

Randomized spectral density estimation applied to neural network optimization

Fabio Matti[†] Hei Yin Lam[†] Haoze He[†] Daniel Kressner[†]

[†]ANCHP, École polytechnique fédérale de Lausanne (EPFL), Lausanne, Switzerland

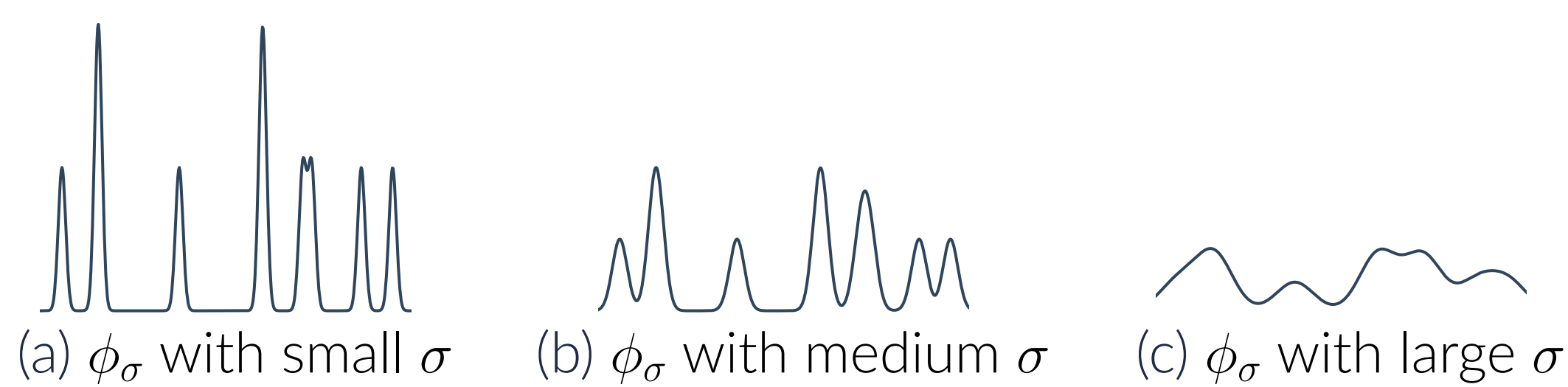
We approximate the **spectral density** of matrices which are only accessed through matrix-vector products. This boils down to estimating the trace of a certain **parameter-dependent** matrix $\mathbf{B}(t)$ for a parameter $t \in [a, b]$. Therefore, we consider three well established **randomized trace estimators** for constant matrices \mathbf{B} when they are straightforwardly applied to parameter-dependent matrices $\mathbf{B}(t)$.

Smoothed spectral density

The **smoothed spectral density** of a symmetric matrix $\mathbf{H} \in \mathbb{R}^{n \times n}$ with eigenvalues $\lambda_1, \dots, \lambda_n \in \mathbb{R}$ is given by the density function

$$\phi_\sigma(t) = \frac{1}{n} \sum_{i=1}^n g_\sigma(t - \lambda_i) = \text{Tr}(g_\sigma(t\mathbf{I}_n - \mathbf{H}))$$

for a Gaussian g_σ of width $\sigma > 0$. This is a **parameter-dependent trace estimation problem** for the positive semi-definite $\mathbf{B}(t) = g_\sigma(t\mathbf{I}_n - \mathbf{H})$.



Fast Hessian-vector products

The **Hessian matrix** $\mathbf{H} \in \mathbb{R}^{n \times n}$ of a neural network parametrized by $\mathbf{w} = (w_1, \dots, w_n)^\top \in \mathbb{R}^n$ and with respect to a loss function \mathcal{L} has entries

$$H_{ij} = \frac{\partial^2 \mathcal{L}}{\partial w_i \partial w_j}.$$

The expansion of the gradient reveals a useful formula for computing **Hessian-vector products** $\mathbf{H}\mathbf{v}$ for any $\mathbf{v} \in \mathbb{R}^n$:

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w} + h\mathbf{v}) = \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) + h\mathbf{H}\mathbf{v} + \mathcal{O}(h^2) \xrightarrow{h \rightarrow 0} \mathbf{H}\mathbf{v} = \left. \frac{\partial}{\partial h} \right|_{h=0} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w} + h\mathbf{v}).$$

Because $f \mapsto \left. \frac{\partial}{\partial h} \right|_{h=0} f(\cdot + h\mathbf{v})$ is a **differential operator**, it can be computed similarly to $\nabla_{\mathbf{w}} \mathcal{L}$ in $\mathcal{O}(n)$ operations with a **backpropagation scheme** [1].

1: Girard-Hutchinson estimator

We define the **Girard-Hutchinson estimator** as

$$\text{Tr}_{\Psi}(\mathbf{B}(t)) = \frac{1}{n_{\Psi}} \sum_{j=1}^{n_{\Psi}} \psi_j^\top \mathbf{B}(t) \psi_j$$

for standard Gaussian $\psi_1, \dots, \psi_{n_{\Psi}}$.

Theorem: Girard-Hutchinson error [2]

If $\mathbf{B}(t) \in \mathbb{R}^{n \times n}$ is symmetric, then

$$\int_a^b |\text{Tr}(\mathbf{B}(t)) - \text{Tr}_{\Psi}(\mathbf{B}(t))| dt \leq \varepsilon \int_a^b \|\mathbf{B}(t)\|_F dt$$

with probability $\geq 1 - \delta$ if $n_{\Psi} = \mathcal{O}(\varepsilon^{-2} \log(\delta^{-1}))$.

Downside: Convergence with n_{Ψ} is slow.

2: Nyström estimator

The **Nyström approximation** is the matrix

$$\widehat{\mathbf{B}}_{\Omega}(t) = (\mathbf{B}(t)\Omega)(\Omega^\top \mathbf{B}(t)\Omega)^\dagger (\mathbf{B}(t)\Omega)^\top.$$

for a standard Gaussian $\Omega \in \mathbb{R}^{n \times n_{\Omega}}$. The error of this approximation depends on the **singular value decay** of $\mathbf{B}(t)$ [2].

We approximate the trace $\text{Tr}(\mathbf{B}(t))$ with

$$\text{Tr}(\widehat{\mathbf{B}}_{\Omega}(t)) = \text{Tr}((\Omega^\top \mathbf{B}(t)\Omega)^\dagger (\Omega \mathbf{B}(t)^2 \Omega^\top))$$

thanks to the cyclic property of the trace.

Downside: Poor approximation for matrices with slow singular value decay.

1 & 2: Nyström++ estimator

The **Nyström++ estimator** computes the Nyström trace estimator and corrects it with the Girard-Hutchinson estimate of the residual

$$\text{Tr}_{\Psi, \Omega}(\mathbf{B}(t)) = \text{Tr}(\widehat{\mathbf{B}}_{\Omega}(t)) + \text{Tr}_{\Psi}(\mathbf{B}(t) - \widehat{\mathbf{B}}_{\Omega}(t))$$

Theorem: Nyström++ error [2]

If $\mathbf{B}(t) \in \mathbb{R}^{n \times n}$ is symmetric positive semi-definite, then

$$\int_a^b |\text{Tr}(\mathbf{B}(t)) - \text{Tr}_{\Psi, \Omega}(\mathbf{B}(t))| dt \leq \varepsilon \int_a^b \text{Tr}(\mathbf{B}(t)) dt$$

with probability $\geq 1 - \delta$ if $n_{\Psi} = n_{\Omega} = \mathcal{O}(\varepsilon^{-1} \log(\delta^{-1}))$.

Chebyshev-Nyström++ algorithm

Expand $g_\sigma(t\mathbf{I}_n - \mathbf{H})$ in the first m **Chebyshev polynomials** T_0, \dots, T_m [3]

$$g_\sigma(t\mathbf{I}_n - \mathbf{H}) \approx \sum_{l=0}^m \mu_l(t) T_l(\mathbf{H}).$$

At each t , the coefficients $\boldsymbol{\mu} = \{\mu_l(t)\}_{l=0}^m$ and the function evaluations $\mathbf{g} = \{g_\sigma(t - \cos(\pi l/m))\}_{l=0}^m$ are related by the **discrete cosine transform (DCT)**

$$\boldsymbol{\mu} = \text{DCT}^{-1}(\mathbf{g}) \iff \mathbf{g} = \text{DCT}(\boldsymbol{\mu}).$$

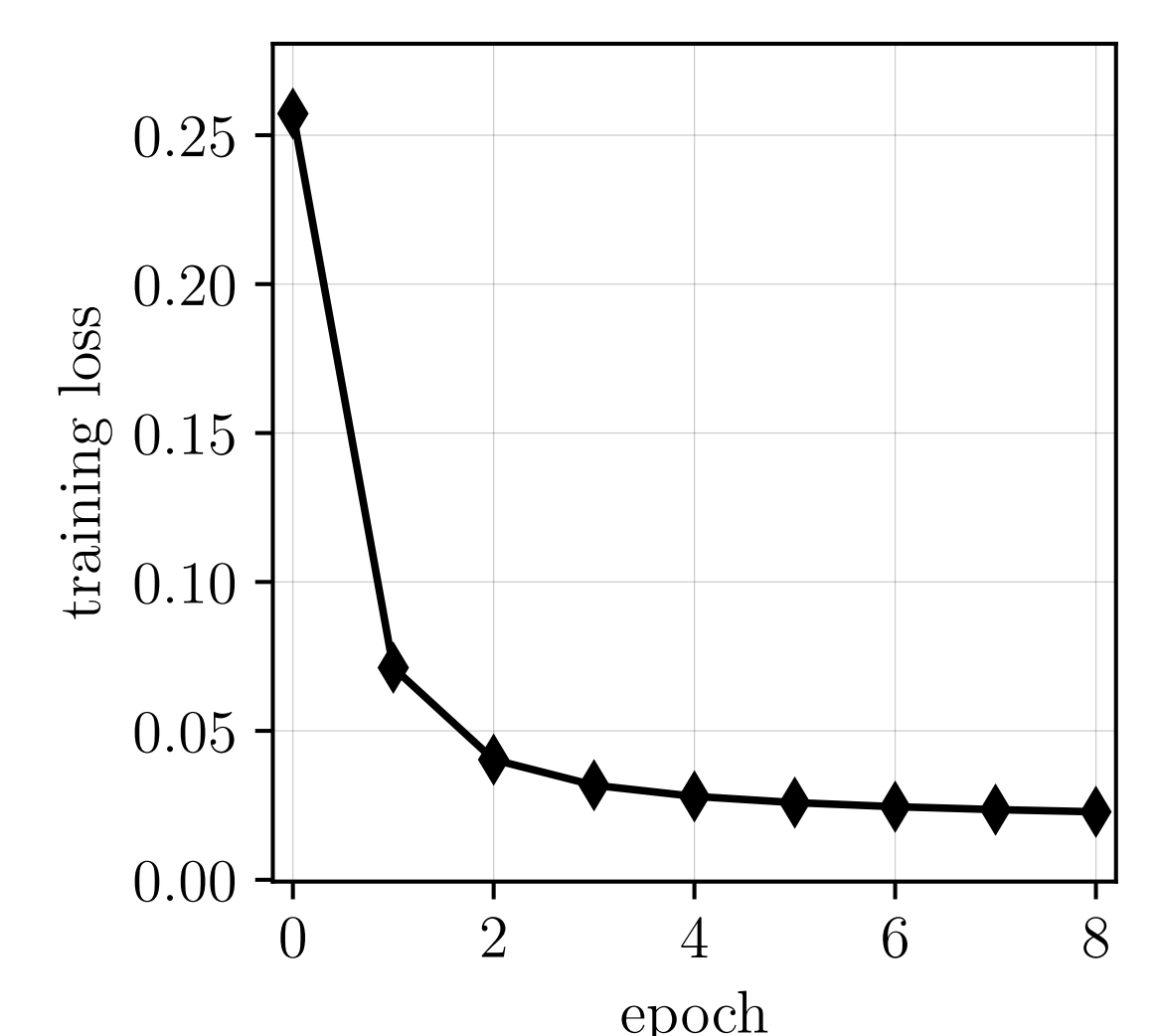
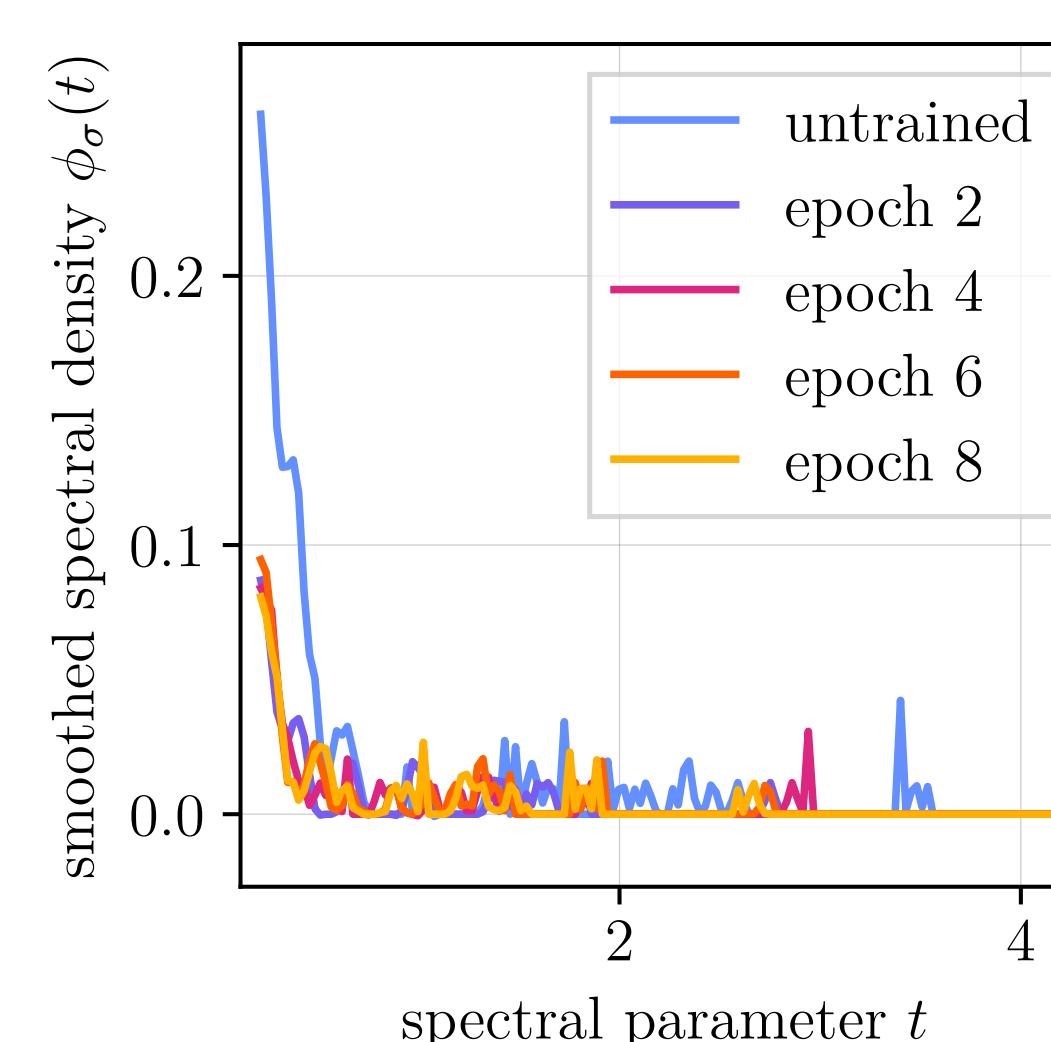
This is used to square Chebyshev expansions in $\mathcal{O}(m \log(m))$ operations:

$$\left(\sum_{l=0}^m \mu_l T_l(\mathbf{H}) \right)^2 = \sum_{l=0}^{2m} \nu_l T_l(\mathbf{H}) \implies \boldsymbol{\nu} = \text{DCT}^{-1}(\text{DCT}(\boldsymbol{\mu})^2).$$

The **Chebyshev-Nyström++ algorithm** applies this observation to the Nyström++ estimator to accurately approximate the spectral density [2].

Application: Sharpness-aware model selection

The smoothed spectral density of the Hessian matrix \mathbf{H} of a neural network is used to monitor its **generalizability**. The presence of large eigenvalues indicates a **sharp minimum**, hence, unfavorable generalization properties.



Conclusion

Despite reusing the same randomization for computing the estimate at each value of the parameter t , the proposed estimators **closely match the corresponding results for constant matrices**.

In practice, the estimators, paired with a **rigorous approach for expanding parameter-dependent matrix functions** in terms of Chebyshev polynomials, can be employed to accurately and efficiently approximate smoothed spectral densities. This can effectively be used to **assess the generalizability of minima** found in a neural network optimization process.

References

- [1] B. A. Pearlmutter, "Fast exact multiplication by the Hessian," *Neural Comput.*, vol. 6, no. 1, pp. 147–160, 1994.
- [2] F. Matti, H. Y. Lam, H. He, and D. Kressner, "Randomized trace estimation for parameter-dependent matrices applied to spectral density approximation," *In preparation*, 2024.
- [3] L. Lin, "Randomized estimation of spectral densities of large matrices made accurate," *Numer. Math.*, vol. 136, no. 1, pp. 183–213, 2017.