Learning*, pp. 25660–25676, 2023.
- Yu Nesterov. Smooth minimization of non-smooth functions. *Mathematical programming*, 103(1):127–152, 2005.
- Laker Newhouse. Softmax is $\frac{1}{2}$ -lipschitz (in a norm that may not matter). <https://www.lakernewhouse.com/assets/writing/softmax-is-0-5-lipschitz.pdf>, 2025. Unpublished note.
- Xianbiao Qi, Jianan Wang, Yihao Chen, Yukai Shi, and Lei Zhang. Lipsformer: Introducing lipschitz continuity to vision transformers. *Proc. International Conference on Learning Representations*, 2023.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9, 2019.
- Marcel Riesz. Sur les maxima des formes bilinéaires et sur les fonctionnelles linéaires. *Acta Mathematica*, 49(3):465–497, 1927.
- R. Tyrrell Rockafellar and Roger J-B Wets. *Variational Analysis*. Springer, 1998.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Elias M Stein. Interpolation of linear operators. *Transactions of the American Mathematical Society*, 83(2):482–492, 1956.
- Richard Sutton and Andrew Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.
- GO Thorin. An extension of a convexity theorem due to m. riesz, kungl. fysiografiska sällskapet i lund förhandlingar, no. 8, 14 (1938), 166–170; medd. *Lunds Univ. Mat. Sem.*, 4:1–5, 1939.
- Lewis Tunstall, Leandro Von Werra, and Thomas Wolf. *Natural language processing with transformers*. "O'Reilly Media, Inc.", 2022.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Proc. Advances in neural information processing systems*, 30, 2017.
- Aladin Virmaux and Kevin Scaman. Lipschitz regularity of deep neural networks: Analysis and efficient estimation. *Proc. Advances in Neural Information Processing Systems*, 31, 2018.
- Xiaojun Xu, Linyi Li, Yu Cheng, Subhabrata Mukherjee, Ahmed Hassan Awadallah, and Bo Li. Certifiably robust transformers with 1-lipschitz self-attention. 2022.
- Nikolay Yudin, Alexander Gaponov, Sergei Kudriashov, and Maxim Rakhuba. Pay attention to attention distribution: A new local lipschitz bound for transformers. *arXiv preprint arXiv:2507.07814*, 2025.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *Proc. Annual Meeting of the Association for Computational Linguistics*, pp. 4791–4800, 2019.

A Appendix

A.1 Proof of Proposition 1

(a) and (b) follow from standard results in matrix analysis (Golub & Van Loan, 2013).

Proof of Proposition 1(c). Let $\mathbf{A} \in \mathbb{R}^{n \times n}$. For $p = 1$ and $p = \infty$, the identity trivially satisfies, and hence fix $1 < p < \infty$. For any $\mathbf{x} \in \mathbb{R}^n$,

$$\|\mathbf{Ax}\|_p^p = \sum_{i=1}^n |(\mathbf{Ax})_i|^p,$$

where $(\mathbf{Ax})_i$ denotes i^{th} element in \mathbf{Ax} . Then, for each i , the triangle inequality gives

$$|(\mathbf{Ax})_i| = \left| \sum_{j=1}^n A_{ij} x_j \right| \leq \sum_{j=1}^n |A_{ij}| |x_j|$$

By Holder's inequality, we obtain for each i ,

$$\sum_{j=1}^n |A_{ij}| |x_j| \leq \left(\sum_{j=1}^n |A_{ij}| \right)^{1-1/p} \left(\sum_{j=1}^n |A_{ij}| |x_j|^p \right)^{1/p}.$$

Thus, we can bound $\|\mathbf{Ax}\|_p^p$ as follows,

$$\begin{aligned} \|\mathbf{Ax}\|_p^p &\leq \sum_{i=1}^n \left(\sum_{j=1}^n |A_{ij}| \right)^{p-1} \left(\sum_{j=1}^n |A_{ij}| |x_j|^p \right) \\ &\leq \left(\max_{1 \leq i \leq n} \sum_{j=1}^n |A_{ij}| \right)^{p-1} \sum_{i=1}^n \sum_{j=1}^n |A_{ij}| |x_j|^p \\ &= \|\mathbf{A}\|_\infty^{p-1} \sum_{j=1}^n |x_j|^p \left(\sum_{i=1}^n |A_{ij}| \right) \\ &\leq \|\mathbf{A}\|_\infty^{p-1} \left(\max_{1 \leq j \leq n} \sum_{i=1}^n |A_{ij}| \right) \sum_{j=1}^n |x_j|^p \\ &= \|\mathbf{A}\|_\infty^{p-1} \|\mathbf{A}\|_1 \|x\|_p^p. \end{aligned}$$

The second and fourth inequalities are obtained by taking the max term out of the summation¹, i.e,

$$\sum_{j=1}^n |A_{ij}| \leq \max_{1 \leq i \leq n} \sum_{j=1}^n |A_{ij}| \text{ for all } i, \quad \sum_{i=1}^n |A_{ij}| \leq \max_{1 \leq j \leq n} \sum_{i=1}^n |A_{ij}| \text{ for all } j \quad (8)$$

The third and fifth equalities follow from (a) and (b). Taking the p -th root and maximizing over $\|\mathbf{x}\|_p = 1$, we get the inequality

$$\|\mathbf{A}\|_p \leq \|\mathbf{A}\|_\infty^{1-1/p} \|\mathbf{A}\|_1^{1/p}.$$

□

¹This inequality will be used in the proof of Proposition 2 as well.

A.2 Proof of Theorem 1

Proof of Theorem 1(a). Note that $\mathbf{J}_{\sigma_1}(\mathbf{x})$ is symmetric for all $\mathbf{x} \in \mathbb{R}^n$. Hence, from Proposition 1(a) and (b), we get that for all $\mathbf{x} \in \mathbb{R}^n$,

$$\|\mathbf{J}_{\sigma_1}(\mathbf{x})\|_1 = \|\mathbf{J}_{\sigma_1}(\mathbf{x})\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |[\mathbf{J}_{\sigma_1}(\mathbf{x})]_{ij}|$$

From Lemma 2, the above optimization problem becomes,

$$\begin{aligned} & \max_{1 \leq i \leq n} s_i(1 - s_i) + \sum_{j \neq i} s_i s_j \\ &= \max_{1 \leq i \leq n} s_i(1 - s_i) - s_i^2 + \sum_{j=1}^n s_i s_j \\ &= \max_{1 \leq i \leq n} 2s_i(1 - s_i), \end{aligned}$$

where $\mathbf{s} = \sigma_1(\mathbf{x}) = (s_1, s_2, \dots, s_n)$. The optimal value for the above problem is upper bounded by $1/2$, and $1/2$ is attained for all $\mathbf{x} \in \mathbb{R}^n$ such that the corresponding output $\mathbf{s} = \sigma_1(\mathbf{x})$ belongs to Δ_n^0 and $s_i = 1/2$ for any $i \in \{1, 2, \dots, n\}$. Thus, we get the following bounds,

$$\|\mathbf{J}_{\sigma_1}(\mathbf{x})\|_1 \leq \frac{1}{2} \quad \text{and} \quad \|\mathbf{J}_{\sigma_1}(\mathbf{x})\|_\infty \leq \frac{1}{2}$$

By Proposition 1(c), we get

$$\|\mathbf{J}_{\sigma_1}(\mathbf{x})\|_p \leq \left(\frac{1}{2}\right)^{\frac{1}{p}} \left(\frac{1}{2}\right)^{1-\frac{1}{p}} = \frac{1}{2}, \quad \text{for all } 1 \leq p \leq \infty.$$

□

Proof of Theorem 1(b). From Theorem 1(a), we get,

$$\sup_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{J}_{\sigma_1}(\mathbf{x})\|_p \leq \frac{1}{2}, \quad \text{for all } 1 \leq p \leq \infty.$$

From Lemma 3(b), for any $\lambda > 0$, we have,

$$\sup_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{J}_{\sigma_\lambda}(\mathbf{x})\|_p \leq \frac{\lambda}{2} \quad \text{for all } 1 < p < \infty.$$

and hence using Lipschitz characterization in Lemma 1,

$$\|\sigma_\lambda(\mathbf{x}) - \sigma_\lambda(\mathbf{y})\|_p \leq \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|_p \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{R}^n \text{ and for all } 1 \leq p \leq \infty.$$

□

A.3 Proof of Proposition 2

Define $\mathbf{M}(\mathbf{s}) \in \mathbb{R}^{n \times n}$ as $\text{diag}(\mathbf{s}) - \mathbf{s}\mathbf{s}^T$. From Lemma 3, for all $1 \leq p \leq \infty$,

$$\sup_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{J}_{\sigma_1}(\mathbf{x})\|_p = \sup_{\mathbf{s} \in \Delta_n^0} \|\mathbf{M}(\mathbf{s})\|_p$$

Proof of Proposition 2(a). As discussed in Theorem 1(a), for $p = 1$ and $p = \infty$, for all $\mathbf{s} \in \Delta_n^0$,

$$\|\mathbf{M}(\mathbf{s})\|_1 = \|\mathbf{M}(\mathbf{s})\|_\infty = \max_{1 \leq i \leq n} 2s_i(1 - s_i) \leq 1/2$$

and $1/2$ is attained if there exists $s \in \Delta_n^\circ$ such that $s_i = 1/2$ for some i and $s_j > 0$ for $j \neq i$.

To prove that $1/2$ is indeed attained, we need to show an example where $1/2$ is attained at a point \mathbf{s} in the interior of the probability simplex Δ_n° .

Let $\mathbf{x} \in \mathbb{R}^n$ such that $x_i = \ln(n-1)$ and $x_j = 0$ for all $j \neq i$, $\mathbf{s} = \sigma_1(\mathbf{x})$ is such that $s_i = 1/2$ and $s_j > 0$ for all $j \neq i$ and hence $\mathbf{s} \in \Delta_n^\circ$. Thus,

$$\max_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{J}_{\sigma_1}(\mathbf{x})\|_1 = \max_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{J}_{\sigma_1}(\mathbf{x})\|_\infty = 1/2$$

□

Proof of Proposition 2(b). We first prove for the case $n > 2$. It suffices to show that for arbitrary $1 < p < \infty$, $\|\mathbf{M}(\mathbf{s})\|_p \neq 1/2$ for any $\mathbf{s} \in \Delta_n^\circ$ and there exists $\mathbf{s} \in \partial\Delta_n$ such that $\|\mathbf{M}(\mathbf{s})\|_p = 1/2$ for $n > 2$ to prove that $1/2$ is indeed approached in limit. We need the following definition of the support of a vector and the corresponding lemma to prove this.

Definition A.1 (Support of a vector). For a vector $\mathbf{s} \in \mathbb{R}^n$, the support of \mathbf{s} is defined as

$$\text{supp}(\mathbf{s}) := \{i \in \{1, \dots, n\} : s_i \neq 0\}.$$

Lemma 4 (Support condition for attaining $1/2$). Let $\mathbf{M}(\mathbf{s}) \in \mathbb{R}^{n \times n}$ be defined by

$$\mathbf{M}(\mathbf{s}) := \text{Diag}(\mathbf{s}) - \mathbf{s}\mathbf{s}^\top,$$

for $\mathbf{s} \in \Delta_n$. Suppose $1 < p < \infty$ and $n > 2$. Then, $\|\mathbf{M}(\mathbf{s})\|_p = \frac{1}{2}$ if and only if \mathbf{s} is a permutation of $(1/2, 1/2, 0, \dots, 0)$

The proof of the lemma is given towards the end of this section. From Lemma 4, there exists no $\mathbf{s} \in \Delta_n$ with $\text{supp}(\mathbf{s}) > 2$ with $\|\mathbf{M}(\mathbf{s})\|_p = 1/2$. Hence, there exists no $\mathbf{s} \in \Delta_n^\circ$ with $\|\mathbf{M}(\mathbf{s})\|_p = 1/2$ since $\text{supp}(\mathbf{s}) = n$ for all $\mathbf{s} \in \Delta_n^\circ$. Also, there exists \mathbf{s} with $\text{supp}(\mathbf{s}) = 2$ such that $\|\mathbf{M}(\mathbf{s})\|_p = 1/2$.

Let $\hat{\mathbf{s}} \in \Delta_n$ with $\text{supp}(\hat{\mathbf{s}}) = 2$ be such that $\|\mathbf{M}(\hat{\mathbf{s}})\|_p = 1/2$. By definition, $\hat{\mathbf{s}} \in \partial\Delta_n$ (boundary of the simplex). Hence, we can obtain a sequence of $\mathbf{s}_k \in \Delta_n^\circ$ (interior of the simplex) such that $\mathbf{s}_k \rightarrow \hat{\mathbf{s}}$ as $k \rightarrow \infty$ (Rockafellar & Wets, 1998). Also since $\mathbf{M}(\mathbf{s})$ is a continuous function, $\mathbf{M}(\mathbf{s}_k) \rightarrow 1/2$ as $k \rightarrow \infty$. For completeness, we show convergence by constructing example sequence in Example 2.

Finally, when $n = 2$, the point $\mathbf{s} = (1/2, 1/2)$ lies in Δ_2° and attains $\|\mathbf{M}(\mathbf{s})\|_p = 1/2$. □

Example 2 (Sequence $(\mathbf{s}_k)_{k \geq 1} \subset \Delta_n^\circ$ such that $\mathbf{M}(\mathbf{s}_k) \rightarrow 1/2$ as $k \rightarrow \infty$). Consider $\hat{\mathbf{s}} \in \partial\Delta_n$ with $\hat{s}_1 = \hat{s}_2 = 1/2$ and all other entries zero. By Lemma 4, $\|\mathbf{M}(\hat{\mathbf{s}})\|_p = 1/2$. To approximate $\hat{\mathbf{s}}$ from the interior of the probability simplex, fix $\epsilon \in (0, 1/2)$ and choose $\delta \in (0, 1/2)$ satisfying $2\delta(1 - \delta) = \epsilon$. Define

$$s_1 = s_2 = \frac{1}{2} - \delta, \quad s_j = \frac{2\delta}{n-2}, \quad j \geq 3,$$

so that $\mathbf{s} \in \Delta_n^\circ$. Consider the vector

$$\mathbf{v} = (2^{-1/p}, -2^{-1/p}, 0, \dots, 0) \in \mathbb{R}^n \quad \text{with} \quad \|\mathbf{v}\|_p = 1,$$

A direct computation shows,

$$[\mathbf{M}(\mathbf{s})\mathbf{v}]_1 = -[\mathbf{M}(\mathbf{s})\mathbf{v}]_2 = \frac{2}{2^{1/p}} \left(\frac{1}{2} - \delta\right)^2 \quad \text{and} \quad [\mathbf{M}(\mathbf{s})\mathbf{v}]_j = 0 \quad \text{for all } j \notin \{1, 2\}.$$

Hence, we get,

$$\|\mathbf{M}(\mathbf{s})\mathbf{v}\|_p = \frac{1}{2} - 2\delta(1 - \delta) = \frac{1}{2} - \epsilon.$$

Thus for every $\epsilon \in (0, 1/2)$, there exists $\mathbf{s} \in \Delta_n^\circ$ with $1/2 - \epsilon \leq \|\mathbf{M}(\mathbf{s})\|_p < 1/2$, proving that $1/2$ is the tight upper bound.

Proof of Lemma 4. Let \mathbf{s} be a permutation of $(1/2, 1/2, 0, \dots, 0)$ i.e $s_i = s_j = 1/2$ for some distinct $i, j \in \{1, 2, \dots, n\}$. $\mathbf{M}(\mathbf{s})$ has values $\mathbf{M}(\mathbf{s})_{ii} = \mathbf{M}(\mathbf{s})_{jj} = 1/4$, $\mathbf{M}(\mathbf{s})_{ij} = \mathbf{M}(\mathbf{s})_{ji} = -1/4$, and zeros elsewhere. To show that $\|\mathbf{M}(\mathbf{s})\|_p = 1/2$, it suffices to construct $\mathbf{v} \in \mathbb{R}^n$ such that $\|\mathbf{v}\|_p = 1$ and $\|\mathbf{M}(\mathbf{s})\mathbf{v}\|_p = 1/2$. Consider \mathbf{v} to be a similar permutation of

$$(2^{-1/p}, -2^{-1/p}, 0, \dots, 0) \in \mathbb{R}^n.$$

i.e. $v_i = 2^{-1/p}$ and $v_j = -2^{-1/p}$. A direct calculation gives $\|\mathbf{M}(\mathbf{s})\mathbf{v}\|_p = 1/2$.

Now we need to prove that if $\|\mathbf{M}(\mathbf{s})\|_p = 1/2$, \mathbf{s} is a permutation of $(1/2, 1/2, 0, \dots, 0)$. By Proposition 1(c), the interpolation inequality gives

$$\|\mathbf{M}(\mathbf{s})\|_p \leq \|\mathbf{M}(\mathbf{s})\|_1^{1/p} \|\mathbf{M}(\mathbf{s})\|_\infty^{1-1/p}.$$

From Proposition 2(a), we know $\|\mathbf{M}(\mathbf{s})\|_1 = \|\mathbf{M}(\mathbf{s})\|_\infty \leq 1/2$. Thus if $\|\mathbf{M}(\mathbf{s})\|_p = 1/2$, equality must hold throughout. This means for all $\mathbf{s} \in \Delta_n$ such that $\|\mathbf{M}(\mathbf{s})\|_p = 1/2$, we have $\|\mathbf{M}(\mathbf{s})\|_1 = \|\mathbf{M}(\mathbf{s})\|_\infty = 1/2$.

We divide into cases according to the support size of \mathbf{s} .

Case (i): $\text{supp}(\mathbf{s}) = 1$. Suppose $s_i = 1$ for some i , and $s_j = 0$ for all $j \neq i$. Then $\mathbf{M}(\mathbf{s}) = 0$. Hence $\|\mathbf{M}(\mathbf{s})\|_p = 0$ for all p . Thus, the value $1/2$ cannot be attained for any \mathbf{s} .

Case (ii): $\text{supp}(\mathbf{s}) = 2$. As discussed in Theorem 1(a), for $p = 1$ and $p = \infty$, for all $\mathbf{s} \in \Delta_n$,

$$\|\mathbf{M}(\mathbf{s})\|_1 = \|\mathbf{M}(\mathbf{s})\|_\infty = \max_{1 \leq i \leq n} 2s_i(1 - s_i) \leq 1/2$$

and $1/2$ is attained only if there exists $\mathbf{s} \in \Delta_n$ such that $s_i = 1/2$ for some i . Thus, if $\text{supp}(\mathbf{s}) = 2$, then \mathbf{s} needs to be permutations of $(1/2, 1/2, 0, \dots, 0)$ for $\|\mathbf{M}(\mathbf{s})\|_1 = \|\mathbf{M}(\mathbf{s})\|_\infty = 1/2$ and we have already shown $\|\mathbf{M}(\mathbf{s})\|_p = 1/2$ for any permutation of $(1/2, 1/2, 0, \dots, 0)$.

Thus, if $\|\mathbf{M}(\mathbf{s})\|_p = 1/2$ with $\text{supp}(\mathbf{s}) = 2$, then $\|\mathbf{M}(\mathbf{s})\|_1 = \|\mathbf{M}(\mathbf{s})\|_\infty = 1/2$ which needs \mathbf{s} to be a permutation of $(1/2, 1/2, 0, \dots, 0)$.

Case (iii): $\text{supp}(\mathbf{s}) > 2$. Suppose, $\hat{\mathbf{s}} \in \Delta_n$ with $\text{supp}(\hat{\mathbf{s}}) > 2$ satisfies $\|\mathbf{M}(\hat{\mathbf{s}})\|_p = 1/2$. By Proposition 1(c),

$$\|\mathbf{M}(\hat{\mathbf{s}})\|_p \leq \|\mathbf{M}(\hat{\mathbf{s}})\|_1^{1/p} \|\mathbf{M}(\hat{\mathbf{s}})\|_\infty^{1-1/p}.$$

Since $\|\mathbf{M}(\hat{\mathbf{s}})\|_1 = \|\mathbf{M}(\hat{\mathbf{s}})\|_\infty = 1/2$ holds, equality must hold in the interpolation inequality in Proposition 1. Thus, it suffices to rule out equality.

We just need to show that one of the inequalities used in the proof of Proposition 1(c) is strict. We show that the inequality in Eq. 8 is strict i.e. there exists $1 \leq i \leq n$, such that,

$$\sum_{j=1}^n |\mathbf{M}(\mathbf{s})_{ij}| < \max_{1 \leq k \leq n} \sum_{j=1}^n |\mathbf{M}(\mathbf{s})_{kj}|$$

Since $\|\mathbf{M}(\hat{\mathbf{s}})\|_1 = \|\mathbf{M}(\hat{\mathbf{s}})\|_\infty = 1/2$, $\hat{s}_{i_1} = 1/2$ for some $1 \leq i_1 \leq n$ and $\hat{s}_j < 1/2$ for all $j \neq i_1$. Since $\text{supp}(\hat{\mathbf{s}}) > 2$, there exists an index $i_2 \in \{1, 2, \dots, n\}$ and $i_2 \neq i_1$, where $\hat{s}_{i_2} < 1/2$. As discussed in Theorem 1, for any row index i , $\sum_{j=1}^n |\mathbf{M}(\hat{\mathbf{s}})_{ij}| = 2\hat{s}_i(1 - \hat{s}_i)$. Hence,

$$\max_{1 \leq i \leq n} \sum_{j=1}^n |\mathbf{M}(\mathbf{s})_{ij}| = \sum_{j=1}^n |\mathbf{M}(\mathbf{s})_{i_1 j}| = 1/2,$$

and for row index i_2 ,

$$\sum_{j=1}^n |\mathbf{M}(\mathbf{s})_{i_2 j}| < 1/2 = \max_{1 \leq k \leq n} \sum_{j=1}^n |\mathbf{M}(\mathbf{s})_{kj}|.$$

This shows that the inequality in Eq. 8 is strict. Hence, the interpolation inequality in Proposition 1(c) is strict for the matrix $\mathbf{M}(\hat{\mathbf{s}})$. This contradicts our assumption that $\|\mathbf{M}(\hat{\mathbf{s}})\|_p = 1/2$.

Combining all cases, we conclude that $\|\mathbf{M}(\mathbf{s})\|_p = 1/2$ if and only if \mathbf{s} is a permutation of $(1/2, 1/2, 0, \dots, 0)$. \square

A.4 Proof of Theorem 3

Let the DSFP map be defined as

$$T(\mathbf{y}) = \sigma_{1/\tau}(\mathbf{A}^\top \sigma_{1/\tau}(-\mathbf{A}\mathbf{y})), \quad \mathbf{y} \in \Delta_m,$$

where $\sigma_{1/\tau}$ denotes the softmax operator with temperature τ . By Theorem 1, the softmax operator is globally Lipschitz with constant $\frac{1}{2\tau}$, uniformly across all ℓ_p norms with $1 \leq p \leq \infty$. Applying the chain rule of Lipschitz constants, we obtain

$$\|T(\mathbf{y}_1) - T(\mathbf{y}_2)\|_p \leq \frac{1}{2\tau} \|\mathbf{A}^\top\|_p \frac{1}{2\tau} \|\mathbf{A}\|_p \|\mathbf{y}_1 - \mathbf{y}_2\|_p.$$

Since $\|\mathbf{A}^\top\|_p = \|\mathbf{A}\|_p$, this yields

$$\|T(\mathbf{y}_1) - T(\mathbf{y}_2)\|_p \leq \frac{\|\mathbf{A}\|_p^2}{4\tau^2} \|\mathbf{y}_1 - \mathbf{y}_2\|_p.$$

Hence, the DSFP map has the Lipschitz constant,

$$c = \frac{\|\mathbf{A}\|_p^2}{4\tau^2}.$$

If $\tau > \|\mathbf{A}\|_p/2$, then $c < 1$ and T is a Banach contraction on the complete metric space $(\Delta_m, \|\cdot\|_p)$. By the Banach fixed-point theorem, T then admits a unique fixed point $\mathbf{y}^* \in \Delta_m$, and for any initialization $\mathbf{y}_0 \in \Delta_m$, the iteration

$$\mathbf{y}_{k+1} \leftarrow (1 - \alpha) \mathbf{y}_k + \alpha T(\mathbf{y}_k), \quad \alpha \in (0, 1],$$

converges linearly in the ℓ_p norm to \mathbf{y}^* . At each step, the row player's strategy is updated as

$$\mathbf{x}_k = \sigma_{1/\tau}(-\mathbf{A}\mathbf{y}_k).$$

Since $\mathbf{y}_k \rightarrow \mathbf{y}^*$ and the softmax map is continuous, it follows that $\mathbf{x}_k \rightarrow \mathbf{x}^* := \sigma_{1/\tau}(-\mathbf{A}\mathbf{y}^*)$. Thus, both players' equilibrium strategies are recovered in the limit.