

Deep Equilibrium Models are Almost Equivalent to Not-so-deep Explicit Models for High-dimensional Gaussian Mixtures

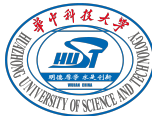
Zenan Ling¹ Longbo Li¹ Zhanbo Feng² Yixuan Zhang³ Fengzhou⁴
Robert C. Qiu¹ Zhengyu Liao^{†1}

¹Huazhong University of Science and Technology ²Shanghai Jiao Tong University

³Hangzhou Dianzi University ⁴Renmin University of China

[†]Correspondence to: zhenyu_liao@hust.edu.cn

2024/07



Deep Equilibrium Models

- **Explicit models:** $z^{(l+1)} = f_{\theta}^{(l)}(z^{(l)}; x)$, for $l = 0, 1, \dots, L - 1$

¹Bai, S. *et al.* *Deep Equilibrium Models*. in *NeurIPS* (2019).

Deep Equilibrium Models

- **Explicit models:** $z^{(l+1)} = f_{\theta}^{(l)}(z^{(l)}; x)$, for $l = 0, 1, \dots, L - 1$
- **Implicit Deep Equilibrium Models (DEQs)¹:** $z^* = f_{\theta}(z^*; x)$

¹Bai, S. *et al.* Deep Equilibrium Models. in *NeurIPS* (2019).

Deep Equilibrium Models

- **Explicit models:** $z^{(l+1)} = f_{\theta}^{(l)}(z^{(l)}; x)$, for $l = 0, 1, \dots, L - 1$
- **Implicit Deep Equilibrium Models (DEQs)¹:** $z^* = f_{\theta}(z^*; x)$
 - a typical “single-layer” implicit model
⇒ an infinite-depth weight-tied model with an input injection

¹Bai, S. *et al.* *Deep Equilibrium Models*. in *NeurIPS* (2019).

Deep Equilibrium Models

- **Explicit models:** $z^{(l+1)} = f_{\theta}^{(l)}(z^{(l)}; x)$, for $l = 0, 1, \dots, L - 1$
- **Implicit Deep Equilibrium Models (DEQs)¹:** $z^* = f_{\theta}(z^*; x)$
 - a typical “single-layer” implicit model
⇒ an infinite-depth weight-tied model with an input injection
 - Memory efficient: DEQs solve an equilibrium point directly and compute gradients with *implicit differentiation*.

¹Bai, S. *et al.* Deep Equilibrium Models. in *NeurIPS* (2019).

Deep Equilibrium Models

- **Explicit models:** $z^{(l+1)} = f_{\theta}^{(l)}(z^{(l)}; x)$, for $l = 0, 1, \dots, L - 1$
- **Implicit Deep Equilibrium Models (DEQs)¹:** $z^* = f_{\theta}(z^*; x)$
 - a typical “single-layer” implicit model
⇒ an infinite-depth weight-tied model with an input injection
 - Memory efficient: DEQs solve an equilibrium point directly and compute gradients with *implicit differentiation*.
 - Universality of DEQs: any deep explicit NN can be reformulated as a “single-layer” DEQ.

¹Bai, S. *et al.* Deep Equilibrium Models. in *NeurIPS* (2019).

Deep Equilibrium Models

- **Explicit models:** $z^{(l+1)} = f_{\theta}^{(l)}(z^{(l)}; x)$, for $l = 0, 1, \dots, L - 1$
- **Implicit Deep Equilibrium Models (DEQs)¹:** $z^* = f_{\theta}(z^*; x)$
 - a typical “single-layer” implicit model
⇒ an infinite-depth weight-tied model with an input injection
 - Memory efficient: DEQs solve an equilibrium point directly and compute gradients with *implicit differentiation*.
 - Universality of DEQs: any deep explicit NN can be reformulated as a “single-layer” DEQ.
 - Remarkable success in various tasks, e.g., NPL, CV [2].

¹Bai, S. *et al.* Deep Equilibrium Models. in *NeurIPS* (2019).

Deep Equilibrium Models

- **Explicit models:** $z^{(l+1)} = f_{\theta}^{(l)}(z^{(l)}; x)$, for $l = 0, 1, \dots, L - 1$
- **Implicit Deep Equilibrium Models (DEQs)¹:** $z^* = f_{\theta}(z^*; x)$
 - a typical “single-layer” implicit model
⇒ an infinite-depth weight-tied model with an input injection
 - Memory efficient: DEQs solve an equilibrium point directly and compute gradients with *implicit differentiation*.
 - Universality of DEQs: any deep explicit NN can be reformulated as a “single-layer” DEQ.
 - Remarkable success in various tasks, e.g., NPL, CV [2].
 - Significant computational overhead: a consequence of root-finding.

¹Bai, S. *et al.* Deep Equilibrium Models. in *NeurIPS* (2019).

Deep Equilibrium Models

- **Explicit models:** $z^{(l+1)} = f_{\theta}^{(l)}(z^{(l)}; x)$, for $l = 0, 1, \dots, L - 1$
- **Implicit Deep Equilibrium Models (DEQs)¹:** $z^* = f_{\theta}(z^*; x)$
 - a typical “single-layer” implicit model
⇒ an infinite-depth weight-tied model with an input injection
 - Memory efficient: DEQs solve an equilibrium point directly and compute gradients with *implicit differentiation*.
 - Universality of DEQs: any deep explicit NN can be reformulated as a “single-layer” DEQ.
 - Remarkable success in various tasks, e.g., NPL, CV [2].
 - Significant computational overhead: a consequence of root-finding.
- **Limited theoretical understanding of DEQs:** the connections and differences between implicit DEQs and explicit models

¹Bai, S. et al. Deep Equilibrium Models. in *NeurIPS* (2019).

Deep Equilibrium Models

- **Explicit models:** $z^{(l+1)} = f_{\theta}^{(l)}(z^{(l)}; x)$, for $l = 0, 1, \dots, L - 1$
- **Implicit Deep Equilibrium Models (DEQs)¹:** $z^* = f_{\theta}(z^*; x)$
 - a typical “single-layer” implicit model
⇒ an infinite-depth weight-tied model with an input injection
 - Memory efficient: DEQs solve an equilibrium point directly and compute gradients with *implicit differentiation*.
 - Universality of DEQs: any deep explicit NN can be reformulated as a “single-layer” DEQ.
 - Remarkable success in various tasks, e.g., NPL, CV [2].
 - Significant computational overhead: a consequence of root-finding.
- **Limited theoretical understanding of DEQs:** the connections and differences between implicit DEQs and explicit models
 - whether general DEQs have advantages over explicit networks, or
 - whether an equivalent explicit NN exists for a given implicit DEQ.

¹Bai, S. et al. Deep Equilibrium Models. in *NeurIPS* (2019).

DEQ Models

Vanilla DEQs

Let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$ denote the input data, consider a vanilla DEQ with output $f(\mathbf{x}_i)$ given by

$$f(\mathbf{x}_i) = \mathbf{a}^\top \mathbf{z}_i^*, \quad (1)$$

where $\mathbf{a} \in \mathbb{R}^m$ and $\mathbf{z}_i^{(*)} \triangleq \lim_{l \rightarrow \infty} \mathbf{z}_i^{(l)} \in \mathbb{R}^m$ with

$$\mathbf{z}_i^{(l)} = \frac{1}{\sqrt{m}} \phi \left(\sigma_a \mathbf{A} \mathbf{z}_i^{(l-1)} + \sigma_b \mathbf{B} \mathbf{x}_i \right) \in \mathbb{R}^m, \text{ for } l \geq 1, \quad (2)$$

for some appropriate initialization $\mathbf{z}_i^{(0)}$. Here, $\mathbf{A} \in \mathbb{R}^{m \times m}$ and $\mathbf{B} \in \mathbb{R}^{m \times p}$ are the DEQ weight parameters, $\sigma_a, \sigma_b \in \mathbb{R}$ are constants, and ϕ is an element-wise activation. Note that \mathbf{z}_i^* can also be determined as the equilibrium point of

$$\mathbf{z}_i^* = \frac{1}{\sqrt{m}} \phi \left(\sigma_a \mathbf{A} \mathbf{z}_i^* + \sigma_b \mathbf{B} \mathbf{x}_i \right). \quad (3)$$

Weight and Activation

Assumption

- (1) **Initialization:** $\mathbf{a} \in \mathbb{R}^m$, $\mathbf{A} \in \mathbb{R}^{m \times m}$ and $\mathbf{B} \in \mathbb{R}^{m \times p}$ are initialized with i.i.d. entries of zero mean, unit variance, and finite fourth-order moment;
- (2) **Weak differentiability:** ϕ is centered and L_1 -Lipschitz, and $\max_{k \in \{0,1,2,3,4\}} |\mathbb{E}[\phi^{(k)}(\xi)]| < \infty$;
- (3) **Variance parameter:** $\sigma_a^2 < 1/(4L_1^2)$ and $\sigma_a^2 < 2/(\mathbb{E}[(\phi^2(\tau\xi))''])$ for $\tau > 0$ and $\xi \sim \mathcal{N}(0, 1)$.

Remark

- (i) hold for commonly-used initialization;
- (ii) hold for commonly-used smooth, e.g., Tanh, and piecewise linear activations, e.g., ReLU and Leaky ReLU;
- (iii) guarantee the existence and uniqueness of the equilibrium point.

GMM Data

Assumption (GMM Data)

For $\mathbf{x}_i \in \mathcal{C}_a$, $\sqrt{p}\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}_a, \mathbf{C}_a)$, $a \in [K]$. For n, p both large that (i) $p = \Theta(n)$ and $n_a = \Theta(n)$; (ii) $\|\boldsymbol{\mu}_a\| = \mathcal{O}(1)$; (iii) for $\mathbf{C}^\circ \equiv \sum_{a=1}^K \frac{n_a}{n} \mathbf{C}_a$ and $\mathbf{C}_a^\circ \equiv \mathbf{C}_a - \mathbf{C}^\circ$, we have $\|\mathbf{C}_a\| = \mathcal{O}(1)$, $\text{tr} \mathbf{C}_a^\circ = \mathcal{O}(\sqrt{p})$ and $\text{tr}(\mathbf{C}_a \mathbf{C}_b) = \mathcal{O}(p)$.

Remark

- (i) Assumptions of GMM are non-trivial, widely studied in LDA, SVM, and NNs;
- (ii) GMM is a universal approximator, can approximate any distribution to an arbitrary error;
- (iii) For large p and n , data generated from generative models, e.g., GANs, behaves as GMM [4].

High-dimensional Statistics

- **Label information:** $\mathbf{J} \equiv [j_1, \dots, j_K] \in \mathbb{R}^{n \times K}$, $[j_a]_i = 1_{\mathbf{x}_i \in \mathcal{C}_a}$;
- **Second-order data fluctuation vector:** $\boldsymbol{\psi} \equiv \{\|\mathbf{x}_i - \mathbb{E}[\mathbf{x}_i]\|^2 - \mathbb{E}[\|\mathbf{x}_i - \mathbb{E}[\mathbf{x}_i]\|^2]\}_{i=1}^n \in \mathbb{R}^n$;
- **Second-order GMM statistics:** $\mathbf{T} = \{\text{tr} \mathbf{C}_a \mathbf{C}_b / p\}_{a,b=1}^K \in \mathbb{R}^{K \times K}$, $\mathbf{t} = \{\text{tr} \mathbf{C}_a^\circ / \sqrt{p}\} \in \mathbb{R}^K$;
- **A fixed point:** $\tau_0 \equiv \sqrt{\text{tr} \mathbf{C}^\circ / p}$, and τ_* be the fixed point to the equation

$$\tau_* = \sqrt{\sigma_a^2 \mathbb{E}[\phi^2(\tau_* \xi)] + \sigma_b^2 \tau_0^2}, \quad \xi \sim \mathcal{N}(0, 1).$$

Implicit CK and NTK

Conjugate Kernels (CKs) and Neural Tangent Kernels (NTKs): an analytical assessment of the convergence and generalization properties of wide NNs.

Implicit CK and NTK

Conjugate Kernels (CKs) and Neural Tangent Kernels (NTKs): an analytical assessment of the convergence and generalization properties of wide NNs.

Implicit-CK [3]

$G^* = \lim_{l \rightarrow \infty} G^{(l)}$, where $G_{ij}^{(l)} = \mathbb{E}_{(u_l, v_l)} [\phi(u_l) \phi(v_l)]$ with $(u_l, v_l) \sim \mathcal{N}\left(0, \begin{bmatrix} \Lambda_{ii}^{(l)} & \Lambda_{ij}^{(l)} \\ \Lambda_{ji}^{(l)} & \Lambda_{jj}^{(l)} \end{bmatrix}\right)$ and $\Lambda_{ij}^{(l)} = \sigma_a^2 G_{ij}^{(l-1)} + \sigma_b^2 \mathbf{x}_i^\top \mathbf{x}_j$, for $l \geq 1$, and $G_{ij}^{(0)} = (\mathbf{z}_i^{(0)})^\top \mathbf{z}_j^{(0)}$.

Implicit CK and NTK

Conjugate Kernels (CKs) and Neural Tangent Kernels (NTKs): an analytical assessment of the convergence and generalization properties of wide NNs.

Implicit-CK [3]

$$\mathbf{G}^* = \lim_{l \rightarrow \infty} \mathbf{G}^{(l)}, \text{ where } \mathbf{G}_{ij}^{(l)} = \mathbb{E}_{(\mathbf{u}_l, \mathbf{v}_l)} [\phi(\mathbf{u}_l) \phi(\mathbf{v}_l)] \text{ with } (\mathbf{u}_l, \mathbf{v}_l) \sim \mathcal{N} \left(0, \begin{bmatrix} \Lambda_{ii}^{(l)} & \Lambda_{ij}^{(l)} \\ \Lambda_{ji}^{(l)} & \Lambda_{jj}^{(l)} \end{bmatrix} \right) \text{ and } \Lambda_{ij}^{(l)} = \sigma_a^2 \mathbf{G}_{ij}^{(l-1)} + \sigma_b^2 \mathbf{x}_i^\top \mathbf{x}_j, \text{ for } l \geq 1, \text{ and } \mathbf{G}_{ij}^{(0)} = (\mathbf{z}_i^{(0)})^\top \mathbf{z}_j^{(0)}.$$

Implicit-NTK

$$\mathbf{K}^* = \lim_{l \rightarrow \infty} \mathbf{K}^{(l)}, \text{ where } \mathbf{K}_{ij}^{(l)} = \sum_{h=1}^{l+1} \left(\mathbf{G}_{ij}^{(h-1)} \prod_{h'=h}^{l+1} \dot{\mathbf{G}}_{ij}^{(h')} \right), \text{ with } \dot{\mathbf{G}}_{ij}^{(l)} = \sigma_a^2 \mathbb{E}_{(\mathbf{u}_l, \mathbf{v}_l)} [\phi'(\mathbf{u}_l) \phi'(\mathbf{v}_l)], \text{ so that } \mathbf{K}_{ij}^* \equiv \mathbf{G}_{ij}^* / (1 - \dot{\mathbf{G}}_{ij}^*).$$

High-dimensional Equivalents

Theorem (High-dimensional approximation of Implicit-CKs [5])

Let Assumptions 1 and 2 hold, and let the activation ϕ be centered s.t., $\mathbb{E}[\phi(\tau_*\xi)] = 0$ for $\xi \sim \mathcal{N}(0, 1)$. It holds that $\|\mathbf{G}^* - \overline{\mathbf{G}}\| = \mathcal{O}(n^{-1/2})$ where

$$\overline{\mathbf{G}} \equiv \alpha_{*,1} \mathbf{X}^\top \mathbf{X} + \mathbf{V} \mathbf{C}_* \mathbf{V}^\top + (\gamma_*^2 - \tau_0^2 \alpha_{*,1}) \mathbf{I}_n,$$

with $\mathbf{V} = [\mathbf{J}/\sqrt{p}, \boldsymbol{\psi}]$ and $\mathbf{C}_* = \begin{bmatrix} \alpha_{*,2} \mathbf{t} \mathbf{t}^\top + \alpha_{*,3} \mathbf{T} & \alpha_{*,2} \mathbf{t} \\ \alpha_{*,2} \mathbf{t}^\top & \alpha_{*,2} \end{bmatrix}$. Non-negative scalars $\gamma_*, \alpha_{*,1}, \alpha_{*,2}, \alpha_{*,3} \geq 0$ are defined, for $\xi \sim \mathcal{N}(0, 1)$, as

$$\begin{aligned} \gamma_* &= \sqrt{\mathbb{E}[\phi^2(\tau_*\xi)]}, \quad \alpha_{*,1} = \frac{\sigma_b^2 \mathbb{E}[\phi'(\tau_*\xi)]^2}{1 - \sigma_a^2 \mathbb{E}[\phi'(\tau_*\xi)]^2}, \\ \alpha_{*,2} &= \frac{\mathbb{E}[\phi''(\tau_*\xi)]^2}{4(1 - \sigma_a^2 \mathbb{E}[\phi'(\tau_*\xi)]^2)} \alpha_{*,4}, \quad \alpha_{*,3} = \frac{\mathbb{E}[\phi''(\tau_*\xi)]^2 (\sigma_a^2 \alpha_{*,1} + \sigma_b^2)^2}{2(1 - \sigma_a^2 \mathbb{E}[\phi'(\tau_*\xi)]^2)} \end{aligned} \quad (4)$$

with $\alpha_{*,4} = (1 - \frac{\sigma_a^2}{2} \mathbb{E}[(\phi^2(\tau_*\xi))''])^{-1} \sigma_b^2$.

High-dimensional Equivalents

Theorem (High-dimensional approximation of Implicit-NTKs)

Under the same settings and notations of Theorem 1, it holds that $\|\mathbf{K}^* - \overline{\mathbf{K}}\| = \mathcal{O}(n^{-1/2})$ where

$$\overline{\mathbf{K}} \equiv \beta_{*,1} \mathbf{X}^\top \mathbf{X} + \mathbf{V} \mathbf{D}_* \mathbf{V}^\top + (\kappa_*^2 - \tau_0^2 \beta_{*,1}) \mathbf{I}_n,$$

with \mathbf{V} defined in Theorem 1, $\mathbf{D}_* = \begin{bmatrix} \beta_{*,2} \mathbf{t} \mathbf{t}^\top + \beta_{*,3} \mathbf{T} & \beta_{*,2} \mathbf{t} \\ \beta_{*,2} \mathbf{t}^\top & \beta_{*,2} \end{bmatrix}$, as well as non-negative scalars $\kappa_*, \beta_{*,1}, \beta_{*,2}, \beta_{*,3} \geq 0$ defined as

$$\begin{aligned} \kappa_* &= \frac{\tau_*}{\sqrt{1 - \sigma_a^2 \mathbb{E}[\phi'(\tau_* \xi)^2]}}, \quad \beta_{*,1} = \frac{\alpha_{*,1}}{1 - \sigma_a^2 \mathbb{E}[\phi'(\tau_* \xi)]^2}, \\ \beta_{*,2} &= \frac{\alpha_{*,2}}{1 - \sigma_a^2 \mathbb{E}[\phi'(\tau_* \xi)]^2}, \quad \beta_{*,3} = \frac{\alpha_{*,3} + \beta_{*,1}(\sigma_a^2 \mathbb{E}[\phi''(\tau_* \xi)]^2 + \sigma_b^2) \alpha_{*,1}}{1 - \sigma_a^2 \mathbb{E}[\phi'(\tau_* \xi)]^2}, \end{aligned}$$

for $\xi \sim \mathcal{N}(0, 1)$.

High-dimensional Equivalents

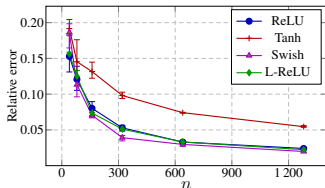


Figure: Evolution of relative spectral norm error $\|\mathbf{G}^* - \overline{\mathbf{G}}\| / \|\mathbf{G}^*\|$ w.r.t. sample size n , for DEQs with different activations and $\sigma_a^2 = 0.2$, on two-class GMM, $p/n = 0.8$, $\mu_a = [\mathbf{0}_{8(a-1)}; 8; \mathbf{0}_{p-8a+7}]$, and $\mathbf{C}_a = (1 + 8(a-1)/\sqrt{p})\mathbf{I}_p$, $a \in \{1, 2\}$.

Remark

As, the “equivalent” Implicit-CK and NTK matrices, $\overline{\mathbf{G}}$ and $\overline{\mathbf{K}}$,

- (1) depend on the input GMM data (\mathbf{X}), their class structure (\mathbf{J}) and higher-order statistics (\mathbf{t} and \mathbf{T}), *explicitly*; and
- (2) are *independent* of the distribution of the weight matrices \mathbf{A} and \mathbf{B} ; and
- (3) depend on σ_a^2 , σ_b^2 , and ϕ *only* via four scalars $\alpha_{*,1-3}$, γ_* , and $\beta_{*,1-3}$, κ_* , *explicitly*.

High-dimensional Equivalents

Theorem (High-dimensional approximation of Explicit-CKs [6])

Consider a fully-connected NN model $\mathbf{x}_i^{(l)} = \frac{1}{\sqrt{m_l}} \sigma_l(\mathbf{W}_l \mathbf{x}_i^{(l-1)})$, for $l = 1, \dots, L$, with GMM input. For the corresponding Explicit-CK matrix $\Sigma^{(l)}$, it holds that $\|\Sigma^{(l)} - \bar{\Sigma}^{(l)}\| = \mathcal{O}(n^{-1/2})$ where

$$\bar{\Sigma}^{(l)} = \tilde{\alpha}_{l,1} \mathbf{X}^\top \mathbf{X} + \mathbf{V} \tilde{\mathbf{C}}_l \mathbf{V}^\top + (\tilde{\tau}_l^2 - \tau_0^2 \tilde{\alpha}_{l,1}) \mathbf{I}_n,$$

with \mathbf{V} defined in Theorem 1, $\tilde{\mathbf{C}}_l = \begin{bmatrix} \tilde{\alpha}_{l,2} \mathbf{t} \mathbf{t}^\top + \tilde{\alpha}_{l,3} \mathbf{T} & \tilde{\alpha}_{l,2} \mathbf{t} \\ \tilde{\alpha}_{l,2} \mathbf{t}^\top & \tilde{\alpha}_{l,2} \end{bmatrix}$. Non-negative scalars $\tilde{\alpha}_{l,1}, \tilde{\alpha}_{l,2}, \tilde{\alpha}_{l,3}$ defined recursively as $\tilde{\alpha}_{0,1} = \tilde{\alpha}_{0,4} = 1$, $\tilde{\alpha}_{0,2} = \tilde{\alpha}_{0,3} = 0$, and

$$\begin{aligned} \tilde{\alpha}_{l,1} &= \mathbb{E}[\sigma'_l(\tilde{\tau}_{l-1} \xi)]^2 \tilde{\alpha}_{l-1,1}, \quad \tilde{\alpha}_{l,2} = \mathbb{E}[\sigma'_l(\tilde{\tau}_{l-1} \xi)]^2 \tilde{\alpha}_{l-1,2} + \frac{1}{4} \mathbb{E}[\sigma''_l(\tilde{\tau}_{l-1} \xi)]^2 \tilde{\alpha}_{l-1,4}, \\ \tilde{\alpha}_{l,3} &= \mathbb{E}[\sigma'_l(\tilde{\tau}_{l-1} \xi)]^2 \tilde{\alpha}_{l-1,3} + \frac{1}{2} \mathbb{E}[\sigma''_l(\tilde{\tau}_{l-1} \xi)]^2 \tilde{\alpha}_{l-1,1}, \end{aligned}$$

with $\tilde{\alpha}_{l,4} = \mathbb{E}[(\sigma_l^2(\tilde{\tau}_{l-1} \xi))''] \tilde{\alpha}_{l-1,4}$, for $\xi \sim \mathcal{N}(0, 1)$.

High-dimensional Equivalence between Implicit and Explicit NNs

Observation

The high-dimensional approximation \overline{G} of the Implicit-CK takes a consistent form with that $(\overline{\Sigma}^{(l)})$ of the Explicit-CK.

High-dimensional Equivalence between Implicit and Explicit NNs

Observation

The high-dimensional approximation \overline{G} of the Implicit-CK takes a consistent form with that $(\overline{\Sigma}^{(l)})$ of the Explicit-CK.

Key idea

Given a DEQ, design activations of an L -layer explicit NN *s.t.* its Explicit-CK $\Sigma^{(L)}$ shares the same coefficients as the Implicit-CK G^* , i.e.,
 $\tilde{\tau}_L = \gamma_*$, $\tilde{\alpha}_{L,i} = \alpha_{*,i}$, $i \in \{1, 2, 3\}$.

High-dimensional Equivalence between Implicit and Explicit NNs

Observation

The high-dimensional approximation \overline{G} of the Implicit-CK takes a consistent form with that $(\overline{\Sigma}^{(l)})$ of the Explicit-CK.

Key idea

Given a DEQ, design activations of an L -layer explicit NN s.t. its Explicit-CK $\Sigma^{(L)}$ shares the same coefficients as the Implicit-CK G^* , i.e.,
 $\tilde{\tau}_L = \gamma_*$, $\tilde{\alpha}_{L,i} = \alpha_{*,i}$, $i \in \{1, 2, 3\}$.

Implicit- versus Explicit-CK

It follows from Theorem 3 that, for the single-hidden-layer ENN, one *must* have $\tilde{\alpha}_{1,2} = \frac{1}{2}\tilde{\alpha}_{1,3}$. On the contrast, $\alpha_{*,2} = \frac{1}{2}\alpha_{*,3}$ does *not* necessarily hold for *all* DEQs. As such, for a given DEQ,

- if $\alpha_{*,2} = \frac{1}{2}\alpha_{*,3}$, a *single-hidden-layer* ENN suffices to match the given DEQ;
- if $\alpha_{*,2} \neq \frac{1}{2}\alpha_{*,3}$, an ENN with (*at least*) *two hidden layers* is required.

High-dimensional Equivalence between Implicit and Explicit NNs

Observation

The high-dimensional approximation $\bar{\mathbf{G}}$ of the Implicit-CK takes a consistent form with that $(\bar{\Sigma}^{(l)})$ of the Explicit-CK.

Key idea

Given a DEQ, design activations of an L -layer explicit NN s.t. its Explicit-CK $\Sigma^{(L)}$ shares the same coefficients as the Implicit-CK \mathbf{G}^* , i.e.,
 $\tilde{\tau}_L = \gamma_*$, $\tilde{\alpha}_{L,i} = \alpha_{*,i}$, $i \in \{1, 2, 3\}$.

Implicit- versus Explicit-CK

It follows from Theorem 3 that, for the single-hidden-layer ENN, one *must* have $\tilde{\alpha}_{1,2} = \frac{1}{2}\tilde{\alpha}_{1,3}$. On the contrast, $\alpha_{*,2} = \frac{1}{2}\alpha_{*,3}$ does *not* necessarily hold for *all* DEQs. As such, for a given DEQ,

- if $\alpha_{*,2} = \frac{1}{2}\alpha_{*,3}$, a *single-hidden-layer* ENN suffices to match the given DEQ;
- if $\alpha_{*,2} \neq \frac{1}{2}\alpha_{*,3}$, an ENN with (*at least*) *two hidden layers* is required.

Remark. Results for NTKs can be similarly obtained with our Theorem 2

Designing Equivalent Explicit NNs via CK matching

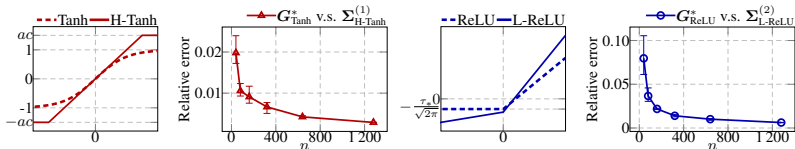


Figure: CKs of implicit DEQs and explicit NNs are close, if the activations of ENNs are constructed according to our Examples.

Example

- For a given Tanh-DEQ, use a *single-hidden-layer* H-Tanh-ENN with

$$\sigma_{\text{H-Tanh}}(x) \equiv ax \cdot 1_{-c \leq x \leq c} + ac \cdot (1_{x \geq c} - 1_{x \leq -c}),$$

$$\text{s.t. } \|G_{\text{Tanh}}^* - \Sigma_{\text{H-Tanh}}^{(1)}\| = \mathcal{O}(n^{-1/2}).$$

- For a given ReLU-DEQ, use a *two-hidden-layer* L-ReLU-ENN with

$$\sigma_{\text{L-ReLU}}^{(l)}(x) \equiv \max(a_l x, b_l x) - \frac{a_l - b_l}{\sqrt{2\pi}} \tilde{\tau}_l, \quad l = 1, 2,$$

$$\text{s.t. } \|G_{\text{ReLU}}^* - \Sigma_{\text{L-ReLU}}^{(2)}\| = \mathcal{O}(n^{-1/2}).$$

Experiments

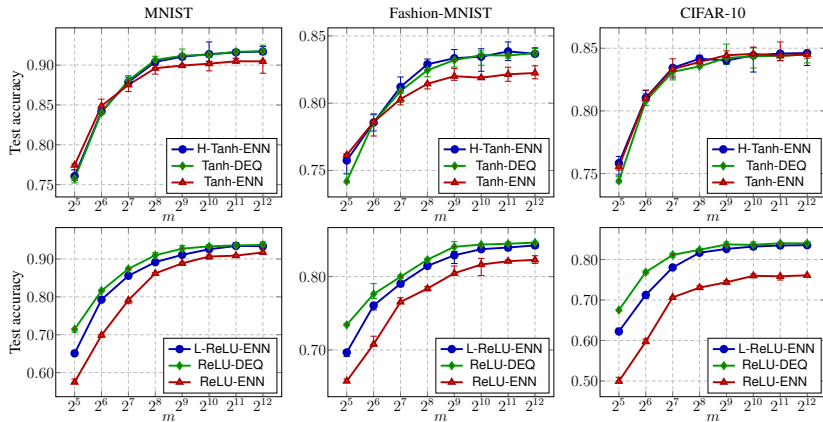


Figure: Classification accuracies of implicit DEQs and explicit models trained with SGD.

Take-away

Take away message:

- for GMM input data, **random matrix theory (RMT)** allows for precise characterization of (the CKs and NTKs of) random explicit and implicit NNs;

Take-away

Take away message:

- for GMM input data, **random matrix theory (RMT)** allows for precise characterization of (the CKs and NTKs of) random explicit and implicit NNs;
- **explicit** connections between implicit and explicit NNs: high-dimensional “**equivalence**”, making implicit NNs explicit and significantly reducing the computational overhead;

Take-away

Take away message:

- for GMM input data, **random matrix theory (RMT)** allows for precise characterization of (the CKs and NTKs of) random explicit and implicit NNs;
- **explicit** connections between implicit and explicit NNs: high-dimensional “**equivalence**”, making implicit NNs explicit and significantly reducing the computational overhead;
- future work: beyond the “lazy” CK/NTK regime \Rightarrow feature learning;

Take-away

Take away message:

- for GMM input data, **random matrix theory (RMT)** allows for precise characterization of (the CKs and NTKs of) random explicit and implicit NNs;
- **explicit** connections between implicit and explicit NNs: high-dimensional “**equivalence**”, making implicit NNs explicit and significantly reducing the computational overhead;
- future work: beyond the “lazy” CK/NTK regime \Rightarrow feature learning;
- future work: extends to another typical implicit models, Neural ODEs, e.g., diffusion models, providing theoretical understanding and accelerating the sampling process.

Reference and Acknowledgement

Reference:

1. Bai, S., Kolter, J. Z. & Koltun, V. *Deep Equilibrium Models*. in *NeurIPS* (2019).
2. Xie, X., Wang, Q., Ling, Z., Liu, G. & Lin, Z. Optimization Induced Equilibrium Networks: An Explicit Optimization Perspective for Understanding Equilibrium Models. *TPAMI* (2022).
3. Ling, Z., Xie, X., Wang, Q., Zhang, Z. & Lin, Z. *Global convergence of over-parameterized deep equilibrium models*. in *AISTATS* (2023).
4. Seddik, M. E. A., Louart, C., Tamaazousti, M. & Couillet, R. *Random matrix theory proves that deep learning representations of gan-data behave as gaussian mixtures*. in *ICML* (2020).
5. Ling, Z. et al. *Deep Equilibrium Models are Almost Equivalent to Not-so-deep Explicit Models for High-dimensional Gaussian Mixtures*. in *ICML* (2024).
6. Gu, L., Du, Y., Yuan, Z., Qiu, R. & Liao, Z. "Lossless" Compression of Deep Neural Networks: A High-dimensional Neural Tangent Kernel Approach. in *NeurIPS* (2022).

Thank you! Q & A?