000 001 002

003 004

005

006 007

009 010

011

016 018 019

021

022 023

024

025

031

033

034 037

038

044 045 047

048 049 050

051 052

DETERMINISTIC BOUNDS AND RANDOM ESTIMATES OF METRIC TENSORS ON NEUROMANIFOLDS

Anonymous authors

Paper under double-blind review

ABSTRACT

The high dimensional parameter space of modern deep neural networks — the neuromanifold — is endowed with a unique metric tensor defined by the Fisher information, estimating which is crucial for both theory and practical methods in deep learning. To analyze this tensor for classification networks, we return to a low dimensional space of probability distributions — the core space — and carefully analyze the spectrum of its Riemannian metric. We extend our discoveries there into deterministic bounds of the metric tensor on the neuromanifold. We introduce an unbiased random estimate of the metric tensor and its bounds based on Hutchinson's trace estimator. It can be evaluated efficiently through a single backward pass, with a standard deviation bounded by the true value up to scaling.

Introduction

Deep learning can be considered as a trajectory through the space of neural networks (neuromanifold; Amari 2016), where each point is a neural network instance with a prescribed architecture but different parameters. This work investigates classifier models in the form $p(y \mid x, \theta)$, where x is the input features, $y \in \{1, \dots, C\}$ is the class labels $(C \ge 2)$, and $\theta \in \Theta$ is the network weights and biases. Given an unlabeled dataset $\mathcal{D}_x = \{x_1, x_2, \cdots\}$, the intrinsic structure of Θ is specified by the Fisher Information Matrix (FIM), defined as:

$$\mathcal{F}(\theta) \coloneqq \sum_{x \in \mathcal{D}_x} \mathbb{E}_{p(y \mid x)} \left[\frac{\partial \log p(y \mid x, \theta)}{\partial \theta} \frac{\partial \log p(y \mid x, \theta)}{\partial \theta^{\top}} \right] = \sum_{x \in \mathcal{D}_x} \mathbb{E}_{p(y \mid x)} \left[\frac{\partial \ell_{xy}}{\partial \theta} \frac{\partial \ell_{xy}}{\partial \theta^{\top}} \right], \quad (1)$$

where $\ell_{xy}(\theta) := \log p(y \mid x, \theta)$ denotes the log-likelihood. This is based on a supervised model $x \to y$. For unsupervised models, one can treat x as constant and apply the same formula. Under regularity conditions, $\mathcal{F}(\theta)$ is a $\dim(\theta) \times \dim(\theta)$ positive semi-definite (psd) matrix varying smoothly with $\theta \in \Theta$. Following Hotelling (1929), and independently Rao (1945), $\mathcal{F}(\theta)$ is used as a metric tensor on Θ , representing a local degenerate inner product. For example, one can measure the intrinsic squared distance between θ and $\theta + d\theta$, where $d\theta$ is a small dynamic on Θ , as $d\theta^{\top} \mathcal{F}(\theta) d\theta$.

The FIM is the unique metric tensor (Čencov, 1982) which underpins the information geometry of the neuromanifold Θ (Amari, 2016). The most widely used application of the FIM is perhaps geometryinspired optimizers such as natural gradient (Amari, 1998), Adam (Kingma & Ba, 2015), and their variants (Martens & Grosse, 2015; Pascanu & Bengio, 2014; Yao et al., 2021; Lin et al., 2021). \mathcal{F} is also applied to regularized fine-tuning (Lodha et al., 2023), transfer learning (Chen et al., 2018), and overcoming catastrophic forgetting (Kirkpatrick et al., 2017). Theoretically, the FIM provides insights due to its connection with the Hessian of the loss landscape and generalization (Hochreiter & Schmidhuber, 1997), and that any f-divergence is locally characterized by the FIM (Blyth, 1994).

Given its deep and broad background, estimating $\mathcal{F}(\theta)$ with some guaranteed quality, even without a specific application pipeline, is an important topic. As a widely used deterministic approximation, the empirical FIM (eFIM, a.k.a. empirical Fisher, see e.g. Le Roux et al. 2007) is given by $\hat{\mathcal{F}}(\theta) :=$ $\sum_{(x,y)\in\mathcal{D}}\left[\frac{\partial\ell_{xy}}{\partial\theta}\frac{\partial\ell_{xy}}{\partial\theta^{\top}}\right], \text{ where } \mathcal{D}=\left\{(x_1,y_1),(x_2,y_2),\cdots\right\} \text{ is a labeled dataset. As another example,}$ the Monte Carlo (MC) estimator $\hat{\mathcal{F}}(\theta)=\frac{1}{m}\sum_{\hat{x},\hat{y}}\frac{\partial\ell_{\hat{x}\hat{y}}}{\partial\theta}\frac{\partial\ell_{\hat{x}\hat{y}}}{\partial\theta^{\top}}, \text{ where } \hat{x},\,\hat{y} \text{ are a set of } m \text{ random}$

¹In the machine learning literature, $\mathcal{F}(\theta)$ is sometimes referred to as a curvature matrix (Martens, 2020) but actually defines a singular semi-Riemannian metric (Sun & Nielsen, 2025) in rigorous terms.

samples drawn from \mathcal{D}_x and $p(y \mid \hat{x})$, respectively, and the symbol $\hat{\mathcal{F}}(\theta)$ is abused for simplicity, gives an unbiased estimate of $\mathcal{F}(\theta)$ up to scaling.

We advance the state of the art in both deterministic and stochastic approaches to computing the FIM, improving accuracy in terms of bound gap and variance. We made the following contributions: ① Envelopes of the FIM in the statistical simplex (space of output probabilities); ② Deterministic bounds of the FIM for classifier networks and their tightness analysis; ③ A novel family of random FIM estimates based on Hutchinson's trick (Hutchinson, 1990; Skorski, 2021), which can be computed efficiently with bounded variance; ④ An empirical study to estimate the FIM of DistilBert (Sanh et al., 2019) to showcase the advantages of Hutchinson's estimate in production settings.

In the rest of this section, we introduce our notations. Section 2 develops fundamental bounds and estimates in low dimensional spaces of probability distributions. Section 3 extends the deterministic bounds into the high dimensional neuromanifold. Section 4 introduces Hutchinson's FIM estimator and discusses its theoretical properties with numerical simulation on DistilBERT (Sanh et al., 2019). Section 5 positions our work into the literature. Section 6 concludes.

NOTATIONS AND CONVENTIONS

We use lowercase letters such as λ or a for both vectors and scalars, which should be distinguished based on context, and capital letters such as A for matrices. All vectors are column vectors. A scalar-vector or vector-scalar derivative such as $\partial \ell/\partial \theta$ yields a gradient vector of the same shape as the vector. A vector-vector derivative such as $\partial z/\partial \theta$ denotes the $\dim(z) \times \dim(\theta)$ Jacobian matrix of the mapping $\theta \to z$. $\|\cdot\|$ denote the Euclidean norm for vectors or Frobenius norm for matrices. $\|\cdot\|_{\sigma}$ denotes the spectral norm (maximum singular value) of matrices. The metric tensors (variants of FIM) are listed in table 1.

Table 1: Metric tensors. Both empirical FIM (2nd column) and Monte Carlo FIM (3rd column) are denoted as $\hat{\mathcal{I}}$ / $\hat{\mathcal{F}}$ for reducing notation overload. We use \mathcal{I} / $\hat{\mathcal{I}}$ / \mathbb{I} for simple low-dimensional statistical manifolds and use \mathcal{F} / $\hat{\mathcal{F}}$ / \mathbb{F} for neuromanifolds. We optionally use superscripts to indicate the associated parameter space. For example, \mathcal{I}^{Δ} and \mathcal{F}^{Δ} denote the metric tensor of the statistical simplex and the space of neural networks with simplex-valued outputs, respectively.

| FIM | empirical FIM | Monte Carlo FIM | Hutchinson FIM |
|-----------------------------------------|-----------------------------------------------------|-----------------------------------------------------|---------------------------------------|
| $\mathcal{I}(z)$ / $\mathcal{F}(heta)$ | $\hat{\mathcal{I}}(z)$ / $\hat{\mathcal{F}}(heta)$ | $\hat{\mathcal{I}}(z)$ / $\hat{\mathcal{F}}(heta)$ | $\mathbb{I}(z)$ / $\mathbb{F}(heta)$ |

2 GEOMETRY OF LOW-DIMENSIONAL CORE SPACES

Consider a classifier network $p(y \mid x, \theta) \coloneqq p(y \mid z(x, \theta))$, where $z(x, \theta)$ is last layer's linear output. Due to the chain rule, we plug $\frac{\partial \ell_{xy}}{\partial \theta} = \left(\frac{\partial z}{\partial \theta}\right)^{\top} \frac{\partial \ell_{xy}}{\partial z}$ into Eq. (1). Then, we can easily arrive at

$$\mathcal{F}(\theta) = \sum_{x \in \mathcal{D}_x} \left(\frac{\partial z}{\partial \theta} \right)^{\top} \cdot \mathcal{I}(z(x, \theta)) \cdot \frac{\partial z}{\partial \theta}, \tag{2}$$

which is in the form of a Gauss-Newton matrix (Martens et al., 2010), or a pullback metric tensor (Sun, $2020)^2$ from a low dimensional statistical manifold with metric $\mathcal{I}(z)$, to the much higher dimensional neuromanifold with metric $\mathcal{F}(\theta)$. In this section, we rediscover the geometrical structure of the low dimensional statistical manifold, which we refer to as the *core space*, or simply the *core*.

In multi-class classification, y (given a feature vector x) follows a category distribution $p(y=i\,|\,x,\theta)=p_i(x,\theta),\,i=1,\cdots,C.$ All possible category distributions over $\{1,\cdots,C\}$ form a closed statistical simplex $\Delta^{C-1}\coloneqq \left\{(p_1,\cdots,p_C): \sum_{i=1}^C p_i=1; \, \forall i,p_i\geq 0\right\}$. The superscript C-1 denotes the dimensionality of Δ and can be omitted. If $p\in \operatorname{int}(\Delta^{C-1})$ (interior of Δ^{C-1}), we can reparameterize $p=\operatorname{SoftMax}(z)$, where $z\in\Re^C$ is the logits. The core Δ^{C-1} is a curved space,

²Strictly speaking, the pullback tensor requires the Jacobian of $\theta \to z$ have full column rank everywhere, which is not satisfied in typical settings of deep neural networks. This leads to singular metric tensors.

where p or z serves as a coordinate system in the sense that different choices of p or z yield different distributions. By Eq. (1), the FIM is:

$$\mathcal{I}^{\Delta}(z) = \mathbb{E}\left[(e_y - p)(e_y - p)^{\top} \right] = \operatorname{diag}(p) - pp^{\top}, \tag{3}$$

where $\operatorname{diag}(\cdot)$ means the diagonal matrix constructed with a given diagonal vector. In below, depending on context, $\operatorname{diag}(\cdot)$ also denotes a diagonal vector extracted from a square matrix. e (without subscripts) denotes a vector of all ones, e_y denotes the one-hot vector with only the y'th bit activated, and e_{ij} denotes the binary matrix with only the ij'th entry set to 1. Note z is a redundant coordinate system as $\dim(z) = C > C - 1$. If $z \in \operatorname{int}(\Delta^{C-1})$, $\mathcal{I}^{\Delta}(z)$ has a one-dimensional kernel: one can easily verify $\mathcal{I}^{\Delta}(z)(te) = 0$ for all $t \in \Re$.

By noting that $\mathcal{I}^{\Delta}(z)$ is a rank-1 perturbation of the diagonal matrix diag (p), we can apply Cauchy's interlacing theorem and study the spectral properties of $\mathcal{I}^{\Delta}(z)$.

Theorem 1 (Spectrum of Simplex FIM). Assume the spectral decomposition $\mathcal{I}^{\Delta}(z) = \sum_{i=1}^{C} \lambda_i v_i v_i^{\top}$, where $\lambda_1 \leq \cdots \leq \lambda_C$. Then $\lambda_1 = 0$; $v_1 = e/\|e\|$; $\sum_{i=1}^{C} \lambda_i = 1 - \|p\|^2$; and

$$\max \left\{ p_i(1-p_i) \right\} \cup \left\{ p_{(C-1)}, \frac{1-\|p\|^2}{C-1} \right\} \le \lambda_C \le \min \left\{ p_{(C)}, 2 \max_i (p_i(1-p_i)), 1-\|p\|^2 \right\},$$

where $p_{(C-1)}$ and $p_{(C)}$ denote the second-largest and the largest elements of p, respectively.

The largest eigenvalue of $\mathcal{I}^{\Delta}(z)$, denoted as λ_C , and its associated eigenvector correspond to the "most informative" direction at any $z \in \Delta^{C-1}$. By Theorem 1, λ_C can be bounded from above and below. The bound gap is at most $\min\{p_{(C)}-p_{(C-1)}, \max_i(p_i(1-p_i))\}$. We have found through numerical simulations that, in practice, the bounds in Theorem 1 are quite tight and can provide an estimate of λ_C within a narrow range. The lemma below gives lower and upper bounds of $\mathcal{I}^{\Delta}(z)$, both with a simpler structure than $\mathcal{I}^{\Delta}(z)$, in the space of psd matrices based on Löwner partial order.

Lemma 2. $\forall z \in \operatorname{int}(\Delta^{C-1})$, assume the spectral decomposition $\mathcal{I}^{\Delta}(z) = \sum_{i=1}^{C} \lambda_i v_i v_i^{\top}$, where $\lambda_1 \leq \cdots \leq \lambda_{C-1} < \lambda_C$. Then, $\lambda_C v_C v_C^{\top} \leq \mathcal{I}^{\Delta}(z) \leq \operatorname{diag}(p)$. Moreover, $\lambda_C v_C v_C^{\top}$ is the best rank-1 representation of $\mathcal{I}^{\Delta}(z)$ in the sense that no rank-1 matrix $B \neq \lambda_C v_C v_C^{\top}$ satisfies $\lambda_C v_C v_C^{\top} \leq B \leq \mathcal{I}^{\Delta}(z)$. Meanwhile, $\operatorname{diag}(p)$ is the best diagonal representation of $\mathcal{I}^{\Delta}(z)$ in the sense that no diagonal matrix $D \neq \operatorname{diag}(p)$ satisfies $\mathcal{I}^{\Delta}(z) \leq D \leq \operatorname{diag}(p)$.

The simplex FIM is upper-bounded by a diagonal matrix and lower bounded by a rank-1 matrix. By Lemma 2, $\lambda_C v_C v_C^{\mathsf{T}}$ is the *lower-envelope* (greatest lower bound) of $\mathcal{I}^{\Delta}(z)$ in rank-1 matrices, and diag (p) is the *upper-envelope* (least upper bound) of $\mathcal{I}^{\Delta}(z)$ in diagonal matrices. If the bounds in Lemma 2 are used as a deterministic estimate of $\mathcal{I}^{\Delta}(z)$, the error can be controlled, as shown below.

Lemma 3. We have
$$\forall z \in \Delta$$
, $\|\mathcal{I}^{\Delta}(z) - \operatorname{diag}(p)\| = \|p\|^2 \ge \frac{1}{C}$; meanwhile, $\|\mathcal{I}^{\Delta}(z) - \lambda_C v_C v_C^\top\| \le \min\left\{1 - \|p\| - p_{(C-1)}, \sqrt{\sum_{i=2}^{C-1} p_{(i)}^2}\right\}$, where $p_{(i)}$ denote the entries of p sorted in ascending order.

Note $\sqrt{\sum_{i=2}^{C-1} p_{(i)}^2}$ is the Euclidean norm of *trimmed p, i.e.* the vector obtained by removing p's smallest and largest elements. By Lemma 3, the upper bound diag (p) always incurs an error of at least 1/C. Depending on p, the lower bound $\lambda_C v_C v_C^{\top}$ can more accurately estimate $\mathcal{I}^{\Delta}(z)$ as the error can go to zero.

Alternatively, one can use random matrices to estimate $\mathcal{I}^{\Delta}(z)$. By Eq. (3), the rank-1 matrix $R(y) = (e_y - p)(e_y - p)^{\top}$ is an unbiased estimator of $\mathcal{I}^{\Delta}(z)$. The eFIM of Δ is given by $\hat{\mathcal{I}}^{\Delta}(z) = R(y)$, where y is a given empirical sample of the distribution specified by z. The lemma below shows the worst case error of using eFIM to estimate $\mathcal{I}^{\Delta}(z)$.

Lemma 4.
$$\forall z \in \Delta^{C-1}$$
, $\exists y \in \{1, \dots, C\}$, such that $||R(y) - \mathcal{I}^{\Delta}(z)|| \ge 1 + ||p||^2 - \lambda_C - 2p_{(1)} \ge 2||p||^2 - 2p_{(1)}$.

The first " \geq " is tighter but the second " \geq is easier to interpret. The term ||p|| can be as large as 1 (when p is close to one-hot). In such cases, using R(y) to estimate $\mathcal{I}^{\Delta}(z)$ may incur significant error if y is adversarially chosen.

In classification tasks with multiple binary labels, we assume $p(y_i = 1 \mid x) = p_i$ $(i = 1, \dots, C)$ and that all dimensions of y are conditional independent given x. All such distributions form a C-dimensional hypercube $\mathcal{C}^C(p) = \{(p_1, \dots, p_C) : \forall i, 0 \leq p_i \leq 1\}$, which is the product space of 1-dimensional simplices. Consider $p_i = \sigma(z_i) \coloneqq 1/(1 + \exp(-z_i))$ for $i = 1, \dots, C$. In this case, the FIM is a diagonal matrix, given by

$$\mathcal{I}^{\mathcal{C}}(z) = \operatorname{diag}((p_1(1-p_1), \cdots, p_C(1-p_C))) = \operatorname{diag}(\sigma'(z_i), \cdots, \sigma'(z_C)). \tag{4}$$

In what follows, unless stated otherwise, our results pertain to the core Δ as it is more commonly used and has a more complex FIM as compared to \mathcal{C} .

3 FIM FOR CLASSIFIER NETWORKS — DETERMINISTIC ANALYSIS

We give a lower and upper bound of $\mathcal{F}^{\Delta}(\theta)$ (Proposition 5) and analyze each bound gap (Propositions 7 and 8). Our bounds result from simple matrix analysis and are more operational than related theoretical bounds such as monotonicity of the FIM under marginalization or coarse-graining (Amari, 2016). Our bounds are novel as they are built on envelopes (tightest bound) in the core and that they depend on the order statistics of the output probability vector.

3.1 DETERMINISTIC LOWER AND UPPER BOUNDS

By Eq. (2), the neuromanifold FIM $\mathcal{F}(\theta)$ is determined by both the core space and the parameteroutput Jacobian $\frac{\partial z}{\partial \theta}$. Similar to Lemma 2, we can have lower and upper bounds of $\mathcal{F}^{\Delta}(\theta)$ in the space of psd matrices (although these bounds are not envelopes as in Lemma 2).

Proposition 5. If $p(y | x, \theta) \in \Delta^{C-1}$ is categorical, then $\forall \theta \in \Theta$, we have

$$\sum_{x \in \mathcal{D}_x} \lambda_C \left(\frac{\partial z}{\partial \theta} \right)^\top v_C v_C^\top \frac{\partial z}{\partial \theta} \preceq \mathcal{F}^\Delta(\theta) \preceq \sum_{x \in \mathcal{D}_x} \sum_{y=1}^C p(y \,|\, x, \theta) \frac{\partial z_y}{\partial \theta} \left(\frac{\partial z_y}{\partial \theta} \right)^\top,$$

where $\lambda_C := \lambda_C(x, \theta)$ and $v_C := v_C(x, \theta)$ denote the largest eigenvalue and its associated eigenvector of $\mathcal{I}(z(x, \theta))$.

Remark. The LHS is a sum of $|\mathcal{D}_x|$ (number of samples in \mathcal{D}_x) matrices of rank-1. Its rank is at most $|\mathcal{D}_x|$. The RHS is a sum of $C|\mathcal{D}_x|$ matrices of rank-1 and potentially has a larger rank.

If $p(y \mid x)$ is in \mathcal{C} , then $\mathcal{I}^{\mathcal{C}}(z(x,\theta))$ is diagonal as in Eq. (4). By Eq. (2), we have $\mathcal{F}^{\mathcal{C}}(\theta) = \sum_{x \in \mathcal{D}_x} \sum_{y=1}^C p_y (1-p_y) \frac{\partial z_y}{\partial \theta} \left(\frac{\partial z_y}{\partial \theta}\right)^{\top}$, which is similar to the upper bound in Proposition 5. In summary, $\mathcal{F}(\theta)$ can be bounded or computed using the Jacobian $\frac{\partial z}{\partial \theta}$ as well as the output probabilities $p(y \mid x, \theta)$. The following analysis depends on the spectral properties of $\frac{\partial z}{\partial \theta}$. Across our formal statements, we denote the singular values of $\frac{\partial z}{\partial \theta}$, sorted in ascending order, as $\sigma_1(x, \theta) \leq \cdots \leq \sigma_C(x, \theta)$. In Proposition 5, by taking the trace on all sides, the trace of the FIM can be bounded from above and below.

Corollary 6. If $p(y | x, \theta) \in \Delta^{C-1}$ is categorical, then it holds for all $\theta \in \Theta$ that

$$\sum_{x \in \mathcal{D}_x} \lambda_C(x, \theta) \sigma_1^2(x, \theta) \leq \sum_{x \in \mathcal{D}_x} \sum_{i=2}^C \lambda_i(x, \theta) \sigma_{C+1-i}^2(x, \theta) \leq \operatorname{tr}(\mathcal{F}^{\Delta}(\theta)) \leq \sum_{x \in \mathcal{D}_x} \sum_{y=1}^C p(y \mid x, \theta) \left\| \frac{\partial z_y}{\partial \theta} \right\|^2.$$

These bounds are useful to get the overall scale of $\mathcal{F}^{\Delta}(\theta)$ without computing its exact value. The proposition below gives the error of the upper bound in Proposition 5 in terms of Frobenius norm.

Proposition 7. We have $\forall \theta \in \Theta$ that

$$\sqrt{\sum_{x \in \mathcal{D}_x} \left\| \left(\frac{\partial z}{\partial \theta} \right)^\top p(x, \theta) \right\|^4} \le \left\| \sum_{x \in \mathcal{D}_x} \sum_{y=1}^C p(y \mid x, \theta) \left(\frac{\partial z_y}{\partial \theta} \right)^\top \frac{\partial z_y}{\partial \theta} - \mathcal{F}^{\Delta}(\theta) \right\| \le \sum_{x \in \mathcal{D}_x} \|p(x, \theta)\|^2 \sigma_C^2(x, \theta),$$

where $p(x, \theta) = \text{SoftMax}(z(x, \theta))$ denotes the output probability vector.

We use Frobenius norm for matrices but it is not difficult to bound the spectral norm using similar techniques. By Proposition 7, the error of the upper bound scales with the 2-norm (maximum singular value) of the parameter-output Jacobian $\frac{\partial z}{\partial \theta}$. Similar to what happens in the core space, using the upper bound of the FIM *always incurs an error*. For example, let p tend to be one-hot, the LHS in Proposition 7 does not vanish but scales with certain rows of $\frac{\partial z}{\partial \theta}$ corresponding to the predicted y. Naturally, we also want to examine the error of the lower bound in Proposition 5, as detailed below.

Proposition 8. We have $\forall \theta \in \Theta$ that

$$\left\| \sum_{x \in \mathcal{D}_x} \lambda_C \left(\frac{\partial z}{\partial \theta} \right)^\top v_C v_C^\top \frac{\partial z}{\partial \theta} - \mathcal{F}^{\Delta}(\theta) \right\| \leq \sum_{x \in \mathcal{D}_x} \sqrt{\sum_{i=2}^{C-1} \sigma_{i+1}^4(x,\theta) p_{(i)}^2(x,\theta)}.$$

Clearly, as p approaches a one-hot vector, all elements in the trimmed vector $p_{(i)}$, for $i=2,\cdots,C-1$, tend to zero, and the error approaches zero since its upper bound on the RHS goes to zero. From this view, the lower bound in Proposition 5 is a better estimate as compared to the upper bound.

Remark. By noting that $0 \le \sigma_i(x,\theta) \le \sigma_C(x,\theta)$, we can relax the bound in Proposition 8 to be comparable to Proposition 7: $\left\|\sum_{x \in \mathcal{D}_x} \lambda_C \left(\frac{\partial z}{\partial \theta}\right)^\top v_C v_C^\top \frac{\partial z}{\partial \theta} - \mathcal{F}^\Delta(\theta)\right\| \le \sum_{x \in \mathcal{D}_x} \sqrt{\sum_{i=2}^{C-1} p_{(i)}^2(x,\theta)} \cdot \sigma_C^2(x,\theta)$. The estimation error of $\sum_{x \in \mathcal{D}_x} \lambda_C \left(\frac{\partial z}{\partial \theta}\right)^\top v_C v_C^\top \frac{\partial z}{\partial \theta}$ is controlled by the norms of the Jacobian and the trimmed probabilities $(p_{(2)}, \cdots, p_{(C-1)})$. The latter is upper bounded by $p_{(C-1)}(x,\theta)$, the second largest probability of each sample x. By comparing with Proposition 7, one can easily observe that Proposition 8 is tighter in general.

3.2 EMPIRICAL FIM (EFIM)

Recall from the introduction, the eFIM $\hat{\mathcal{F}}(\theta)$ gives a biased, deterministic estimate of $\mathcal{F}(\theta)$. Intuitively, when the network is trained, computations based on the given labels are close to the expectation w.r.t. $p(y \mid x)$, and the eFIM is expected to approximate $\mathcal{F}(\theta)$ well. However, the bias of $\hat{\mathcal{F}}(\theta)$ can be enlarged if y is set adversarially. By simple derivations, $\hat{\mathcal{F}}(\theta) = \sum_{x \in \mathcal{D}_x} \left(\frac{\partial z}{\partial \theta}\right)^\top \cdot R(y) \cdot \frac{\partial z}{\partial \theta}$. Observe that it is similar to Eq. (2), except $\mathcal{I}(z(x,\theta))$ is replaced by its empirical counterpart R(y). If the neural network output is in the simplex, the error of eFIM can be bounded, as stated below.

Proposition 9.
$$\forall \theta \in \Theta$$
, $\forall y$, we have $\|\mathcal{F}^{\Delta}(\theta) - \hat{\mathcal{F}}^{\Delta}(\theta)\|_{\sigma} \leq \sum_{x \in \mathcal{D}_x} (1 + \|p(x, \theta)\|^2) \sigma_C^2(x, \theta)$.

Here we need to switch to the spectral norm $\|\cdot\|_{\sigma}$ to get a simple expression of the upper bound. The approximation error in terms of the spectral norm is controlled by the spectral norm of the parameter-output Jacobian. The error by Frobenius norm is even larger. The bound is loose as compared to Propositions 7 and 8.

We have found in Lemma 4 that using R(y) to approximate $\mathcal{I}^{\Delta}(z)$ suffers from a large error if y is chosen in a tricky way. The same principle applies to using $\hat{\mathcal{F}}(\theta)$ to approximate $\mathcal{F}(\theta)$.

Proposition 10. $\forall \theta \in \Theta, \forall x, \exists y, such that$

$$\left\| \left(\frac{\partial z}{\partial \theta} \right)^{\top} \mathcal{I}^{\Delta}(z(x,\theta)) \frac{\partial z}{\partial \theta} - \left(\frac{\partial z}{\partial \theta} \right)^{\top} R(y) \frac{\partial z}{\partial \theta} \right\|_{\sigma} \ge \sigma_{1}^{2}(x,\theta) \left| 1 + \| p(x,\theta) \|^{2} - \lambda_{C}(x,\theta) - 2p_{(1)}(x,\theta) \right|.$$

In the above inequality, the LHS is the error of $\hat{\mathcal{F}}(\theta)$ for one single $x \in \mathcal{D}_x$. Therefore, when y is set unfavorably, the eFIM suffers from an approximation error that scales with the smallest singular value of $\frac{\partial z}{\partial \theta}$. Among all the investigated deterministic approximations, the lower bound in Proposition 5 provides the smallest guaranteed error but is relatively expensive to compute. We solve the computational issues in the next section.

4 HUTCHINSON'S ESTIMATE OF THE FIM

4.1 LIMITATIONS OF MONTE CARLO ESTIMATES

We show that the quality of the MC estimate $\mathcal{F}(\theta)$ can be arbitrarily bad. Consider the single neuron model $z = \theta x$ for binary classification, where z, θ, x are all scalars, and θ is close to zero.

Then $p \approx \frac{1}{2}$ is a fair Bernoulli distribution. $\mathcal{I}(z) = p(1-p) \approx \frac{1}{4}$. The Jacobian is simply $\frac{\partial z}{\partial \theta} = x$. and $\mathcal{F}(\theta) = \mathbb{E}_{p(x)} \left[\frac{\partial z}{\partial \theta} \mathcal{I}(z) \frac{\partial z}{\partial \theta} \right] \approx \frac{1}{4} \mathbb{E}_{p(x)}[x^2]$. A basic MC estimator takes the form $\hat{\mathcal{F}}(\theta) = \frac{1}{4m} \sum_{i=1}^m x_i^2$, where x_i 's are independently and identically distributed according to p(x). Its variance is $\mathrm{Var}(\hat{\mathcal{F}}) = \frac{1}{4m} \left[\mathbb{E}_{p(x)}(x^4) - \mathbb{E}_{p(x)}^2(x^2) \right]$. We let p(x) be a heavy tailed distribution, e.g. Student's t-distribution with $\nu > 4$ degrees of freedom, so that $\mathrm{Var}(\hat{\mathcal{F}})$ is large while $\mathcal{F}(\theta)$ is small. Then $\mathbb{E}_{p(x)}(x^2) = \frac{\nu}{\nu-2}$ and $\mathbb{E}_{p(x)}(x^4) = \frac{3\nu^2}{(\nu-2)(\nu-4)}$. The ratio $\frac{\mathbb{E}_{p(x)}(x^4)}{(\mathbb{E}_{p(x)}x^2)^2} = \frac{3(\nu-2)}{\nu-4}$ can be arbitrarily large when $\nu \to 4^+$. Therefore the coefficient of variation (CV) $\mathrm{Std}(\hat{\mathcal{F}})/F(\theta)$ is unbounded. Throughout our analysis, the CV is a key indicator of the quality of a FIM estimator, as a bounded CV for a random variable X ensures the random estimator's probability mass within $[0,\alpha\mu]$, where $\alpha > 1$ and $\mu \geq 0$ is the mean of X. If $\mathrm{CV} = \frac{\mathrm{Std}X}{\mu} \leq K$, then by Cantelli inequality, we have $\mathbb{P}(X \geq \alpha\mu) = \mathbb{P}(X \geq \mu + (\alpha-1)\mu) \leq \mathbb{P}(X \geq \mu + \frac{\alpha-1}{K}\mathrm{Std}X) \leq \left(1 + \left(\frac{\alpha-1}{K}\right)^2\right)^{-1}$. The general case is more complicated, but follows a similar idea. The variance of MC estimators depends on the 4th moment of the Jacobian $\frac{\partial z}{\partial \theta}$ w.r.t. p(x) while the mean value $\mathcal{F}(\theta)$ only depends on the 2nd moment of $\frac{\partial z}{\partial \theta}$. The ratio of the variance and $\mathcal{F}^2(\theta)$, or the CV $\mathrm{Std}(\hat{\mathcal{F}})/\mathcal{F}(\theta)$, is unbounded without further assumption on p(x). One can increase the number of samples m to reduce variance. However, this is computationally expensive especially in online settings.

4.2 HUTCHINSON'S ESTIMATE

In light of the challenges of MC estimates, we introduce a new way to get an unbiased estimate of the FIM. First, compute the scalar-valued function

$$\mathfrak{h}(\mathcal{D}_x, \theta) := \sum_{x \in \mathcal{D}_x} \sum_{y=1}^C \sqrt{\tilde{p}(y \mid x, \theta)} \ell_{xy}(\theta) \xi_{xy}, \tag{5}$$

where ξ_{xy} is a standard multivariate Gaussian vector of size $C|\mathcal{D}|$ or a Rademacher vector, and $\tilde{p}(y|x,\theta)$ has the same value as $p(y|x,\theta)$ but is non-differentiable, meaning its gradient is always zero, preventing error from back-propagating through $\tilde{p}(y|x,\theta)$. This \tilde{p} can be implemented by Tensor.detach() in PyTorch (Paszke et al., 2019) or similar functions in other auto-differentiation (AD) frameworks. Second, the gradient vector $\frac{\partial \mathfrak{h}}{\partial \theta} = \sum_{x \in \mathcal{D}_x} \sum_{y=1}^C \sqrt{p(y|x,\theta)} \frac{\partial \ell_{xy}}{\partial \theta} \xi_{xy}$ can be evaluated via AD, e.g. by \mathfrak{h} .backward() in Pytorch. Third, the random psd matrix $\mathbb{F}(\theta) := \frac{\partial \mathfrak{h}}{\partial \theta} \frac{\partial \mathfrak{h}}{\partial \theta^{-1}}$, which we refer to as the "Hutchinson's estimate" (of the FIM), can be used to estimate $\mathcal{F}(\theta)$. By straightforward derivations,

$$\mathbb{E}_{p(\xi)}\left(\mathbb{F}(\theta)\right) = \sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} \sum_{x' \in \mathcal{D}_{x}} \sum_{y'=1}^{C} \sqrt{p(y \mid x, \theta)} \sqrt{p(y' \mid x', \theta)} \frac{\partial \ell_{xy}}{\partial \theta} \frac{\partial \ell_{x'y'}}{\partial \theta} \mathbb{E}_{p(\xi)}\left[\xi_{xy} \xi_{x'y'}\right] = \mathcal{F}(\theta). \quad (6)$$

The last "=" is because $\mathbb{E}_{p(\xi)}(\xi_{xy}\xi_{x'y'})=1$ if x=x' and y=y', and $\mathbb{E}_{p(\xi)}(\xi_{xy}\xi_{x'y'})=0$ otherwise. Considering $\frac{\partial \mathfrak{h}}{\partial \theta}$ as an implicit representation of the FIM, its **computational cost** is \oplus evaluating the \mathfrak{h} function, \oplus the backward pass to compute the gradient of \mathfrak{h} . The cost is same as evaluating the gradient of the loss $-\sum_{x\in\mathcal{D}_x}\sum_{y=1}^C\ell_{xy}(\theta)$, noting that \mathfrak{h} is the log-likelihood randomly flipped by a Gaussian/Rademacher vector. Moreover, \mathfrak{h} can reuse the logits already computed during the forward pass. Therefore $\frac{\partial \mathfrak{h}}{\partial \theta}$ requires merely one additional backward pass, making it practical for large scale networks. In summary, $\mathbb{F}(\theta)$ is a *universal estimator* of $\mathcal{F}(\theta)$ for general statistical model, which is independent of neural network architectures and applicable to non-neural network models. Hutchinson's estimate has guaranteed quality, as formally established below.

Proposition 11. $\mathbb{E}_{p(\xi)}(\mathbb{F}(\theta)) = \mathcal{F}(\theta)$. If $p(\xi)$ is standard multivariate Gaussian, then $\operatorname{Var}(\mathbb{F}_{ii}(\theta)) = 2\mathcal{F}_{ii}(\theta)^2$; if $p(\xi)$ is standard multivariate Rademacher, $\operatorname{Var}(\mathbb{F}_{ii}(\theta)) = 2\mathcal{F}_{ii}(\theta)^2 - 2\sum_{x \in \mathcal{D}_x} \sum_{y=1}^C p^2(y \mid x) (\frac{\partial \ell_{xy}}{\partial \theta_i})^4$.

It is known that Rademacher distribution yields smaller variance for Hutchinson's estimator compared to the Gaussian distribution. In what follows, $p(\xi)$ is Rademacher by default. By Proposition 11, $\operatorname{Std}(\mathbb{F}_{ii}(\theta)) \leq \sqrt{2}\mathcal{F}_{ii}(\theta)$. Thus the CV $\operatorname{Std}(\mathbb{F}_{ii}(\theta))/\mathcal{F}_{ii}(\theta)$ is bounded by $\sqrt{2}$. We only investigate

the diagonal of Hutchinson's estimate because the diagonal FIM is widely used, but our results can be readily extended to off-diagonal entries.

Remark. Taking trace on both sides of $\mathbb{E}_{p(\xi)}(\mathbb{F}(\theta)) = \mathcal{F}(\theta)$, we get $\mathbb{E}_{p(\xi)}(\|\frac{\partial \mathfrak{h}}{\partial \theta}\|^2) = \operatorname{tr}(\mathcal{F}(\theta))$. The squared Euclidean-norm of $\frac{\partial \mathfrak{h}}{\partial \theta}$ is an unbiased estimate of the trace of the FIM. This is useful for computing related regularizers (Peebles et al., 2020).

Note that a sample of the random matrix $\mathbb{F}(\theta)$ is always rank-1: $\mathrm{rank}\,\mathbb{F}(\theta)=1\leq \mathrm{rank}\,\mathcal{F}(\theta)$, but the expectation of $\mathbb{F}(\theta)$ has the same rank as $\mathcal{F}(\theta)$. Ideally, one can compute the numerical average of more than one $\mathbb{F}(\theta)$ samples to reduce variance and recover the rank, each requiring a separate backward pass. Due to computational constraints in deep learning practice, much fewer (e.g.,1) samples are used. Instead, accumulated statistics along the learning path $\theta_1\to\theta_2\to\cdots$ can be used to maintain a (exponential) moving average of $\mathbb{F}(\theta_i)$. The underlying assumption is that θ_1,θ_2,\cdots connected by small learning steps lie close to one another in the parameter space. Therefore, averaging $\mathbb{F}(\theta_i)$ provides a reasonable approximation of the local FIM with sufficient rank.

4.3 DIAGONAL CORE

For multi-label classification, and for computing the upper bound in Proposition 5, the core matrix is diagonal, in the form $\mathcal{I}^{\mathrm{DG}}(z(x,\theta)) = \mathrm{diag}\,(\zeta_1(x,\theta),\cdots,\zeta_C(x,\theta))$, and the associated FIM is $\mathcal{F}^{\mathrm{DG}}(\theta) = \sum_{x \in \mathcal{D}_x} \left(\frac{\partial z}{\partial \theta}\right)^{\top} \cdot \mathcal{I}^{\mathrm{DG}}(z(x,\theta)) \cdot \frac{\partial z}{\partial \theta}$. In the former case, $\zeta_y(x,\theta) = p(y\,|\,x,\theta)(1-p(y\,|\,x,\theta))$; in the latter case, $\zeta_y(x,\theta) = p(y\,|\,x,\theta)$. Here, the tensor superscript — e.g., "DG" for diagonal or "LR" for low-rank — indicates the parametric form of the core FIM, in contrast to denoting the core space as in \mathcal{I}^{Δ} . We define the scalar valued function

$$\mathfrak{h}^{\mathrm{DG}}(\theta) \coloneqq \sum_{x \in \mathcal{D}_x} \sum_{y=1}^C \sqrt{\tilde{\zeta}_y(x,\theta)} z_y(x,\theta) \xi_{xy},\tag{7}$$

where ξ_{xy} are standard Rademacher samples that are independent across all x and y. Similar to the derivation steps in section 1, we first compute the random vector $\frac{\partial \mathfrak{h}^{\mathrm{DG}}}{\partial \theta}$ through AD, and then compute $\mathbb{F}^{\mathrm{DG}}(\theta) \coloneqq \frac{\partial \mathfrak{h}^{\mathrm{DG}}}{\partial \theta} \frac{\partial \mathfrak{h}^{\mathrm{DG}}}{\partial \theta^{\mathrm{T}}}$ (or its diagonal blocks) to estimate $\mathcal{F}^{\mathrm{DG}}(\theta)$.

For computing the upper bound in Proposition 5, $\tilde{\zeta}_y(x,\theta) = \tilde{p}_y(x,\theta)$, then we find that Eq. (5) and Eq. (7) are similar. The only difference is that, the "raw" logits z_y in Eq. (7) is replaced by $\ell_{xy}(\theta) = z_y - \log \sum_y \exp(z_y)$ in Eq. (5). Compared to $\frac{\partial z}{\partial \theta}$, the gradient $\frac{\partial \ell_{xy}}{\partial \theta} = \frac{\partial z_y}{\partial \theta} - \sum_y p(y \mid x,\theta) \frac{\partial z_y}{\partial \theta}$ is centered. Due to their computational similarity, in practice, one should use Eq. (5) instead of Eq. (7) and get an unbiased estimate of $\mathcal{F}^{\Delta}(\theta)$. Eq. (7) is useful when the dimensions of y are conditional independent given x, e.g. for computing $\mathcal{F}^{\mathcal{C}}(\theta)$.

4.4 LOW RANK CORE

By Lemma 2, the FIM of the core space Δ has a rank-1 lower bound $\mathcal{I}^{\Delta}(z) \succeq \mathcal{I}^{\mathrm{LR}}(z) \coloneqq \lambda_C v_C v_C^{\top}$. By Proposition 5, $\mathcal{F}^{\Delta}(\theta) \succeq \mathcal{F}^{\mathrm{LR}}(\theta) \coloneqq \sum_{x \in \mathcal{D}_x} \lambda_C(x,\theta) \left(\frac{\partial z}{\partial \theta}\right)^{\top} v_C(x,\theta) v_C^{\top}(x,\theta) \frac{\partial z}{\partial \theta}$. We define

$$\mathfrak{h}^{LR}(\theta) = \sum_{x \in \mathcal{D}_x} \sqrt{\tilde{\lambda}_C(x,\theta)} \tilde{v}_C^{\top}(x,\theta) z(x,\theta) \xi_x, \tag{8}$$

where ξ_x are independent standard Rademacher samples. For computing $\mathfrak{h}^{LR}(\theta)$, we only need $|\mathcal{D}_x|$ Rademacher samples, as compared to $C|\mathcal{D}_x|$ samples for computing $\mathfrak{h}(\theta)$ and $\mathfrak{h}^{DG}(\theta)$. Correspondingly, $\mathbb{F}^{LR}(\theta) := \frac{\partial \mathfrak{h}^{LR}}{\partial \theta} \frac{\partial \mathfrak{h}^{LR}}{\partial \theta^{\top}}$ is used to estimate $\mathcal{F}^{LR}(\theta)$. Note that $\mathfrak{h}, \mathfrak{h}^{DG}$ and \mathfrak{h}^{LR} can be computed solely based on the neural network output logits $z(x,\theta)$ for each $x \in \mathcal{D}_x$.

We remain to solve $\lambda_C(x,\theta)$ and $v_C(x,\theta)$ for each $x \in \mathcal{D}_x$. They can be conveniently computed based on the power iteration. By Eq. (3), starting from a random unit vector v_C^0 , we compute

$$v_C^{t+1} = \frac{\mathcal{I}^{\Delta}(z)v_C^t}{\|\mathcal{I}^{\Delta}(z)v_C^t\|} = \frac{p \circ v_C^t - p^\top v_C^t p}{\|p \circ v_C^t - p^\top v_C^t p\|},$$

4.5 NUMERICAL SIMULATIONS

 To provide intuition, we compute the diagonal FIM of DistilBERT (Sanh et al., 2019), pretrained by Hugging Face (Wolf et al., 2020) ³ combined with a randomly initialized classification head (two dense layers) for AG News (Zhang et al., 2015) topic classification (C=4 classes). More detailed quantitative results and another representative case is provided in section B, where DistilBERT is finetuned on the Stanford Sentiment Treebank v2 (SST-2) (Socher et al., 2013), and the FIM is computed in regions of Θ corresponding to a more confident model. Figure 1 shows the normalized density plots of $\mathbb{F}_{ii}^{\overline{\mathrm{DG}}}(\theta)$ (Hutchinson's estimate of the upper bound in Proposition 5), $\mathbb{F}_{ii}(\theta)$ (Hutchinson's unbiased estimate), $\mathbb{F}_{ii}^{LR}(\theta)$ (Hutchinson's estiamte of the lower bound in Proposition 5), and the empirical FIM $\hat{\mathcal{F}}_{ii}(\theta)$. All estimators use the first 128 data samples to compute the FIM. All Hutchinson estimators use 10 samples for variance reduction. Due to the pathological structure (Karakida et al., 2021) of the FIM, all densities exhibit a spike near zero and become sparse on large Fisher information values. For example, all layers have more than 20% of their parameters with $\mathbb{F}_{ii} < 10^{-5}$. The visualization is smoothed out on a logarithmic y-axis. The mean values of these densities are reflected on the low-right corner of the subplots (up to a scaling factor). Across the layers, the classification head has the largest scale of Fisher information and the embedding layer has the lowest scale. In general, the deeper layers (close to the input) have smaller values of \mathbb{F}_{ii} . The scale of \mathbb{F}_{ii}^{DG} appears larger than \mathbb{F}_{ii} , which in turn is larger than \mathbb{F}_{ii}^{LR} . This makes sense as the expected values of \mathbb{F}_{ii}^{DG} and \mathbb{F}_{ii}^{LR} are upper and lower bounds of the expected values of \mathbb{F}_{ii} , respectively. The scale of $\hat{\mathcal{F}}_{ii}$ is not informative as the others regarding \mathcal{F}_{ii} because it is biased. The classification head is not trained and hence has large gradient values, leading to large values of $\hat{\mathcal{F}}_{ii}$.

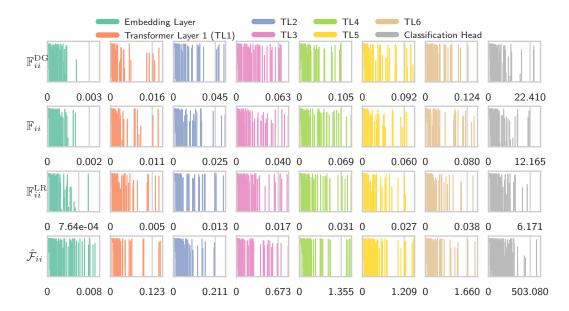


Figure 1: Density plots (based on kernel density estimation with a small bandwidth) of diagonal FIM elements based on different approximations (rows) across different layers (columns) of DistilBERT on the AG News dataset. The four rows, from top to bottom, represent Hutchinson's estimates $\mathbb{F}^{DG}(\theta)$, $\mathbb{F}(\theta)$, $\mathbb{F}^{LR}(\theta)$, and the eFIM $\hat{\mathcal{F}}(\theta)$. The columns are arranged from layers close to the input (left) to those near the output (right). In each subplot, the maximum value of the x-axis (number on the bottom right corner) shows the mean value of the FIM multiplied by 2,000. The y-axis means probability density in log-scale.

³Available as distilbert-base-uncased in the Hugging Face library.

In the experiments, we find that the computational speeds of \mathbb{F}^{DG} and \mathbb{F} are very similar. Computing \mathbb{F}^{LR} is slightly more expensive, with a computational overhead scaling with the number of classes C, because \mathbb{F}^{LR} requires the power iteration to compute λ_C and v_C . \mathbb{F} is unbiased, whereas both \mathbb{F}^{DG} and \mathbb{F}^{LR} are biased. For \mathbb{F}^{LR} , the Hutchinson probe has a lower dimensionality, leading to a better bias-variance trade-off than \mathbb{F} and \mathbb{F}^{DG} in our experiments.

5 RELATED WORK

A prominent application of Fisher information in deep learning is the natural gradient (Amari, 1998) and its variants. The Adam optimizer (Kingma & Ba, 2015) uses the empirical diagonal FIM. Efforts have been made to obtain more accurate approximations of $\mathcal{F}(\theta)$ at the expense of higher computational cost, such as modeling the diagonal blocks of $\mathcal{F}(\theta)$ with Kronecker product (Martens, 2020) of component-wise FIM (Ollivier, 2015; Sun & Nielsen, 2017), or computing $\mathcal{F}(\theta)$ through low rank approximations (Le Roux et al., 2007; Botev et al., 2017). The FIM can be alternatively defined on a sub-model (Sun & Nielsen, 2017) instead of the global mapping $x \to y$ or based on α -embeddings of a parametric family (Nielsen, 2017). AdaHessian (Yao et al., 2021) uses Hutchinson probes to approximate the diagonal Hessian.

From theoretical perspectives, the quality of Kronecker approximation is discussed (Martens & Grosse, 2015) with its error bounded. It is well known that the eFIM differs from $\mathcal{F}(\theta)$ (Pascanu & Bengio, 2014; Martens, 2020; Kunstner et al., 2020) and leads to distinct optimization paths. The accuracy of two different MC approximations of $\mathcal{F}(\theta)$ is analyzed (Guo & Spall, 2019; Soen & Sun, 2021; 2024; Sun & Spall, 2021), which lie in the framework of MC information geometry (Nielsen & Hadjeres, 2019). By our analysis, the Hutchinson's estimate $\mathbb{F}(\theta)$ has unique advantages over both MC and the eFIM. Notably, the MC estimate in section 4.1 needs to compute $\frac{\partial \ell_{\hat{x}\hat{y}}}{\partial \theta}$ for each $x \in \mathcal{D}_x$, while $\mathbb{F}(\theta)$ only needs to evaluate one gradient vector $\frac{\partial h}{\partial \theta}$. Our bounds improves over existing bounds, e.g. those of $\mathcal{F}(\theta)$ (Soen & Sun, 2024), through carefully analyzing the core space.

The Hutchinson's stochastic trace estimator is used to estimate the trace of the FIM (Jastrzebski et al., 2021), or the FIM for Gaussian processes (Stein et al., 2013; Geoga et al., 2020) where the FIM entries are in the form of a trace. Closely related to this is computations around the Hessian, where Hutchinson's trick is applied to compute the Hessian trace (Hu et al., 2024), or the principal curvature (Böttcher & Wheeler, 2024), or related regularizers (Peebles et al., 2020). The Hessian trace estimator is implemented in deep learning libraries (Dangel et al., 2020; Yao et al., 2020) and usually relies on the Hessian-vector product. As a natural yet important next step, our estimators leverage both Hutchinson's trick and AD's interfaces, avoid the need for expensive Hessian computations/approximations, and are well-suited in scalable settings. In Eq. (6), we perform a double contraction of a high dimensional tensor indexed by x, y, x', y', i and j (i and j are indices of the FIM) and thereby obtain an unbiased estimator of the full metric tensor $\mathcal{F}(\theta)$ including its substructures and trace. Our estimator can be applied to different classification networks regardless of the network architecture.

6 Conclusion

We explore the FIM $\mathcal F$ of classifier networks, focusing on the case of multi-class classification. We provide deterministic lower and upper bounds of the FIM based on related bounds in the low dimensional core space. We discover a new family of random estimators $\mathbb F$ based on Hutchinson's trace estimator. Their estimate has guaranteed quality with bounded variance and can be computed efficiently through auto-differentiation. The proposed $\mathbb F$ is readily integrated into deep learning libraries (Dangel et al., 2020; Yao et al., 2020) for efficiently evaluating the FIM or the Hessian. Our analysis in the core space gives insights and useful tools for information geometry where the simplex is widely used. As a limitation, the results here address novel computation of $\mathcal F$ but are not directly piped into a downstream application that uses the proposed $\mathbb F$. For example, new deep learning optimizers based on the proposed $\mathbb F$, are not developed here and left as future work. Advanced variance reduction techniques (Meyer et al., 2021) that could improve our proposed random estimator $\mathbb F(\theta)$ remain to be investigated.

486 ETHICS STATEMENT 487 488 The authors have read the ICLR Code of Ethics the confirm that this research fully complies with the 489 Code of Ethics. 490 491 REPRODUCIBILITY STATEMENT 492 493 The authors confirm that all assumptions and proofs of the theoretical developments are provided 494 in the main text and the appendix. The code to compute the proposed Hutchinson's estimate of the 495 Fisher information matrix will be released upon acceptance. 496 497 498 THE USE OF LARGE LANGUAGE MODELS (LLMS) 499 500 The authors acknowledge that LLMs are used for editing purpose (grammar, wording, and translation). 501 LLMs are not used to develop the core results. 502 REFERENCES 504 505 Shun-ichi Amari. Natural gradient works efficiently in learning. Neural Comput., 10(2):251–276, 506 1998. 507 Shun-ichi Amari. Information Geometry and Its Applications, volume 194 of Applied Mathematical 508 Sciences. Springer-Verlag, Berlin, 2016. 509 510 Stephen Blyth. Local divergence and association. Biometrika, 81(3):579–584, 1994. 511 512 Aleksandar Botev, Hippolyt Ritter, and David Barber. Practical Gauss-Newton optimisation for 513 deep learning. In Doina Precup and Yee Whye Teh (eds.), Proceedings of the 34th International Conference on Machine Learning, volume 70 of Proceedings of Machine Learning Research, pp. 514 557-565. PMLR, 2017. 515 516 Lucas Böttcher and Gregory Wheeler. Visualizing high-dimensional loss landscapes with Hessian 517 directions. Journal of Statistical Mechanics: Theory and Experiment, 2024(2):023401, 2024. 518 519 N. N. Čencov. Statistical Decision Rules and Optimal Inference. Translations of mathematical 520 monographs. American Mathematical Society, 1982. 521 Shixing Chen, Caojin Zhang, and Ming Dong. Coupled end-to-end transfer learning with generalized 522 Fisher information. In 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 523 pp. 4329-4338, 2018. 524 Felix Dangel, Frederik Kunstner, and Philipp Hennig. BackPACK: Packing more into backprop. In 526 International Conference on Learning Representations (ICLR), 2020. 527 528 Christopher J. Geoga, Mihai Anitescu, and Michael L. Stein. Scalable Gaussian process computations 529 using hierarchical matrices. Journal of Computational and Graphical Statistics, 29(2):227-237, 2020. 530 531 Shenghan Guo and James C. Spall. Relative accuracy of two methods for approximating observed 532 Fisher information. In Data-Driven Modeling, Filtering and Control: Methods and applications, 533 pp. 189-211. IET Press, London, 2019. 534 Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. Neural Comput., 9(1):1–42, January 1997. 536 ISSN 0899-7667. doi: 10.1162/neco.1997.9.1.1. URL https://doi.org/10.1162/neco. 537 1997.9.1.1. 538

Harold Hotelling. Spaces of statistical parameters. American Mathematical Society Meeting, 1929.

(unpublished. Presented orally by O. Ore during the meeting).

Zheyuan Hu, Zekun Shi, George Em Karniadakis, and Kenji Kawaguchi. Hutchinson trace estimation for high-dimensional and high-order physics-informed neural networks. *Computer Methods in Applied Mechanics and Engineering*, 424:116883, 2024. ISSN 0045-7825. URL https://www.sciencedirect.com/science/article/pii/S0045782524001397.

- M.F. Hutchinson. A stochastic estimator of the trace of the influence matrix for Laplacian smoothing splines. *Communications in Statistics Simulation and Computation*, 19(2):433–450, 1990.
- Stanislaw Jastrzebski, Devansh Arpit, Oliver Astrand, Giancarlo B Kerg, Huan Wang, Caiming Xiong, Richard Socher, Kyunghyun Cho, and Krzysztof J Geras. Catastrophic Fisher explosion: Early phase Fisher matrix impacts generalization. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 4772–4784. PMLR, 18–24 Jul 2021.
- Ryo Karakida, Shotaro Akaho, and Shun-ichi Amari. Pathological spectra of the Fisher information metric and its variants in deep neural networks. *Neural Computation*, 33(8):2274–2307, 2021.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017. doi: 10.1073/pnas.1611835114.
- Frederik Kunstner, Lukas Balles, and Philipp Hennig. Limitations of the empirical Fisher approximation for natural gradient descent. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 4133–4144. Curran Associates, Inc., 2020.
- Nicolas Le Roux, Pierre-Antoine Manzagol, and Yoshua Bengio. Topmoumoute online natural gradient algorithm. In *Advances in neural information processing systems*, volume 20, pp. 849–856. Curran Associates, Inc., 2007.
- Wu Lin, Frank Nielsen, Khan Mohammad Emtiyaz, and Mark Schmidt. Tractable structured natural-gradient descent using local parameterizations. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 6680–6691. PMLR, 18–24 Jul 2021.
- Abhilasha Lodha, Gayatri Belapurkar, Saloni Chalkapurkar, Yuanming Tao, Reshmi Ghosh, Samyadeep Basu, Dmitry Petrov, and Soundararajan Srinivasan. On surgical fine-tuning for language encoder. In *EMNLP* 2023, pp. 1–9. EMNLP, 2023.
- James Martens. New insights and perspectives on the natural gradient method. *Journal of Machine Learning Research*, 21(146):1–76, 2020.
- James Martens and Roger Grosse. Optimizing neural networks with Kronecker-factored approximate curvature. In *International Conference on Machine Learning (ICML)*, pp. 2408–2417. PMLR, 2015.
- James Martens et al. Deep learning via Hessian-free optimization. In *International Conference on Machine Learning (ICML)*, volume 27, pp. 735–742, 2010.
- Raphael A. Meyer, Cameron Musco, Christopher Musco, and David P. Woodruff. Hutch++: Optimal stochastic trace estimation. In *Proceedings of the 4th Symposium on Simplicity in Algorithms (SOSA)*, pp. 142–155, 2021.
- Frank Nielsen. The α -representations of the Fisher information matrix, 2017. https://franknielsen.github.io/blog/alpha-FIM/index.html.
- Frank Nielsen and Gaëtan Hadjeres. Monte Carlo information-geometric structures. In Frank Nielsen (ed.), *Geometric Structures of Information*, pp. 69–103. Springer International Publishing, Cham, 2019.

Yann Ollivier. Riemannian metrics for neural networks I: feedforward networks. *Information and Inference: A Journal of the IMA*, 4(2):108–153, 2015.

- Razvan Pascanu and Yoshua Bengio. Revisiting natural gradient for deep networks. In *International Conference on Learning Representations*, 2014.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pp. 8024–8035. Curran Associates, Inc., 2019. https://pytorch.org/.
- William Peebles, John Peebles, Jun-Yan Zhu, Alexei A. Efros, and Antonio Torralba. The Hessian penalty: A weak prior for unsupervised disentanglement. In *Proceedings of (ECCV) European Conference on Computer Vision*, pp. 581 597, August 2020.
- C. R. Rao. Information and the accuracy attainable in the estimation of statistical parameters. *Bulletin of the Calcutta Mathematical Society*, 37(3):81–91, 1945.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *CoRR*, abs/1910.01108, 2019. URL http://arxiv.org/abs/1910.01108. http://arxiv.org/abs/1910.01108.
- Maciej Skorski. Modern analysis of Hutchinson's trace estimator. In 2021 55th Annual Conference on Information Sciences and Systems (CISS), pp. 1–5, 2021.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Y. Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1631–1642, 2013.
- Alexander Soen and Ke Sun. On the variance of the Fisher information for deep learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, pp. 5708–5719. Curran Associates, Inc., 2021.
- Alexander Soen and Ke Sun. Trade-Offs of diagonal Fisher information matrix estimators. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 37, pp. 5870–5912. Curran Associates, Inc., 2024.
- Michael L. Stein, Jie Chen, and Mihai Anitescu. Stochastic approximation of score functions for Gaussian processes. *The Annals of Applied Statistics*, 7(2):1162 1191, 2013.
- Ke Sun. Information geometry for data geometry through pullbacks. In *Deep Learning through Information Geometry (Workshop at NeurIPS 2020)*, 2020.
- Ke Sun and Frank Nielsen. Relative Fisher information and natural gradient for learning large modular models. In *International Conference on Machine Learning (ICML)*, volume 70 of *Proceedings of Machine Learning Research*, pp. 3289–3298. PMLR, 2017.
- Ke Sun and Frank Nielsen. A geometric modeling of Occam's razor in deep learning. *Information Geometry*, 2025. Special Issue: Half a Century of Information Geometry, Part 2. Formerly titled "Lightlike neuromanifolds, Occam's razor and deep learning".
- Shiqing Sun and James C. Spall. Connection of diagonal Hessian estimates to natural gradients in stochastic optimization. In *Proceedings of the 55th Annual Conference on Information Sciences and Systems (CISS)*, 2021.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. In Qun Liu and David Schlangen (eds.), *Empirical Methods in Natural Language Processing (EMNLP): System Demonstrations*, pp. 38–45. Association for Computational Linguistics, 2020. https://huggingface.co.

Zhewei Yao, Amir Gholami, Kurt Keutzer, and Michael W Mahoney. PyHessian: Neural networks through the lens of the Hessian. In *IEEE international conference on big data* (*Big Data*), pp. 581–590. IEEE, IEEE Computer Society, 2020.

Zhewei Yao, Amir Gholami, Sheng Shen, Mustafa Mustafa, Kurt Keutzer, and Michael Mahoney. AdaHessian: An adaptive second order optimizer for machine learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 10665–10673, 2021.

Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. Character-level convolutional networks for text classification. *CoRR*, abs/1509.01626, 2015.

A FURTHER ANALYSIS IN THE CORE SPACE

The lemma below gives the average error (variance) of using R(y) to estimate $\mathcal{I}^{\Delta}(z)$, where y is a random variable distributed according to $p(y \mid z)$.

Lemma 12. The element-wise variance of the random matrix R(y), denoted by $Var(R_{ij})$, is given by

$$\operatorname{Var}(R_{ij}) = \begin{cases} p_i(1-p_i)(1-4p_i(1-p_i)) & \text{if } i=j; \\ p_ip_j(p_i+p_j-4p_ip_j) & \text{otherwise.} \end{cases}$$

 $\forall i, j, \operatorname{Var}(R_{ij}) \leq 1/16$. For both diagonal and off-diagonal entries, the coefficient of variation (CV) $\operatorname{Std}(R_{ij})/|\mathcal{I}_{ij}^{\Delta}(z)|$ can be arbitrarily large, where $\operatorname{Std}(\cdot)$ means standard deviation.

By Lemma 12, when using the rank-1 matrix R(y) as an estimator of $\mathcal{I}^{\Delta}(z)$, the absolute error is bounded, but the relative error given by the CV is unbounded. One may alternatively use the rank-2 random matrix $R'(y) = e_{yy} - pp^{\top}$ to estimate $\mathcal{I}^{\Delta}(z)$. Obviously we have $\mathbb{E}(R'(y)) = \operatorname{diag}(p) - pp^{T} = \mathcal{I}^{\Delta}(z)$ and thus R'(y) is unbiased. The variance appears only on the diagonal while all off-diagonal entries are deterministic with zero-variance. This R'(y) is not used in our developments but is of theoretical interest.

B EXPERIMENTS ON SST-2

We compute the diagonal FIM of DistilBERT (Sanh et al., 2019), which is fine-tuned on the Stanford Sentiment Treebank v2 (SST-2) (Socher et al., 2013) for binary sentiment classification. The model is available as distilbert-base-uncased-finetuned-sst-2-english in the Hugging Face library (Wolf et al., 2020). The density of diagonal FIM entries are shown in fig. 2. There are two differences with the AG News experiment in the main text: (1) The number of classes has reduced to C=2; (2) The model is already fine-tuned and the Fisher information is evaluated on a different region in the parameter space compared to the AG News case. Note \mathbb{F}^{LR}_{ii} is very close to and sometimes larger than the value of \mathbb{F}_{ii} . This is because when C=2, the core matrix is already rank-1. And \mathbb{F} and \mathbb{F}^{LR} are essentially different (unbiased) estimators of \mathcal{F} . The scale of the upper bound \mathbb{F}^{DG}_{ii} is much larger than \mathbb{F}_{ii} showing that the bound is loose. All numerical results presented here are performed on a MacBook Pro with Apple M1 CPU and 16GB RAM.

Estimating the underlying true (diagonal) FIM for real DNNs is a challenging problem by itself due to the pathological spectra (Karakida et al., 2021) of $\mathcal F$ and that many of its diagonal entries are close to zero and may require higher precision. We simulate the ground truth $\mathcal F(\theta)$ given by Hutchinson's estimate using the numerical average of 32 different $\mathbb F(\theta)$, each independently computed via a Rademacher vector and an associated backward pass. Because $\mathbb F(\theta)$ is unbiased, the empirical average would uncover the true $\mathcal F(\theta)$ by the law of large numbers. The Monte Carlo FIM estimator is computationally infeasible in this realistic DNN case and less helpful to compute $\mathcal F(\theta)$ — it requires a separate backward pass for each sample in a mini-batch.

Tables 2 and 3 report the mean absolute error (MAE) of the investigated estimators as compared to the "ground truth" for DistilBERT on AG News and SST-2. One can easily observe that \mathbb{F} is more accurate than the empirical FIM $\hat{\mathcal{F}}$. Among the estimators, \mathbb{F}^{LR} gives more accurate estimates than \mathbb{F}^{DG} , which is consistent with our theoretical analysis.

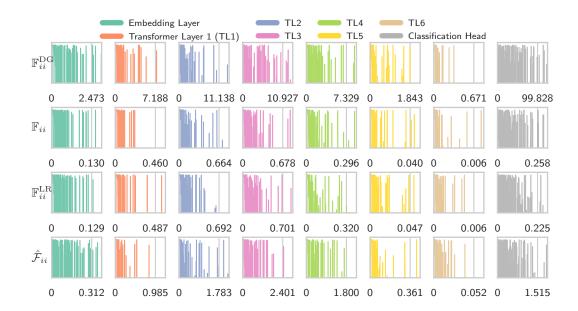


Figure 2: Density plots of diagonal FIM elements based on different approximations (rows) across different layers (columns) on DistilBERT fine-tuned on the SST-2 dataset. The maximum value of the x-axis (number on the bottom right corner) shows the mean value of the FIM multiplied by 2,000. The y-axis means probability density in log-scale.

Table 2: Estimation error of the diagonal FIM for DistilBERT on AG News with two significant figures. The unit of the error in MAE is 10^{-5} .

| | Embedding | TL1 | TL2 | TL3 | TL4 | TL5 | TL6 | Classification Head |
|----------------------------|-----------|------|-----|-----|-----|-----|-----|---------------------|
| \mathbb{F} | 0.11 | 0.73 | 1.4 | 2.2 | 3.9 | 4.9 | 6.5 | 910 |
| $\hat{\mathcal{F}}$ | 0.34 | 6.1 | 13 | 24 | 55 | 77 | 120 | 25,000 |
| \mathbb{F}^{DG} | 0.13 | 0.81 | 1.4 | 2.0 | 3.2 | 3.8 | 5.1 | 660 |
| \mathbb{F}^{LR} | 0.083 | 0.61 | 1.1 | 1.7 | 3.1 | 3.8 | 5.2 | 820 |

Table 3: Estimation error of the diagonal FIM for DistilBERT (fine-tuned) on SST-2 with two significant figures. The unit of the error in MAE is 10^{-5} .

| | Embedding | TL1 | TL2 | TL3 | TL4 | TL5 | TL6 | Classification Head |
|----------------------------|-----------|-----|-----|-----|-----|-----|------|---------------------|
| \mathbb{F} | 2.5 | 8.4 | 15 | 15 | 9.2 | 1.6 | 0.24 | 9.4 |
| $\hat{\mathcal{F}}$ | 19 | 51 | 88 | 110 | 110 | 17 | 2.5 | 62 |
| \mathbb{F}^{DG} | 110 | 290 | 440 | 490 | 360 | 74 | 22 | 3,300 |
| $\mathbb{F}^{	ext{LR}}$ | 1.2 | 7.2 | 11 | 12 | 9.4 | 1.5 | 0.22 | 9.8 |

C ACCURACY OF HUTCHINSON'S ESTIMATE ON DIAGONAL AND LOW RANK CORES

In this section, we show that Hutchinson's estimates $\mathbb{F}^{DG}(\theta)$ and $\mathbb{F}^{LR}(\theta)$ are both unbiased with bounded variances.

Proposition 13. The random matrix $\mathbb{F}^{DG}(\theta)$ is an unbiased estimator of $\mathcal{F}^{DG}(\theta)$. The variance of its diagonal elements is $\operatorname{Var}\left(\mathbb{F}^{DG}_{ii}(\theta)\right) = 2(\mathcal{F}^{DG}_{ii}(\theta))^2 - 2\sum_{x \in \mathcal{D}_x} \sum_{y=1}^C \zeta_y^2(x,\theta)(\frac{\partial z_y}{\partial \theta_i})^4$.

Proposition 14. $\mathbb{F}^{LR}(\theta)$ is an unbiased estimate of $\mathcal{F}^{LR}(\theta)$; the variance of its diagonal elements is $\operatorname{Var}\left(\mathbb{F}^{LR}_{ii}(\theta)\right) = 2(\mathcal{F}^{LR}_{ii}(\theta))^2 - 2\sum_{x \in \mathcal{D}_x} \lambda_C^2(x,\theta) \left(v_C^\top(x,\theta) \frac{\partial z}{\partial \theta_i}\right)^4$.

We have $\operatorname{Std}(\mathbb{F}^{\operatorname{DG}}_{ii}(\theta))/\mathcal{F}^{\operatorname{DG}}_{ii}(\theta) \leq \sqrt{2}$ by Proposition 13, and at the same time, we have $\operatorname{Std}(\mathbb{F}^{\operatorname{LR}}_{ii}(\theta))/\mathcal{F}^{\operatorname{LR}}_{ii}(\theta) \leq \sqrt{2}$ by Proposition 14. Their estimation quality is guaranteed.

D PROOF OF THEOREM 1

Proof. We already know the closed form FIM

$$\mathcal{I}^{\Delta}(z) = \operatorname{diag}(p) - pp^{\top}.$$

Therefore

$$\mathcal{I}^{\Delta}(z)e = (\operatorname{diag}(p) - pp^{\top})e = p - \left(\sum_{i=1}^{C} p_i\right)p = p - p = 0.$$

Therefore $te, t \in \Re$ is a one-dimensional kernel of $\mathcal{I}^{\Delta}(z)$. Since $\mathcal{I}^{\Delta}(z) \succeq 0$, we must have $\lambda_1 = 0$, and $v_1 = e/\|e\|$.

To show the sum of the eigenvalues of $\mathcal{I}^{\Delta}(z)$, we have

$$\sum_{i=1}^{C} \lambda_i = \operatorname{tr}(\mathcal{I}^{\Delta}(z)) = \operatorname{tr}(\operatorname{diag}(p)) - \operatorname{tr}(pp^{\top}) = 1 - \operatorname{tr}(p^{\top}p) = 1 - p^{\top}p = 1 - \|p\|^2.$$

In below, we consider the maximum eigenvalue λ_C . We know that

$$\lambda_C = \sup_{\|u\|=1} u^{\top} \mathcal{I}^{\Delta}(z) u.$$

Therefore

$$\forall i, \quad \lambda_C \ge e_i \mathcal{I}^{\Delta}(z) e_i = \mathcal{I}_{ii}^{\Delta}(z) = p_i (1 - p_i).$$

Therefore $\lambda_C \ge \max_i p_i (1 - p_i)$. At the same time, because $\lambda_1 = 0$, we have

$$\sum_{i=1}^{C} \lambda_i = \lambda_2 + \lambda_3 + \dots + \lambda_C \le (C-1)\lambda_C.$$

Therefore

$$\lambda_C \ge \frac{\sum_{i=1}^C \lambda_i}{C-1} = \frac{1 - \|p\|^2}{C-1}.$$

Because

$$\operatorname{diag}(p) = \mathcal{I}^{\Delta}(z) + pp^{\top}.$$

By the Cauchy's interlacing theorem, we have

$$\lambda_{C-1} \le p_{(C-1)} \le \lambda_C \le p_{(C)}.$$

It remains to prove the upper bounds of λ_C . First, we have

$$\lambda_{C} = \sup_{\|u\|=1} u^{\top} \mathcal{I}^{\Delta}(z) u. = \sup_{\|u\|=1} \left(\sum_{i=1}^{C} p_{i} u_{i}^{2} - (p^{\top} u)^{2} \right)$$

$$\leq \sup_{\|u\|=1} \sum_{i=1}^{C} p_{i} u_{i}^{2} = \max_{i} p_{i} = p_{(C)},$$

which has just been proved using Cauchy's interlacing theorem.

By the Gershgorin circle theorem, λ_C must lie in one of the Gershgorin discs, given by the closed intervals

$$\left[p_i(1-p_i) - \sum_{j \neq i} p_i p_j, \ p_i(1-p_i) + \sum_{j \neq i} p_i p_j \right], \quad i = 1, \dots, C.$$

Therefore

$$\lambda_C \le \max_i \left(p_i (1 - p_i) + \sum_{j \ne i} p_i p_j \right)$$

= $\max_i \left(p_i (1 - p_i) + p_i (1 - p_i) \right) = 2 \max_i p_i (1 - p_i).$

Because $\mathcal{I}^{\Delta}(z) \succeq 0$,

$$\lambda_C \le \sum_{i=1}^C \lambda_i = 1 - \|p\|^2.$$

The statement follows immediately by combining the above lower and upper bounds of λ_C .

E PROOF OF LEMMA 2

Proof. Because $\mathcal{I}^{\Delta}(z) \succeq 0$. All its eigenvalues are greater or equal to 0. We have

$$\mathcal{I}^{\Delta}(z) - \lambda_C v_C v_C^{\top} = \sum_{i=1}^{C-1} \lambda_i v_i v_i^{\top} \succeq 0.$$

To show that $\lambda_C v_C v_C^{\top}$ is the best rank-1 representation. Assume that $\exists u \neq 0$, such that $\mathcal{I}^{\Delta}(z) \succeq uu^{\top} \succeq \lambda_C v_C v_C^{\top}$. Then

$$v_C^{\top} \mathcal{I}^{\Delta}(z) v_C = \lambda_C \ge (v_C^{\top} u)^2 \ge \lambda_C.$$

Therefore

$$v_C^{\top} u = \pm \sqrt{\lambda_C}.$$

Assume that $u = \sum_{i=1}^{C} \alpha_i v_i$, then $\alpha_C = v_C^{\top} u = \pm \sqrt{\lambda_C}$. Moreover, we have

$$\lambda_C \ge \frac{u^\top}{\|u\|} \mathcal{I}^{\Delta}(z) \frac{u}{\|u\|} \ge \frac{u^\top}{\|u\|} u u^\top \frac{u}{\|u\|} = \|u\|^2 = \sum_{i=1}^C \alpha_i^2.$$

Therefore $\forall i \neq C, \, \alpha_i = 0$. In summary, $u = \pm \sqrt{\lambda_C} v_C$. Hence, $uu^\top = \lambda_C v_C v_C^\top$.

We have

$$\operatorname{diag}(p) - \mathcal{I}^{\Delta}(z) = \operatorname{diag}(p) - (\operatorname{diag}(p) - pp^{\top}) = pp^{\top} \succeq 0.$$

Therefore $\mathrm{diag}\,(p)\succeq\mathcal{I}^{\Delta}(z).$ Assume that $\mathrm{diag}\,(q)$ satisfies

$$\mathcal{I}^{\Delta}(z) \leq \operatorname{diag}(q) \leq \operatorname{diag}(p)$$
.

Then

$$\operatorname{diag}(p) - \mathcal{I}^{\Delta}(z) = pp^{\top} \succeq \operatorname{diag}(q) - \mathcal{I}^{\Delta}(z) \succeq 0.$$

Therefore

$$\operatorname{diag}(q) - \mathcal{I}^{\Delta}(z) = \beta p p^{\top} (\beta \le 1).$$

Consequently,

$$\operatorname{diag}(q) = \mathcal{I}^{\Delta}(z) + \beta p p^{\top} = \operatorname{diag}(p) - p p^{\top} + \beta p p^{\top} = \operatorname{diag}(p) + (\beta - 1) p p^{\top}.$$

Therefore all off-diagonal entries of $(\beta-1)pp^{\top}$ are zero. We must have $\beta=1$ and thus diag (q)= diag (p).

F PROOF OF LEMMA 3

Proof.

$$\|\lambda_C v_C v_C^{\top} - \mathcal{I}^{\Delta}(z)\| = \|\sum_{i=1}^{C-1} \lambda_i v_i v_i^{\top}\| = \sqrt{\sum_{i=1}^{C-1} \lambda_i^2} \le \sqrt{(\sum_{i=1}^{C-1} \lambda_i)^2}$$
$$= \sum_{i=1}^{C-1} \lambda_i = \operatorname{tr}(\mathcal{I}^{\Delta}(z)) - \lambda_C = 1 - \|p\|^2 - \lambda_C.$$

By Theorem 1, we have $\lambda_C \geq p_{(C-1)}$. Therefore

$$\|\lambda_C v_C v_C^\top - \mathcal{I}^{\Delta}(z)\| \le 1 - \|p\|^2 - p_{(C-1)}.$$

By Cauchy's interlacing theorem (see our proof of Theorem 1), we have

$$\forall i \in \{1, \dots, C-1\}, \quad \lambda_i \leq p_{(i)}.$$

Hence

$$\|\lambda_C v_C v_C^\top - \mathcal{I}^\Delta(z)\| = \sqrt{\sum_{i=1}^{C-1} \lambda_i^2} = \sqrt{\sum_{i=2}^{C-1} \lambda_i^2} \le \sqrt{\sum_{i=2}^{C-1} p_{(i)}^2}.$$

The statement follows immediately by combining the above upper bounds.

G Proof of Lemma 4

Proof. The spectrum of R(y) is

$$0 \le \dots \le 0 \le ||e_y - p||^2$$
.

The spectrum of $\mathcal{I}^{\Delta}(z)$, by our assumption, is

$$\lambda_1 \leq \cdots \leq \lambda_{C-1} \leq \lambda_C$$
.

By Hoffman-Wielandt inequality, we have $\forall z \in \Delta^{C-1}, y \in \{1, \dots, C\}$

$$||R(y) - \mathcal{I}^{\Delta}(z)|| \ge \sqrt{\sum_{i=1}^{C-1} \lambda_i^2 + (\lambda_C - ||e_y - p||^2)^2}$$

$$\ge |\lambda_C - ||e_y - p||^2|$$

$$= |\lambda_C - e_y^{\top} e_y - p^{\top} p + 2e_y^{\top} p|$$

$$= |\lambda_C - 1 - ||p||^2 + 2p_y|$$

$$= \max\{\lambda_C - 1 - ||p||^2 + 2p_y, \ 1 + ||p||^2 - \lambda_C - 2p_y\}.$$

By Theorem 1, we have $\lambda_C \leq 1 - \|p\|^2$. One can choose y so that $p_y = p_{(1)}$, then

$$||R(y) - \mathcal{I}^{\Delta}(z)|| \ge 1 + ||p||^2 - \lambda_C - 2p_{(1)}$$

$$\ge 1 + ||p||^2 - (1 - ||p||^2) - 2p_{(1)}$$

$$= 2||p||^2 - 2p_{(1)}.$$

H Proof of Lemma 12

Proof. We first look at the diagonal entries of R. We have

$$R_{ii} = ([y = i] - p_i)^2 = \begin{cases} (1 - p_i)^2 & \text{if } y = i; \\ p_i^2 & \text{otherwise.} \end{cases}$$

Therefore

$$\mathbb{E}(R_{ii}) = p_i(1 - p_i)^2 + (1 - p_i)p_i^2 = p_i(1 - p_i) = \mathcal{I}_{ii}^{\Delta}(z).$$

This shows that R_{ii} is an unbiased estimator of the diagonal entries of $\mathcal{I}^{\Delta}(z)$. We have

$$\mathbb{E}(R_{ii}^2) = p_i (1 - p_i)^4 + (1 - p_i) p_i^4 = p_i (1 - p_i) \left[(1 - p_i)^3 + p_i^3 \right]$$
$$= p_i (1 - p_i) \left[(1 - p_i)^2 - p_i (1 - p_i) + p_i^2 \right].$$

Therefore

$$Var(R_{ii}) = \mathbb{E}(R_{ii}^2) - (\mathbb{E}(R_{ii}))^2$$

$$= p_i(1 - p_i) \left[(1 - p_i)^2 - p_i(1 - p_i) + p_i^2 \right] - p_i^2(1 - p_i)^2$$

$$= p_i(1 - p_i) \left[(1 - p_i)^2 - 2p_i(1 - p_i) + p_i^2 \right]$$

$$= p_i(1 - p_i)(1 - 4p_i(1 - p_i))$$

$$= \mathcal{I}_{ii}^{\Delta}(z)(1 - 4\mathcal{I}_{ii}^{\Delta}(z))$$

$$= -4 \left(\mathcal{I}_{ii}^{\Delta}(z) - \frac{1}{8} \right)^2 + \frac{1}{16} \le \frac{1}{16}.$$

The coefficient of variation (CV)

$$\frac{\sqrt{\operatorname{Var}(R_{ii})}}{\mathcal{I}_{ii}^{\Delta}(z)} = \sqrt{\frac{\mathcal{I}_{ii}^{\Delta}(z)(1 - 4\mathcal{I}_{ii}^{\Delta}(z))}{\mathcal{I}_{ii}^{\Delta}(z)^2}} = \sqrt{\frac{1}{\mathcal{I}_{ii}^{\Delta}(z)} - 4}$$

is unbounded. As $\mathcal{I}_{ii}^{\Delta}(z) \to 0$, the CV can take arbitrarily large value.

Next, we consider the off-diagonal entries of R. For $i \neq j$, we have

$$R_{ij} = ([[y = i]] - p_i)([[y = j]] - p_j)$$

= $p_i p_j - [[y = i]] p_j - [[y = j]] p_i$.

Hence.

$$\mathbb{E}(R_{ij}) = p_i p_j - p_j p_j - p_j p_i = -p_i p_j = \mathcal{I}_{ij}^{\Delta}(z).$$

At the same time,

$$\mathbb{E}(R_{ij}^{2}) = \mathbb{E}(p_{i}p_{j} - [y = i]p_{j} - [y = j]p_{i})^{2}$$

$$= p_{i}^{2}p_{j}^{2} + \mathbb{E}([y = i]p_{j}^{2} + [y = j]p_{i}^{2} - 2[y = i]p_{i}p_{j}^{2} - 2[y = j]p_{i}^{2}p_{j})$$

$$= p_{i}^{2}p_{j}^{2} + p_{i}p_{j}^{2} + p_{j}p_{i}^{2} - 2p_{i}^{2}p_{j}^{2} - 2p_{i}^{2}p_{j}^{2}$$

$$= p_{i}p_{j}^{2} + p_{i}^{2}p_{j} - 3p_{i}^{2}p_{j}^{2}$$

$$= p_{i}p_{j}(p_{i} + p_{j} - 3p_{i}p_{j}).$$

Therefore

$$Var(R_{ij}) = \mathbb{E}(R_{ij}^2) - (\mathbb{E}(R_{ij}))^2$$

$$= p_i p_j (p_i + p_j - 3p_i p_j) - p_i^2 p_j^2$$

$$= p_i p_j (p_i + p_j - 4p_i p_j)$$

$$\leq p_i p_j (1 - 4p_i p_j)$$

$$= -4 \left(p_i p_j - \frac{1}{8} \right)^2 + \frac{1}{16} \leq \frac{1}{16}.$$

The coefficient of variation

$$\frac{\sqrt{\text{Var}(R_{ij})}}{|\mathcal{I}_{ij}^{\Delta}(z)|} = \sqrt{\frac{p_i p_j (p_i + p_j - 4p_i p_j)}{p_i^2 p_j^2}} = \sqrt{\frac{1}{p_i} + \frac{1}{p_j} - 4}$$

is unbounded. As either $p_i \to 0$, or $p_j \to 0$, the CV can take arbitrarily large value.

I Proof of Proposition 5

Proof. By Lemma 2, we have

$$\lambda_C v_C v_C^{\top} \leq \mathcal{I}^{\Delta}(z) \leq \operatorname{diag}(p)$$
.

Therefore

$$\forall x, \theta \quad \left(\frac{\partial z}{\partial \theta}\right)^{\top} \lambda_C v_C v_C^{\top} \frac{\partial z}{\partial \theta} \preceq \left(\frac{\partial z}{\partial \theta}\right)^{\top} \mathcal{I}^{\Delta}(z(x, \theta)) \frac{\partial z}{\partial \theta} \preceq \left(\frac{\partial z}{\partial \theta}\right)^{\top} \operatorname{diag}(p) \frac{\partial z}{\partial \theta}.$$

Therefore

$$\forall \theta \quad \sum_{x \in \mathcal{D}_x} \lambda_C \left(\frac{\partial z}{\partial \theta} \right)^\top v_C v_C^\top \frac{\partial z}{\partial \theta} \preceq \sum_{x \in \mathcal{D}_x} \left(\frac{\partial z}{\partial \theta} \right)^\top \mathcal{I}^\Delta(z(x,\theta)) \frac{\partial z}{\partial \theta} \preceq \sum_{x \in \mathcal{D}_x} \sum_{i=1}^C p_i \frac{\partial z_i}{\partial \theta} \frac{\partial z_i}{\partial \theta^\top}.$$

J Proof of Corollary 6

Proof. We first prove the upper bound. By Proposition 5, we have

$$\mathcal{F}^{\Delta}(\theta) \preceq \sum_{x \in \mathcal{D}_x} \sum_{i=1}^{C} p_i \frac{\partial z_i}{\partial \theta} \frac{\partial z_i}{\partial \theta^{\top}}.$$

Taking trace on both sides, we get

$$\operatorname{tr}(\mathcal{F}^{\Delta}(\theta)) \leq \sum_{x \in \mathcal{D}_{x}} \sum_{i=1}^{C} p_{i} \operatorname{tr}\left(\frac{\partial z_{i}}{\partial \theta} \frac{\partial z_{i}}{\partial \theta^{\top}}\right)$$

$$= \sum_{x \in \mathcal{D}_{x}} \sum_{i=1}^{C} p_{i} \operatorname{tr}\left(\frac{\partial z_{i}}{\partial \theta^{\top}} \frac{\partial z_{i}}{\partial \theta}\right)$$

$$= \sum_{x \in \mathcal{D}_{x}} \sum_{i=1}^{C} p_{i} \frac{\partial z_{i}}{\partial \theta^{\top}} \frac{\partial z_{i}}{\partial \theta}$$

$$= \sum_{x \in \mathcal{D}_{x}} \sum_{i=1}^{C} p_{i} \left\|\frac{\partial z_{i}}{\partial \theta}\right\|^{2}.$$

The lower bound is not straightforward from Proposition 5. By Eq. (2), we have

$$\operatorname{tr}(\mathcal{F}^{\Delta}(\theta)) = \sum_{x \in \mathcal{D}_x} \operatorname{tr}\left[\left(\frac{\partial z}{\partial \theta}\right)^{\top} \mathcal{I}^{\Delta}(z) \frac{\partial z}{\partial \theta}\right] = \sum_{x \in \mathcal{D}_x} \operatorname{tr}\left[\frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta}\right)^{\top} \mathcal{I}^{\Delta}(z)\right].$$

Note that $\frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta} \right)^{\top}$ is a $C \times C$ matrix with sorted eigenvalues $\sigma_1^2(x,\theta) \leq \cdots \leq \sigma_C^2(x,\theta)$. By Theorem 1, $\mathcal{I}^{\Delta}(z)$ is another $C \times C$ matrix with sorted eigenvalues $0 = \lambda_1(x,\theta) \leq \cdots \leq \lambda_C(x,\theta)$. Applying the Von Neumann trace inequality, we get

$$\operatorname{tr}(\mathcal{F}^{\Delta}(\theta)) \ge \sum_{x \in \mathcal{D}_x} \sum_{i=2}^{C} \lambda_i(x, \theta) \sigma_{C-i+1}^2(x, \theta) \ge \sum_{x \in \mathcal{D}_x} \lambda_C(x, \theta) \sigma_1^2(x, \theta).$$

The last " \geq " is because all terms $\lambda_i(x,\theta)\sigma_{C-i+1}^2(x,\theta)$ are non-negative.

K Proof of Proposition 7

Proof. Denote the singular values of $\frac{\partial z}{\partial \theta}$ as $0 \leq \sigma_1 \leq \cdots \leq \sigma_C$. Then the eigenvalues of the $C \times C$ Hermitian matrix $\frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta} \right)^{\top}$ is $\sigma_1^2 \leq \cdots \leq \sigma_C^2$.

To prove the upper bound, we have

$$\begin{split} & \left\| \sum_{x \in \mathcal{D}_x} \sum_{i=1}^C p_i \left(\frac{\partial z_i}{\partial \theta} \right)^\top \frac{\partial z_i}{\partial \theta} - \mathcal{F}^{\Delta}(\theta) \right\| \\ &= \left\| \sum_{x \in \mathcal{D}_x} \left(\frac{\partial z}{\partial \theta} \right)^\top \left(\operatorname{diag}\left(p \right) - \operatorname{diag}\left(p \right) + p p^\top \right) \frac{\partial z}{\partial \theta} \right\| \\ &= \left\| \sum_{x \in \mathcal{D}_x} \left(\frac{\partial z}{\partial \theta} \right)^\top p p^\top \frac{\partial z}{\partial \theta} \right\| \\ &\leq \sum_{x \in \mathcal{D}_x} \sqrt{\operatorname{tr} \left[\left(\frac{\partial z}{\partial \theta} \right)^\top p p^\top \frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta} \right)^\top p p^\top \frac{\partial z}{\partial \theta} \right]} \\ &= \sum_{x \in \mathcal{D}_x} \sqrt{\operatorname{tr} \left[p^\top \frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta} \right)^\top p p^\top \frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta} \right)^\top p \right]} \\ &\leq \sum_{x \in \mathcal{D}_x} \sqrt{\left[p^\top \frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta} \right)^\top p \right]^2} \\ &= \sum_{x \in \mathcal{D}_x} p^\top \frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta} \right)^\top p} \\ &= \sum_{x \in \mathcal{D}_x} \|p\|^2 \cdot \frac{p^\top}{\|p\|} \frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta} \right)^\top \frac{p}{\|p\|} \\ &\leq \sum_{x \in \mathcal{D}_x} \|p\|^2 \sigma_C^2. \end{split}$$

Now we are ready to prove the lower bound. From the above, we have

$$\left\| \sum_{x \in \mathcal{D}_x} \sum_{i=1}^C p_i \left(\frac{\partial z_i}{\partial \theta} \right)^\top \frac{\partial z_i}{\partial \theta} - \mathcal{F}^{\Delta}(\theta) \right\| = \left\| \sum_{x \in \mathcal{D}_x} \left(\frac{\partial z}{\partial \theta} \right)^\top p p^\top \frac{\partial z}{\partial \theta} \right\|.$$

Denote $\omega(x)\coloneqq \left(\frac{\partial z}{\partial \theta}\right)^{\top} p$. Then

$$\left\| \sum_{x \in \mathcal{D}_x} \sum_{i=1}^C p_i \left(\frac{\partial z_i}{\partial \theta} \right)^\top \frac{\partial z_i}{\partial \theta} - \mathcal{F}^{\Delta}(\theta) \right\| = \left\| \sum_{x \in \mathcal{D}_x} \omega(x) \omega(x)^\top \right\|$$

$$= \sqrt{\operatorname{tr} \left(\left(\sum_{x \in \mathcal{D}_x} \omega(x) \omega(x)^\top \right)^2 \right)}$$

$$\geq \sqrt{\sum_{x \in \mathcal{D}_x} (\omega(x)^\top \omega(x))^2}$$

$$= \sqrt{\sum_{x \in \mathcal{D}_x} \|\omega(x)\|^4}.$$

The last ">" is due to

$$\operatorname{tr}\left(\omega(x)\omega(x)^{\top}\omega(x')\omega(x')^{\top}\right) = \operatorname{tr}\left(\omega(x')^{\top}\omega(x)\omega(x)^{\top}\omega(x')\right) = (\omega(x')^{\top}\omega(x))^{2} \geq 0.$$

L PROOF OF PROPOSITION 8

 Proof. We can first have a loose bound:

$$\left\| \sum_{x \in \mathcal{D}_{x}} \lambda_{C} \left(\frac{\partial z}{\partial \theta} \right)^{\top} v_{C} v_{C}^{\top} \frac{\partial z}{\partial \theta} - \mathcal{F}^{\Delta}(\theta) \right\|$$

$$= \left\| \sum_{x \in \mathcal{D}_{x}} \lambda_{C} \left(\frac{\partial z}{\partial \theta} \right)^{\top} v_{C} v_{C}^{\top} \frac{\partial z}{\partial \theta} - \sum_{x \in \mathcal{D}_{x}} \left(\frac{\partial z}{\partial \theta} \right)^{\top} \mathcal{I}^{\Delta}(z) \frac{\partial z}{\partial \theta} \right\|$$

$$= \left\| \sum_{x \in \mathcal{D}_{x}} \left(\frac{\partial z}{\partial \theta} \right)^{\top} \left(\sum_{i=1}^{C-1} \lambda_{i} v_{i} v_{i}^{\top} \right) \frac{\partial z}{\partial \theta} \right\|$$

$$\leq \left\| \sum_{x \in \mathcal{D}_{x}} p_{(C-1)} \left(\frac{\partial z}{\partial \theta} \right)^{\top} \frac{\partial z}{\partial \theta} \right\| \quad \text{(Due to that } \sum_{i=1}^{C-1} \lambda_{i} v_{i} v_{i}^{\top} \preceq p_{(C-1)} I \text{)}$$

$$\leq \sum_{x \in \mathcal{D}_{x}} p_{(C-1)} \left\| \frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta} \right)^{\top} \right\|.$$

The eigenvalues of $\left(\frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta}\right)^{\top}\right)^2$ are $\sigma_1^4 \leq \cdots \leq \sigma_C^4$. We have

$$\begin{split} & \left\| \left(\frac{\partial z}{\partial \theta} \right)^{\top} \left(\sum_{i=1}^{C-1} \lambda_{i} v_{i} v_{i}^{\top} \right) \frac{\partial z}{\partial \theta} \right\|^{2} \\ &= \operatorname{tr} \left[\left(\frac{\partial z}{\partial \theta} \right)^{\top} \left(\sum_{i=1}^{C-1} \lambda_{i} v_{i} v_{i}^{\top} \right) \frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta} \right)^{\top} \left(\sum_{i=1}^{C-1} \lambda_{i} v_{i} v_{i}^{\top} \right) \frac{\partial z}{\partial \theta} \right] \\ &= \operatorname{tr} \left[\left(\frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta} \right)^{\top} \left(\sum_{i=1}^{C-1} \lambda_{i} v_{i} v_{i}^{\top} \right) \right)^{2} \right] \\ &\leq \operatorname{tr} \left[\left(\frac{\partial z}{\partial \theta} \left(\frac{\partial z}{\partial \theta} \right)^{\top} \right)^{2} \left(\sum_{i=1}^{C-1} \lambda_{i}^{2} v_{i} v_{i}^{\top} \right) \right] \quad \text{Due to } \operatorname{tr}(AB)^{2} \leq \operatorname{tr}(A^{2}B^{2}) \\ &\leq \sum_{i=1}^{C-1} \sigma_{i+1}^{4} \lambda_{i}^{2}. \end{split}$$

The last " \leq " is due to Von Neumann's trace inequality, and that the smallest two eigenvalues of the matrix $\sum_{i=1}^{C-1} \lambda_i^2 v_i v_i^{\mathsf{T}}$ are both zero. We also have the Cauchy interlacing

$$\lambda_2 \le p_{(2)} \le \lambda_3 \le p_{(3)} \le \dots \le \lambda_{C-1} \le p_{(C-1)}.$$

To sum up,

$$\left\| \sum_{x \in \mathcal{D}_x} \lambda_C \left(\frac{\partial z}{\partial \theta} \right)^\top v_C v_C^\top \frac{\partial z}{\partial \theta} - \mathcal{F}^{\Delta}(\theta) \right\|$$

$$\leq \sum_{x \in \mathcal{D}_x} \left\| \left(\frac{\partial z}{\partial \theta} \right)^\top \left(\sum_{i=1}^{C-1} \lambda_i v_i v_i^\top \right) \frac{\partial z}{\partial \theta} \right\|$$

$$\leq \sum_{x \in \mathcal{D}_x} \sqrt{\sum_{i=2}^{C-1} \sigma_{i+1}^4 \lambda_i^2}$$

$$\leq \sum_{x \in \mathcal{D}_x} \sqrt{\sum_{i=2}^{C-1} \sigma_{i+1}^4 p_{(i)}^2}.$$

If one relax $\forall i \in \{2, \dots, C-1\}$, $p_{(i)} \leq p_{(C-1)}$, then we get the loose bound proved earlier.

M Proof of Proposition 9

Proc

$$\|\mathcal{F}(\theta) - \hat{\mathcal{F}}^{\Delta}(\theta)\|_{\sigma} = \left\| \sum_{x \in \mathcal{D}_{x}} \left(\frac{\partial z}{\partial \theta} \right)^{\top} \cdot \mathcal{I}(z(x,\theta)) \cdot \frac{\partial z}{\partial \theta} - \sum_{x \in \mathcal{D}_{x}} \left(\frac{\partial z}{\partial \theta} \right)^{\top} \left(e_{y} - p)(e_{y} - p)^{\top} \frac{\partial z}{\partial \theta} \right\|_{\sigma}$$

$$= \left\| \sum_{x \in \mathcal{D}_{x}} \left(\frac{\partial z}{\partial \theta} \right)^{\top} \left[\operatorname{diag}(p) - pp^{\top} - (e_{y} - p)(e_{y} - p)^{\top} \right] \frac{\partial z}{\partial \theta} \right\|_{\sigma}$$

$$\leq \sum_{x \in \mathcal{D}_{x}} \left\| \left(\frac{\partial z}{\partial \theta} \right)^{\top} \left[\operatorname{diag}(p) - pp^{\top} - (e_{y} - p)(e_{y} - p)^{\top} \right] \frac{\partial z}{\partial \theta} \right\|_{\sigma}$$

$$\leq \sum_{x \in \mathcal{D}_{x}} \left\| \frac{\partial z}{\partial \theta} \right\|_{\sigma} \left\| \operatorname{diag}(p) - pp^{\top} - (e_{y} - p)(e_{y} - p)^{\top} \right\|_{\sigma} \left\| \frac{\partial z}{\partial \theta} \right\|_{\sigma}$$

$$= \sum_{x \in \mathcal{D}_{x}} \sigma_{C}^{2} \left\| \operatorname{diag}(p) - pp^{\top} - (e_{y} - p)(e_{y} - p)^{\top} \right\|_{\sigma}.$$

Now we examine the matrix diag $(p) - pp^{\top} - (e_y - p)(e_y - p)^{\top}$. By Theorem 1, the spectrum of diag $(p) - pp^{\top}$ is

$$\lambda_1 = 0 \le \lambda_2 \le \dots \le \lambda_C$$
.

By Cauchy interlacing theorem, the spectrum of diag $(p) - pp^{\top} - (e_y - p)(e_y - p)^{\top}$, given by $\lambda'_1, \dots, \lambda'_C$, must satisfy

$$\lambda_1' \le \lambda_1 = 0 \le \lambda_2' \le \lambda_2 \le \dots \le \lambda_C' \le \lambda_C.$$

with at least one eigenvalue that is not positive: $\lambda_1' \leq 0$. Therefore

$$\|\operatorname{diag}(p) - pp^{\top} - (e_y - p)(e_y - p)^{\top}\|_{\sigma} \le \max\{-\lambda_1', \lambda_C\}.$$

We also have

$$\begin{split} \lambda_1' &= \inf_{u:||u||=1} u^\top \left[\operatorname{diag}(p) - pp^\top - (e_y - p)(e_y - p)^\top \right] u \\ &\geq \inf_{u:||u||=1} - u^\top \left[(e_y - p)(e_y - p)^\top \right] u \\ &= -(e_y - p)^\top (e_y - p) \\ &= -(1 + p^\top p - 2p_y) \\ &= 2p_y - 1 - ||p||^2. \end{split}$$

Therefore

$$\begin{aligned} \left\| \operatorname{diag}(p) - pp^{\top} - (e_y - p)(e_y - p)^{\top} \right\|_{\sigma} &\leq \max\{1 + \|p\|^2 - 2p_y, \lambda_C\} \\ &\leq \max\{1 + \|p\|^2 - 2p_y, 1 - \|p\|^2\} \\ &< 1 + \|p\|^2. \end{aligned}$$

In summary,

$$\|\mathcal{F}(\theta) - \hat{\mathcal{F}}^{\Delta}(\theta)\|_{\sigma} \le \sum_{x \in \mathcal{D}_x} \sigma_C^2 (1 + \|p\|^2).$$

N PROOF OF PROPOSITION 10

Proof.

$$\begin{split} & \left\| \left(\frac{\partial z}{\partial \theta} \right)^{\top} \cdot \mathcal{I}^{\Delta}(z(x,\theta)) \cdot \frac{\partial z}{\partial \theta} - \left(\frac{\partial z}{\partial \theta} \right)^{\top} \cdot \hat{\mathcal{I}}^{\Delta}(z(x,\theta)) \cdot \frac{\partial z}{\partial \theta} \right\|_{\sigma} \\ & \geq \left\| \left(\frac{\partial z}{\partial \theta} \right)^{\top} \cdot \left[\mathcal{I}^{\Delta}(z(x,\theta)) - \hat{\mathcal{I}}^{\Delta}(z(x,\theta)) \right] \cdot \frac{\partial z}{\partial \theta} \right\|_{\sigma} \\ & = \left\| \left(\frac{\partial z}{\partial \theta} \right)^{\top} \cdot \left[\operatorname{diag}(p) - pp^{\top} - (e_{y} - p)(e_{y} - p)^{\top} \right] \cdot \frac{\partial z}{\partial \theta} \right\|_{\sigma} \\ & = \sup_{u: \|u\| = 1} \left| \left(\frac{\partial z}{\partial \theta} u \right)^{\top} \cdot \left[\operatorname{diag}(p) - pp^{\top} - (e_{y} - p)(e_{y} - p)^{\top} \right] \cdot \left(\frac{\partial z}{\partial \theta} u \right) \right| \\ & \geq \sup_{v: \|v\| = 1} \left| \sigma_{(1)} v \cdot \left[\operatorname{diag}(p) - pp^{\top} - (e_{y} - p)(e_{y} - p)^{\top} \right] \cdot \sigma_{(1)} v \right| \\ & \geq \sigma_{(1)}^{2} \left\| \operatorname{diag}(p) - pp^{\top} - (e_{y} - p)(e_{y} - p)^{\top} \right\|_{\sigma} \\ & \geq \sigma_{(1)}^{2} \left| \left(\frac{e_{y} - p}{\|e_{y} - p\|} \right)^{\top} \left((e_{y} - p)(e_{y} - p)^{\top} - \lambda_{C} \right) \frac{e_{y} - p}{\|e_{y} - p\|} \right| \\ & = \sigma_{(1)}^{2} \left| \left| 1 + \|p\|^{2} - \lambda_{C} - 2p_{y} \right|. \end{split}$$

We choose $p_y = p_{(1)}$, therefore $\exists y$, such that

$$\begin{split} & \left\| \left(\frac{\partial z}{\partial \theta} \right)^\top \cdot \mathcal{I}^{\Delta}(z(x,\theta)) \cdot \frac{\partial z}{\partial \theta} - \left(\frac{\partial z}{\partial \theta} \right)^\top \cdot \hat{\mathcal{I}}^{\Delta}(z(x,\theta)) \cdot \frac{\partial z}{\partial \theta} \right\|_{\sigma} \\ \ge & \sigma_{(1)}^2 \left| 1 + \|p\|^2 - \lambda_C - 2p_{(1)} \right|. \end{split}$$

O Proof of Proposition 11

Proof. From the derivations in the main text, we already know that $\mathbb{E}_{p(\xi)} \mathbb{I}(\theta) = \mathcal{I}(\theta)$. To show the estimator variance, we first consider the case when $p(\xi)$ is a standard multivariate Gaussian distribution. First we note that both $\mathfrak{h}(\mathcal{D}_x,\theta)$ and $\partial \mathfrak{h}/\partial \theta_i$ are in the form of a sum of independent Gaussian random variables. Hence,

$$\frac{\partial \mathfrak{h}}{\partial \theta_i} = \sum_{x \in \mathcal{D}_x} \sum_{y=1}^C \sqrt{p(y \mid x, \theta)} \frac{\partial \ell_{xy}}{\partial \theta_i} \xi_{xy} \sim G\left(0, \sum_{x \in \mathcal{D}_x} \sum_{y=1}^C p(y \mid x, \theta) \left(\frac{\partial \ell_{xy}}{\partial \theta_i}\right)^2\right).$$

Therefore

$$\begin{split} & \underset{p(\xi)}{\mathbb{E}} \left(\frac{\partial \mathfrak{h}}{\partial \theta_i} \right)^2 = \sum_{x \in \mathcal{D}_x} \sum_{y=1}^C p(y \, | \, x, \theta) \left(\frac{\partial \ell_{xy}}{\partial \theta_i} \right)^2 = \mathcal{I}_{ii}(\theta); \\ & \underset{p(\xi)}{\mathbb{E}} \left(\frac{\partial \mathfrak{h}}{\partial \theta_i} \right)^4 = 3\mathcal{I}_{ii}^2(\theta). \end{split}$$

Therefore

$$\operatorname{Var}(\mathbb{I}(\theta_i)) = \underset{p(\xi)}{\mathbb{E}} \left(\frac{\partial \mathfrak{h}}{\partial \theta_i} \right)^4 - \mathcal{I}_{ii}^2(\theta) = 2\mathcal{I}_{ii}^2(\theta).$$

We now consider that $p(\xi)$ is Rademacher.

$$\operatorname{Var}(\mathbb{I}(\theta_{i})) = \underset{p(\xi)}{\mathbb{E}} \left(\frac{\partial \mathfrak{h}}{\partial \theta_{i}}\right)^{4} - \left(\mathbb{E}\left(\frac{\partial \mathfrak{h}}{\partial \theta_{i}}\right)^{2}\right)^{2}$$

$$= \underset{p(\xi)}{\mathbb{E}} \left(\frac{\partial \mathfrak{h}}{\partial \theta_{i}}\right)^{4} - \mathcal{I}_{ii}^{2}(\theta)$$

$$= \underset{p(\xi)}{\mathbb{E}} \left(\sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} \sqrt{p(y \mid x, \theta)} \frac{\partial \ell_{xy}}{\partial \theta_{i}} \xi_{xy}\right)^{4} - \mathcal{I}_{ii}^{2}(\theta)$$

$$= \underset{x \in \mathcal{D}_{x}}{\sum} \sum_{y=1}^{C} p^{2}(y \mid x, \theta) \left(\frac{\partial \ell_{xy}}{\partial \theta_{i}}\right)^{4}$$

$$+ 3 \underset{(x,y) \neq (x',y')}{\sum} p(y \mid x, \theta) \left(\frac{\partial \ell_{xy}}{\partial \theta_{i}}\right)^{2} p(y' \mid x', \theta) \left(\frac{\partial \ell_{x'y'}}{\partial \theta_{i}}\right)^{2} - \mathcal{I}_{ii}^{2}(\theta).$$

Note that

$$\begin{split} \mathcal{I}_{ii}^{2}(\theta) &= \left(\sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} p(y \mid x, \theta) \left(\frac{\partial \ell_{xy}}{\partial \theta_{i}}\right)^{2}\right)^{2} \\ &= \sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} p^{2}(y \mid x, \theta) \left(\frac{\partial \ell_{xy}}{\partial \theta_{i}}\right)^{4} + \sum_{(x,y) \neq (x',y')} p(y \mid x, \theta) \left(\frac{\partial \ell_{xy}}{\partial \theta_{i}}\right)^{2} p(y' \mid x', \theta) \left(\frac{\partial \ell_{x'y'}}{\partial \theta_{i}}\right)^{2}. \end{split}$$

Hence,

$$\operatorname{Var}(\mathbb{I}(\theta_{i})) = 3\mathcal{I}_{ii}^{2}(\theta) - 2\sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} p^{2}(y \mid x, \theta) \left(\frac{\partial \ell_{xy}}{\partial \theta_{i}}\right)^{4} - \mathcal{I}_{ii}^{2}(\theta)$$
$$= 2\mathcal{I}_{ii}^{2}(\theta) - 2\sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} p^{2}(y \mid x, \theta) \left(\frac{\partial \ell_{xy}}{\partial \theta_{i}}\right)^{4}.$$

P Proof of Proposition 13

Proof.

$$\mathbb{E}_{p(\xi)}(\mathbb{F}^{\mathrm{DG}}(\theta)) = \mathbb{E}_{p(\xi)}\left(\frac{\partial \mathfrak{h}^{\mathrm{DG}}}{\partial \theta} \frac{\partial \mathfrak{h}^{\mathrm{DG}}}{\partial \theta^{\top}}\right) \\
= \mathbb{E}_{p(\xi)}\left(\sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} \sqrt{\zeta_{y}(x,\theta)} \frac{\partial z_{y}}{\partial \theta} \xi_{xy} \sum_{x' \in \mathcal{D}_{x}} \sum_{y'=1}^{C} \sqrt{\zeta_{y'}(x',\theta)} \frac{\partial z_{y'}}{\partial \theta^{\top}} \xi_{x'y'}\right) \\
= \sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} \sum_{x' \in \mathcal{D}_{x}} \sum_{y'=1}^{C} \sqrt{\zeta_{y}(x,\theta)} \sqrt{\zeta_{y'}(x',\theta)} \frac{\partial z_{y}}{\partial \theta} \frac{\partial z_{y'}}{\partial \theta^{\top}} \mathbb{E}_{p(\xi)}\left(\xi_{xy}\xi_{x'y'}\right) \\
= \sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} \zeta_{y}(x,\theta) \frac{\partial z_{y}}{\partial \theta} \frac{\partial z_{y}}{\partial \theta^{\top}} \\
= \sum_{x \in \mathcal{D}_{x}} \left(\frac{\partial z}{\partial \theta}\right)^{\top} \mathcal{I}^{\mathrm{DG}}(z(x,\theta)) \frac{\partial z}{\partial \theta} \\
= \mathcal{F}^{\mathrm{DG}}(\theta).$$

Therefore,

$$\mathbb{E}_{p(\xi)} \left(\mathbb{F}_{ii}^{\mathrm{DG}}(\theta) \right) = \mathbb{E}_{p(\xi)} \left(\frac{\partial \mathfrak{h}^{\mathrm{DG}}}{\partial \theta_{i}} \right)^{2} = \sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} \zeta_{y}(x,\theta) \left(\frac{\partial z_{y}}{\partial \theta_{i}} \right)^{2} = \mathcal{F}_{ii}^{\mathrm{DG}}(\theta).$$

$$\mathbb{E}_{p(\xi)} \left(\frac{\partial \mathfrak{h}^{\mathrm{DG}}}{\partial \theta_{i}} \right)^{4} = \mathbb{E}_{p(\xi)} \left(\sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} \sqrt{\zeta_{y}(x,\theta)} \frac{\partial z_{y}}{\partial \theta_{i}} \xi_{xy} \right)^{4}$$

$$= \sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} \zeta_{y}^{2}(x,\theta) \left(\frac{\partial z_{y}}{\partial \theta_{i}} \right)^{4} + 3 \sum_{(x,y) \neq (x',y')} \zeta_{y}(x,\theta) \left(\frac{\partial z_{y}}{\partial \theta_{i}} \right)^{2} \zeta_{y'}(x',\theta) \left(\frac{\partial z_{y'}}{\partial \theta_{i}} \right)^{2}$$

$$= 3(\mathcal{F}_{ii}^{\mathrm{DG}}(\theta))^{2} - 2 \sum_{x \in \mathcal{D}_{x}} \sum_{y=1}^{C} \zeta_{y}^{2}(x,\theta) \left(\frac{\partial z_{y}}{\partial \theta_{i}} \right)^{4}.$$

Hence,

$$\begin{aligned} \operatorname{Var}(\mathbb{F}_{ii}^{\operatorname{DG}}(\theta)) &= \underset{p(\xi)}{\mathbb{E}} \left(\frac{\partial \mathfrak{h}^{\operatorname{DG}}}{\partial \theta_i} \right)^4 - (\mathcal{F}_{ii}^{\operatorname{DG}}(\theta))^2 \\ &= 2(\mathcal{F}_{ii}^{\operatorname{DG}}(\theta))^2 - 2 \sum_{x \in \mathcal{D}_x} \sum_{y=1}^C \zeta_y^2(x, \theta) \left(\frac{\partial z_y}{\partial \theta_i} \right)^4. \end{aligned}$$

O Proof of Proposition 14

Proof. The proof is similar to Proposition 13 and is also based on the Hutchinson's trick.

$$\mathbb{E}_{p(\xi)} (\mathbb{F}^{LR}(\theta))$$

$$= \mathbb{E}_{p(\xi)} \left(\frac{\partial \mathfrak{h}^{LR}}{\partial \theta} \frac{\partial \mathfrak{h}^{LR}}{\partial \theta^{\top}} \right)$$

$$= \mathbb{E}_{p(\xi)} \left(\sum_{x \in \mathcal{D}_x} \sqrt{\lambda_C(x, \theta)} \left(\frac{\partial z}{\partial \theta} \right)^{\top} v_C(x, \theta) \xi_x \sum_{x' \in \mathcal{D}_x} \sqrt{\lambda_C(x', \theta)} v_C(x', \theta)^{\top} \left(\frac{\partial z}{\partial \theta} \right) \xi_{x'} \right)$$

$$= \sum_{x \in \mathcal{D}_x} \lambda_C(x, \theta) \left(\frac{\partial z}{\partial \theta} \right)^{\top} v_C(x, \theta) v_C(x, \theta)^{\top} \left(\frac{\partial z}{\partial \theta} \right)$$

$$= \mathcal{F}^{LR}(\theta).$$

Therefore

$$\mathbb{E}_{p(\xi)}(\mathbb{F}_{ii}^{LR}(\theta)) = \sum_{x \in \mathcal{D}_x} \lambda_C(x, \theta) \left(\left(\frac{\partial z}{\partial \theta_i} \right)^\top v_C(x, \theta) \right)^2 = \mathcal{F}_{ii}^{LR}(\theta);$$

$$\mathbb{E}_{p(\xi)} \left(\frac{\partial \mathfrak{h}^{LR}}{\partial \theta_i} \right)^4 = \mathbb{E}_{p(\xi)} \left(\sum_{x \in \mathcal{D}_x} \sqrt{\lambda_C(x, \theta)} v_C^\top(x, \theta) \frac{\partial z}{\partial \theta_i} \xi_x \right)^4$$

$$= \sum_{x \in \mathcal{D}_x} \lambda_C^2(x, \theta) \left(v_C^\top(x, \theta) \frac{\partial z}{\partial \theta_i} \right)^4$$

$$+ 3 \sum_{x \neq x'} \lambda_C(x, \theta) \left(v_C^\top(x, \theta) \frac{\partial z}{\partial \theta_i} \right)^2 \lambda_C(x', \theta) \left(v_C^\top(x', \theta) \frac{\partial z}{\partial \theta_i} \right)^2$$

$$= 3(\mathcal{F}_{ii}^{LR}(\theta))^2 - 2 \sum_{x \in \mathcal{D}_x} \lambda_C^2(x, \theta) \left(v_C^\top(x, \theta) \frac{\partial z}{\partial \theta_i} \right)^4.$$

1350 Hence,

$$\begin{aligned} \operatorname{Var}\left(\mathbb{F}_{ii}^{\operatorname{LR}}(\theta)\right) &= \underset{p(\xi)}{\mathbb{E}} \left(\frac{\partial \mathfrak{h}^{\operatorname{LR}}}{\partial \theta_i}\right)^4 - (\mathcal{F}_{ii}^{\operatorname{LR}}(\theta))^2 \\ &= 2(\mathcal{F}_{ii}^{\operatorname{LR}}(\theta))^2 - 2 \sum_{x \in \mathcal{D}_x} \lambda_C^2(x,\theta) \left(v_C^\intercal(x,\theta) \frac{\partial z}{\partial \theta_i}\right)^4. \end{aligned}$$