

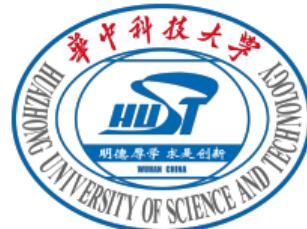
Random Matrix Theory and Its Applications in ML

@ Fudan University

Zhenyu Liao

School of Electronic Information and Communications
Huazhong University of Science and Technology

May 25, 2024



Outline

1 Introduction and Motivation

- Sample covariance matrix
- Curse of dimensionality in high-dimensional classification

2 An Introduction Deep Learning for Statisticians/Mathematicians

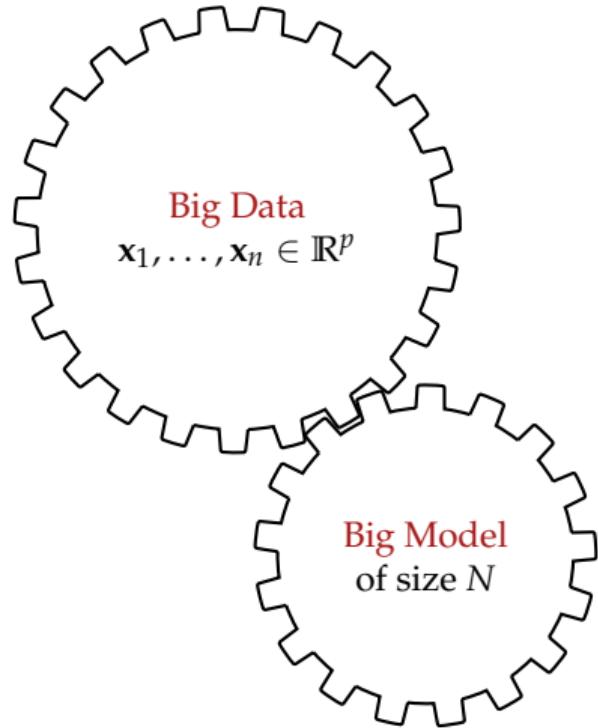
3 Important Theoretical Questions for DL

4 Random (and Not-so Random) Matrix Theory in DL

- Shallow and deep NN with random weights
- NN with nonrandom weights

5 Conclusion

Motivation: understanding large-dimensional machine learning



- ▶ **Big Data era:** exploit large n, p, N
- ▶ **counterintuitive** phenomena **different** from classical asymptotics statistics
- ▶ complete **change** of understanding of many methods in statistics and machine learning
- ▶ **Random Matrix Theory (RMT)** provides the tools!

Sample covariance matrix in the large n, p regime

- ▶ **Problem:** estimate covariance $\mathbf{C} \in \mathbb{R}^{p \times p}$ from n data samples $\mathbf{x}_1, \dots, \mathbf{x}_n$ with $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$,
- ▶ Maximum likelihood sample covariance matrix with **entry-wise** convergence

$$\hat{\mathbf{C}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \in \mathbb{R}^{p \times p}, \quad [\hat{\mathbf{C}}]_{ij} \rightarrow [\mathbf{C}]_{ij}$$

almost surely as $n \rightarrow \infty$: optimal for $n \gg p$ (or, for p “small”).

- ▶ In the regime $n \sim p$, conventional wisdom breaks down: for $\mathbf{C} = \mathbf{I}_p$ with $n < p$, $\hat{\mathbf{C}}$ has at least $p - n$ **zero eigenvalues**:

$$\boxed{\|\hat{\mathbf{C}} - \mathbf{C}\| \not\rightarrow 0, \quad n, p \rightarrow \infty} \Rightarrow \text{eigenvalue mismatch and not consistent!}$$

- ▶ due to $\|\mathbf{A}\|_\infty \leq \|\mathbf{A}\| \leq p\|\mathbf{A}\|_\infty$ for $\mathbf{A} \in \mathbb{R}^{p \times p}$ and $\|\mathbf{A}\|_\infty \equiv \max_{ij} |\mathbf{A}_{ij}|$.

When is one in the random matrix regime? Almost always!

What about $n = 100p$? For $\mathbf{C} = \mathbf{I}_p$, as $n, p \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$: MP law

$$\mu(dx) = (1 - c^{-1})^+ \delta(x) + \frac{1}{2\pi cx} \sqrt{(x - E_-)^+ (E_+ - x)^+} dx$$

where $E_- = (1 - \sqrt{c})^2$, $E_+ = (1 + \sqrt{c})^2$ and $(x)^+ \equiv \max(x, 0)$. **Close match!**

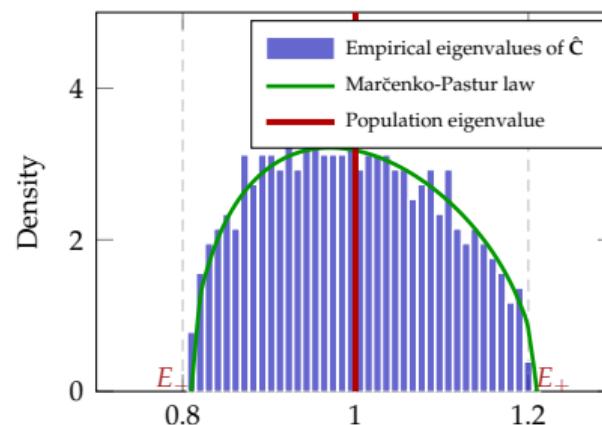


Figure: Eigenvalue distribution of $\hat{\mathbf{C}}$ versus Marčenko-Pastur law, $p = 500, n = 50\,000$.

- ▶ eigenvalues span on $[E_- = (1 - \sqrt{c})^2, E_+ = (1 + \sqrt{c})^2]$.
- ▶ for $n = 100p$, on a range of $\pm 2\sqrt{c} = \pm 0.2$ around the population eigenvalue 1.

Classical large- n asymptotic analysis mostly fails today

- ▶ large- n intuition, and many existing popular methods in statistics, biology, finance, signal processing, telecommunication, and machine learning (ML), must **fail** even with $n = 100p$!
- ▶ **RMT** as a flexible and powerful tool to **understand** and **recreate** these methods
- ▶ in essence, “**increasing** complexity of the system models employed in above fields demand **low** complexity analysis”
- ▶ as motivating examples, how RMT can be applied to assess **machine learning** method such as principle component analysis (PCA), and kernel spectral clustering

“Curse of dimensionality”: loss of relevance of Euclidean distance

- Binary Gaussian mixture classification $\mathbf{x} \in \mathbb{R}^p$:

$$\mathcal{C}_1 : \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_1, \mathbf{C}_1), \text{ versus } \mathcal{C}_2 : \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_2, \mathbf{C}_2);$$

- Neyman-Pearson test: classification is possible **only** when

$$\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\| \geq C_{\boldsymbol{\mu}}, \text{ or } \|\mathbf{C}_1 - \mathbf{C}_2\| \geq C_{\mathbf{C}} \cdot p^{-1/2}$$

for some constants $C_{\boldsymbol{\mu}}, C_{\mathbf{C}} > 0$ [CLM18].

- In this **non-trivial** setting, for $\mathbf{x}_i \in \mathcal{C}_a, \mathbf{x}_j \in \mathcal{C}_b$:

$$\max_{1 \leq i \neq j \leq n} \left\{ \frac{1}{p} \|\mathbf{x}_i - \mathbf{x}_j\|^2 - \frac{2}{p} \operatorname{tr} \mathbf{C}^\circ \right\} \rightarrow 0$$

as $n, p \rightarrow \infty$ (i.e., $n \sim p$), for $\mathbf{C}^\circ \equiv \frac{1}{2}(\mathbf{C}_1 + \mathbf{C}_2)$, regardless of the classes $\mathcal{C}_a, \mathcal{C}_b$!

⁰Romain Couillet, Zhenyu Liao, and Xiaoyi Mai. “Classification asymptotics in the random matrix regime”. In: *2018 26th European Signal Processing Conference (EUSIPCO)*. IEEE. 2018, pp. 1875–1879

Loss of relevance of Euclidean distance: visual representation

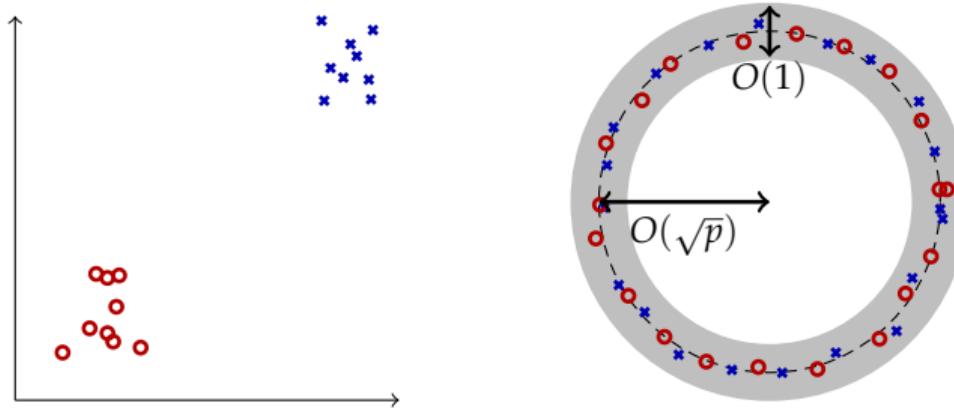
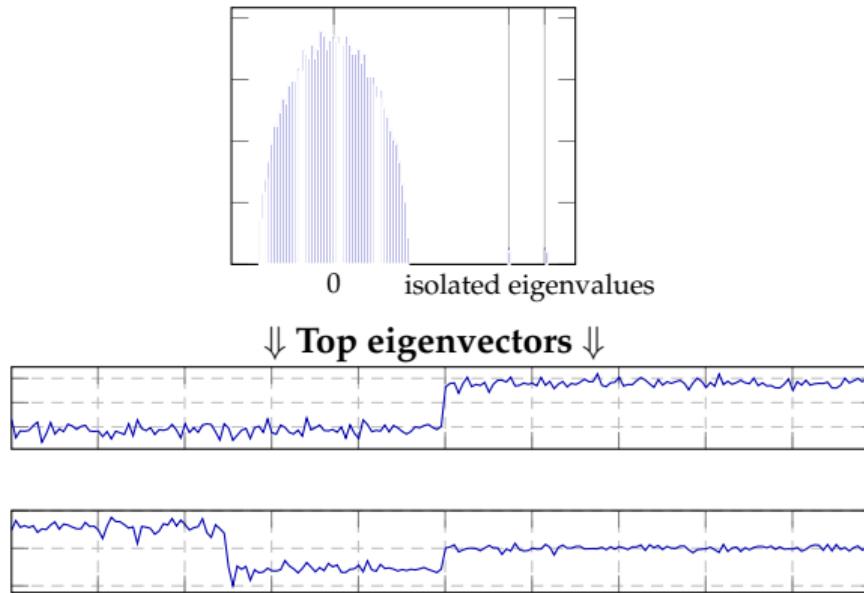


Figure: Visual representation of classification in (left) small and (right) large dimensions.

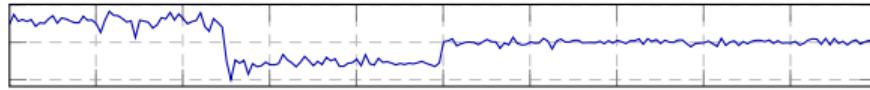
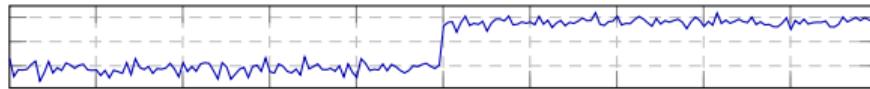
⇒ Direct consequence to various **distance-based** machine learning methods
(e.g., kernel spectral clustering)!

Reminder on kernel spectral clustering

Two-step classification of n data points with distance kernel $\mathbf{K} \equiv \{f(\|\mathbf{x}_i - \mathbf{x}_j\|^2/p)\}_{i,j=1}^n$:

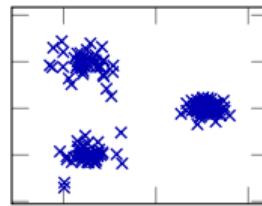


Reminder on kernel spectral clustering



↓ K-dimensional representation ↓

Eig. 2



Eig. 1



EM or k-means clustering

Cluster Gaussian data $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbf{R}^p$ into \mathcal{C}_1 or \mathcal{C}_2 , with second top eigenvectors \mathbf{v}_2 of heat kernel $\mathbf{K}_{ij} = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2p)$, small and large dimensional data.

(a) $p = 5, n = 500$

$$\mathbf{K} = \begin{bmatrix} & \mathcal{C}_1 & \mathcal{C}_2 \\ \mathcal{C}_1 & & \\ & & \mathcal{C}_2 \end{bmatrix}$$

$$\mathbf{v}_2 = [\text{blue wavy line}]$$

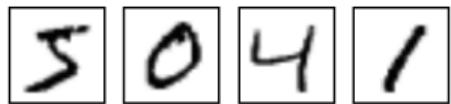
(b) $p = 250, n = 500$

$$\mathbf{K} = \begin{bmatrix} & \mathcal{C}_1 & \mathcal{C}_2 \\ \mathcal{C}_1 & & \\ & & \mathcal{C}_2 \end{bmatrix}$$

$$\mathbf{v}_2 = [\text{blue wavy line}]$$

Kernel matrices for large dimensional real-world data

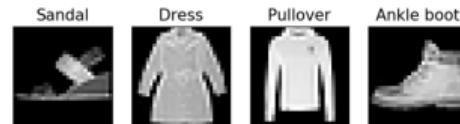
(a) MNIST



$$\mathbf{K} = \begin{bmatrix} & & & \\ & \text{[A 28x28 grid of handwritten digits]} & & \\ & & & \end{bmatrix}$$

$$\mathbf{v}_2 = \begin{bmatrix} \text{[A 1x28 vector of blue noise]} \end{bmatrix}$$

(b) Fashion-MNIST



$$\mathbf{K} = \begin{bmatrix} & & & \\ & \text{[A 28x28 grid of fashion items]} & & \\ & & & \end{bmatrix}$$

$$\mathbf{v}_2 = \begin{bmatrix} \text{[A 1x28 vector of blue noise]} \end{bmatrix}$$

A RMT viewpoint of large kernel matrices

- ▶ “local” **linearization** of *nonlinear* kernel matrices in large dimensions, e.g., Gaussian kernel matrix $\mathbf{K}_{ij} = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2p)$ with $\mathbf{C}_1 = \mathbf{C}_2 = \mathbf{I}_p$ (e.g., $\mathcal{C}_1 : \mathbf{x}_i = \boldsymbol{\mu}_1 + \mathbf{z}_i$ versus $\mathcal{C}_2 : \mathbf{x}_j = \boldsymbol{\mu}_2 + \mathbf{z}_j$) so that

$$\|\mathbf{x}_i - \mathbf{x}_j\|^2/p \rightarrow 2, \text{ and } \mathbf{K} = \exp\left(-\frac{2}{2}\right) \left(\mathbf{1}_n \mathbf{1}_n^\top + \frac{1}{p} \mathbf{Z}^\top \mathbf{Z} \right) + g(\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|) \frac{1}{p} \mathbf{j} \mathbf{j}^\top + * + o_{\|\cdot\|}(1)$$

with Gaussian $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_n] \in \mathbb{R}^{p \times n}$ and class-information $\mathbf{j} = [\mathbf{1}_{n/2}; -\mathbf{1}_{n/2}]$,

- ▶ **accumulated effect** of small “hidden” statistical information ($\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|$ in this case)

A RMT viewpoint of large kernel matrices

Therefore

► entry-wise:

$$\mathbf{K}_{ij} = \exp(-1) \left(1 + \underbrace{\frac{1}{p} \mathbf{z}_i^\top \mathbf{z}_j}_{O(p^{-1/2})} \right) \pm \underbrace{\frac{1}{p} g(\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|)}_{O(p^{-1})} + *, \text{ so that } \frac{1}{p} g(\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|) \ll \frac{1}{p} \mathbf{z}_i^\top \mathbf{z}_j,$$

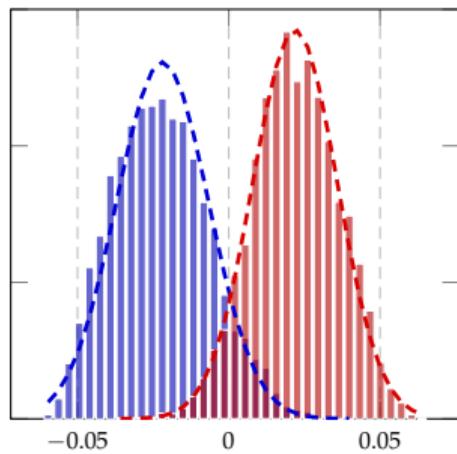
► spectrum-wise:

- $\|\mathbf{K} - \exp(-1)\mathbf{1}_n \mathbf{1}_n^\top\| \not\rightarrow 0$;
- $\|\frac{1}{p} \mathbf{Z}^\top \mathbf{Z}\| = O(1)$ and $\|g(\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|) \frac{1}{p} \mathbf{j} \mathbf{j}^\top\| = O(1)$!

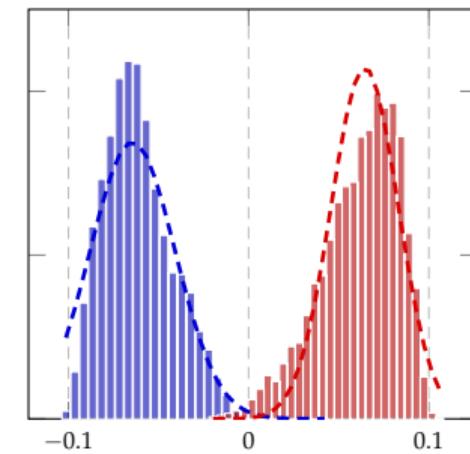
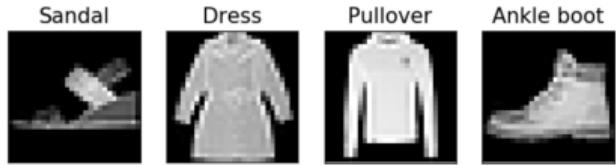
► **Same** phenomenon as the sample covariance example: $[\hat{\mathbf{C}} - \mathbf{C}]_{ij} \rightarrow 0 \not\Rightarrow \|\hat{\mathbf{C}} - \mathbf{C}\| \rightarrow 0$!

⇒ With RMT, we **understand** kernel spectral clustering for large dimensional data!

Some more numerical results



(a) MNIST



(b) Fashion-MNIST

Question: what are deep neural networks?



Deep Learning (DL) \approx multilayered neural network (NN) is becoming the **most** popular machine learning (ML) model, but

- ▶ what is machine learning?
- ▶ what is a deep neural network (DNN)?
- ▶ how is such a network trained?
- ▶ is there any theory for DL, and if yes, how far is the theory from practice?

Credit: most materials in this part are borrowed from [HH19].

¹Catherine F. Higham and Desmond J. Higham. "Deep Learning: An Introduction for Applied Mathematicians". In: *SIAM Review* 61.4 (Jan. 2019), pp. 860–891

Example: binary classification of points in \mathbb{R}^2

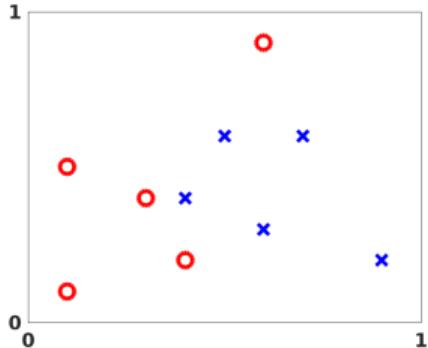


Figure: Labeled data points $x \in \mathbb{R}^2$. Circles denote points in class \mathcal{C}_1 . Crosses denote points in class \mathcal{C}_2 .

- ▶ build a model/**function** f (from above historical data) that takes any points $x \in \mathbb{R}^2$ and returns \mathcal{C}_1 or \mathcal{C}_2
- ▶ **logistic regression:**
$$f(x) = \sigma(\mathbf{w}^T x + b)$$
 for $\mathbf{w} \in \mathbb{R}^2$ and $b \in \mathbb{R}$ to be determined, and **sigmoid** function $\sigma(t) = \frac{1}{1+e^{-t}}$

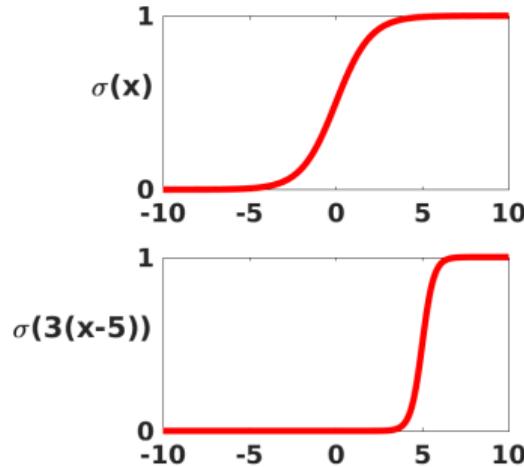


Figure: Sigmoid function.

- ▶ “learn” or estimate parameters \mathbf{w}, b from data/samples, by minimizing some **cost function** (e.g., negative likelihood, MSE)
- ▶ predict $x \in \mathcal{C}_1$ if $f(x) < 1/2$ and $x \in \mathcal{C}_2$ otherwise.

Neural networks are nothing but “cascaded” logistic regressors

- ▶ logistic regression $f(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b) \in \mathbb{R}$ for $\mathbf{w} \in \mathbb{R}^2$, $b \in \mathbb{R}$ extends to

$$f(\mathbf{x}) = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) \in \mathbb{R}^N \quad \mathbf{W} \in \mathbb{R}^{N \times 2}, \mathbf{b} \in \mathbb{R}^N \quad (1)$$

and $\sigma(\cdot)$ applied entry-wise: this is **one layer** of a DNN

- ▶ repeat this to make the network **deep**, with possibly different **width** in each layer
- ▶ $\sigma(\mathbf{W}_2\mathbf{x} + \mathbf{b}_2) \in \mathbb{R}^2$, $\sigma(\mathbf{W}_3\sigma(\mathbf{W}_2\mathbf{x} + \mathbf{b}_2) + \mathbf{b}_3) \in \mathbb{R}^3$
- ▶ $f_{4L-NN}(\mathbf{x}) = \sigma(\mathbf{W}_4\sigma(\mathbf{W}_3\sigma(\mathbf{W}_2\mathbf{x} + \mathbf{b}_2) + \mathbf{b}_3) + \mathbf{b}_4) \in \mathbb{R}^2$

Define the **label**/target output as

$$\mathbf{y}(\mathbf{x}_i) = \begin{cases} \begin{bmatrix} 1 \\ 0 \end{bmatrix} & \mathbf{x}_i \in \mathcal{C}_1, \\ \begin{bmatrix} 0 \\ 1 \end{bmatrix} & \mathbf{x}_i \in \mathcal{C}_2. \end{cases} \quad (2)$$

the MSE cost function writes $\text{Cost}(\mathbf{W}_2, \mathbf{W}_3, \mathbf{W}_4, \mathbf{b}_2, \mathbf{b}_3, \mathbf{b}_4) = \frac{1}{10} \sum_{i=1}^{10} \|\mathbf{y}(\mathbf{x}_i) - f_{4L-NN}(\mathbf{x}_i)\|^2$

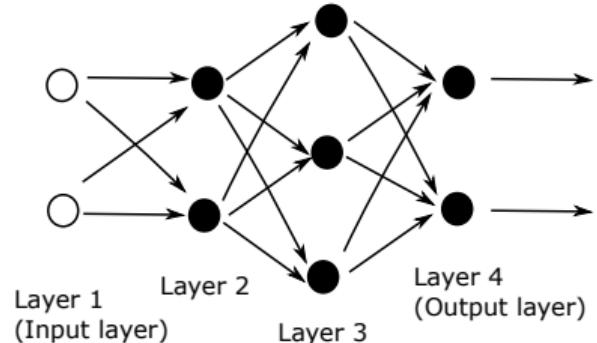


Figure: A network with four layers.

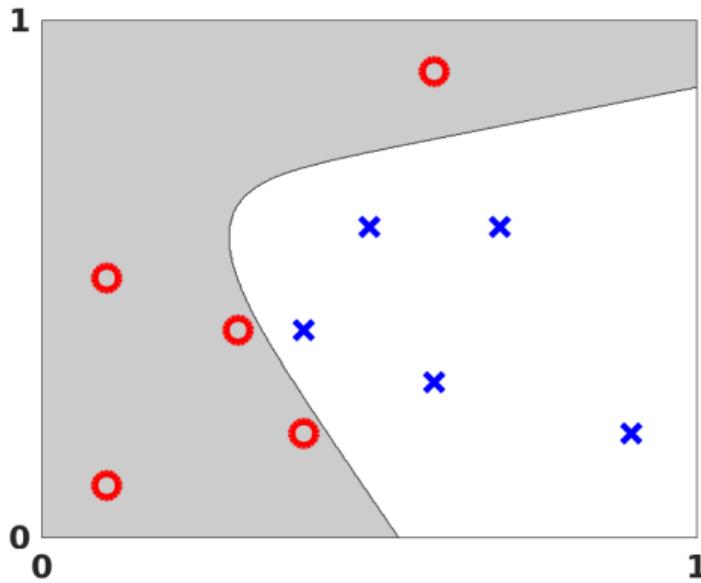


Figure: Visualization of output from a multilayered neural network applied to the data.

- ▶ from training to test!

General formulation and gradient decent training of DNN

We can define the network in a **layer-by-layer** fashion:

$$\mathbf{a}_0 = \mathbf{x} \in \mathbb{R}^{N_0}, \quad \boxed{\mathbf{a}_\ell = \sigma(\mathbf{W}_\ell \mathbf{a}_{\ell-1} + \mathbf{b}_\ell) \in \mathbb{R}^{N_\ell}}, \quad \ell = 1, \dots, L,$$

with weights $\mathbf{W}_\ell \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$ and bias $\mathbf{b} \in \mathbb{R}^{N_\ell}$ at layer ℓ .

- ▶ \mathbf{W}_ℓ s and \mathbf{b}_ℓ s obtained by minimizing **cost function** on a given training set $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ of size n :

$$\text{Cost} = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{y}_i - \mathbf{a}_L(\mathbf{x}_i)\|^2. \tag{3}$$

- ▶ update using (stochastic) gradient descent, for parameter P ,

$$P(t+1) = P(t) - \eta \nabla_P \text{Cost}(P(t)). \tag{4}$$

Matlab code to train a simple NN

```
%%%%% DATA %%%%%%
x1 = [0.1,0.3,0.1,0.6,0.4,0.6,0.5,0.9,0.4,0.7]; x2 = [0.1,0.4,0.5,0.9,0.2,0.3,0.6,0.2,0.4,0.6]; y = [ones(1,5) zeros(1,5); zeros(1,5) ones(1,5)];

% Initialize weights and biases
W2 = 0.5*randn(2,2); W3 = 0.5*randn(3,2); W4 = 0.5*randn(2,3); b2 = 0.5*randn(2,1); b3 = 0.5*randn(3,1); b4 = 0.5*randn(2,1);

% Forward and Back propagate
eta = 0.05; % learning rate
Niter = 1e6; % number of SG iterations
savecost = zeros(Niter,1); % value of cost function at each iteration
for counter = 1:Niter
    k = randi(10); % choose a training point at random
    x = [x1(k); x2(k)];
    % Forward pass
    a2 = activate(x,W2,b2); a3 = activate(a2,W3,b3); a4 = activate(a3,W4,b4);
    % Backward pass
    delta4 = a4.*((1-a4).*(a4-y(:,k))); delta3 = a3.*((1-a3).*(W4'*delta4)); delta2 = a2.*((1-a2).*(W3'*delta3));
    % Gradient step
    W2 = W2 - eta*delta2*x'; W3 = W3 - eta*delta3*a2'; W4 = W4 - eta*delta4*a3'; b2 = b2 - eta*delta2; b3 = b3 - eta*delta3; b4 = b4 - eta*delta4;
    % Monitor progress
    newcost = cost(W2,W3,W4,b2,b3,b4) % display cost to screen
    savecost(counter) = newcost;
end

% Show decay of cost function
semilogy([1:1e4:Niter],savecost(1:1e4:Niter))

function costval = cost(W2,W3,W4,b2,b3,b4)
    costvec = zeros(10,1);
    for i = 1:10
        x =[x1(i);x2(i)];
        a2 = activate(x,W2,b2); a3 = activate(a2,W3,b3); a4 = activate(a3,W4,b4);
        costvec(i) = norm(y(:,i) - a4,2);
    end
    costval = norm(costvec,2)^2;
end % of nested function
```

Matlab code to train a simple NN

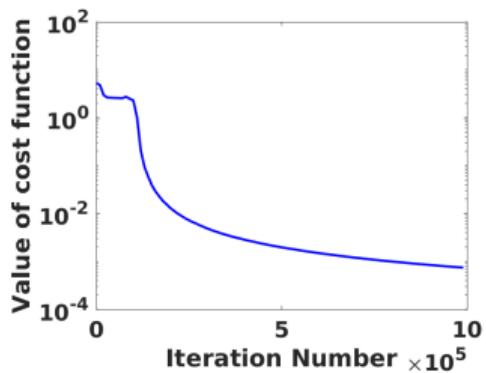


Figure: Vertical axis shows a scaled value of the cost function. Horizontal axis shows the iteration number. Here we used the stochastic gradient descent to train the aforementioned simple network.

Some commonly used tricks in DNN

- ▶ **stochastic gradient descent**: sample (without replacement) a mini-batch for gradient $\frac{1}{B} \sum_{i=1}^B \nabla_p \text{Cost}(\mathbf{x}_i)$
- ▶ **convolution neural network (CNN)**: repeatedly apply small linear **kernel**, or filter, across portions of input data, making weight matrices **sparse** and highly **structured**

$$\begin{bmatrix} 1 & -1 & & & & \\ & 1 & -1 & & & \\ & & 1 & -1 & & \\ & & & 1 & -1 & \\ & & & & 1 & -1 \\ & & & & & 1 & -1 \end{bmatrix} \in \mathbb{R}^{5 \times 6}. \quad (5)$$

- ▶ different choice of activation and/or cost function:
 - rectified linear unit, or **ReLU**, activation: $\sigma(t) = \max(t, 0)$
 - **cross-entropy** cost function:
- ▶ dropout, batch normalization, and other types of normalization, etc.
- ▶ use of **tensors** instead of vectors or matrices for input data or intermediate representations

What do we care about DL, from a theoretical perspective?

What does **deep learning theory** care about and why?

- ▶ **theoretical guarantee:** explanation of **when** and **why** DL works in some cases, and not in others
- ▶ theory-guided **design principles** for more efficient DNN (e.g., better performant, less computational demand, more novel ideas on how to make DL work better, etc.)
- ▶ **too many** “tuning” hyperparameters in DNN design: number of layers, operator, width, and activation in each layer, different tricks, etc.
- ▶ for safety-related applications (e.g., self driving, healthcare), we need theory-supported DL that
 - ① allows us to **combine** domain knowledge in DNN design
 - ② can be used **safely**

A (too) brief review of DL theory

From an approximation theoretical perspective:

- ▶ **universal approximation theorem:** for any (somewhat regular, e.g., Lebesgue p -integrable) function of interest $f: \mathbb{R}^{p \times K}$ and given $\varepsilon > 0$, there exists a **fully-connected ReLU network** F with **width** at least m such that $\int_{\mathbb{R}^p} \|f(\mathbf{x}) - F(\mathbf{x})\|^p d\mathbf{x} < \varepsilon$.
- ▶ different type of input space, e.g., $\mathbf{x} = [x_1, \dots, x_p] \subset [0, 1]^p$, function or data on graph?
- ▶ how activation, width, depth, etc. come into play, in particular, **depth versus width?**
- ▶ **LIMITATION:** do not provide a construction for the network, but that such a construction is **possible**

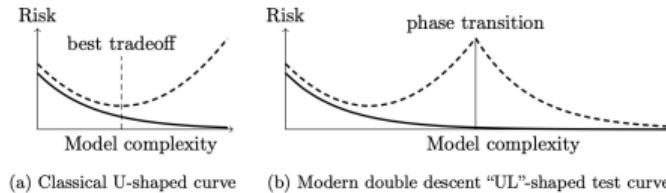
From an optimization perspective:

- ▶ DNN training involves **non-convex** (and possibly non-smooth) optimization: challenging!
- ▶ **empirically** simple (stochastic) gradient descent seems to work well, **WHY?**
- ▶ GUESS: DL landscape has nice properties?
- ▶ e.g., how to converge better and faster?
- ▶ **IMPORTANT:** pure optimization deals **only** with training, and **NOT** test/**generalization**

A (too) brief review of DL theory

From a statistical perspective:

- ▶ **generalization theory**: for which type of **data**, and by using which **ML model** (trained with which **algorithm**), can we get a high probability error bound of which **metric**
- ▶ Rademacher complexity (distribution-dependent in general), PAC-Bayes bound, etc.
- ▶ **Question**: why DL models **generalize so well** despite high model complexity (i.e., **over-parameterized**)?
 - ① nice property of the (over-parameterized) DL model: Neural Tangent Kernel [JGH18]
 - ② inductive bias due to algorithm: Double Descent or Benign Overfitting [BMR21]



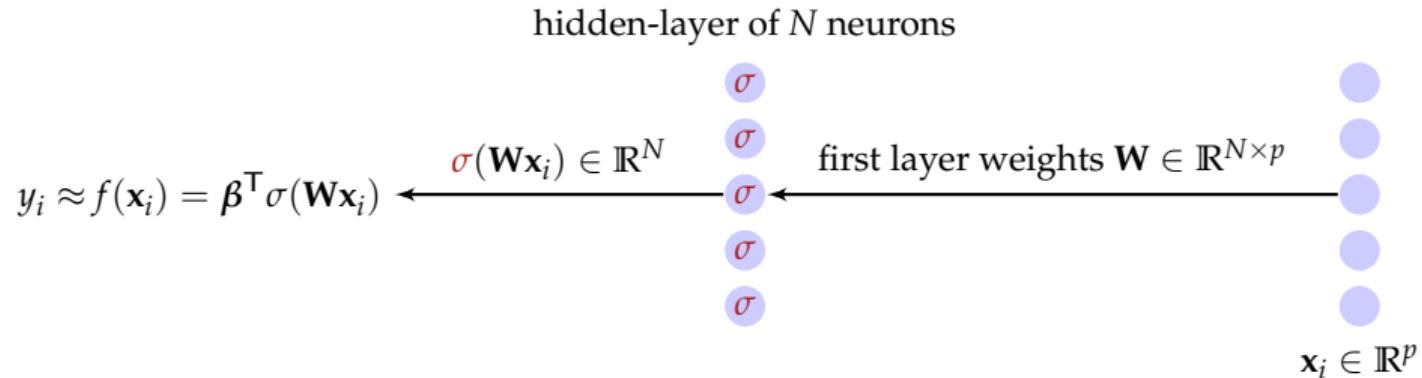
A Good DL theory should cover **both optimization and generalization!**

²Arthur Jacot, Franck Gabriel, and Clément Hongler. "Neural Tangent Kernel: Convergence and Generalization in Neural Networks". In: *Advances in Neural Information Processing Systems*. Vol. 31. NIPS'18. Curran Associates, Inc., 2018, pp. 8571–8580

³Peter L. Bartlett, Andrea Montanari, and Alexander Rakhlin. "Deep Learning: A Statistical Viewpoint". In: *Acta Numerica* 30 (May 2021), pp. 87–201

- ▶ kernel $K(\cdot, \cdot) : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$, **similarity** measure between input data points in \mathbb{R}^p
- ▶ examples include:
 - linear kernel $K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y}$, cosine kernel = $\frac{\mathbf{x}^\top \mathbf{y}}{\|\mathbf{x}\| \cdot \|\mathbf{y}\|}$, Gaussian (RBF) kernel = $\exp(-\|\mathbf{x} - \mathbf{y}\|^2 / \gamma^2)$
 - kernel induced by NN: $K(\mathbf{x}, \mathbf{y}) = \sigma(\mathbf{Wx})^\top \sigma(\mathbf{Wy})$, parameterized by the network (e.g., weights and activations)
- ▶ PS: kernels are widely studied in the ML literature, we know quite a lot (reproducing kernel Hilbert space, RKHS, etc.)

Example of a two-layer NN model



- Given training set $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with $\mathbf{x}_i \in \mathbb{R}^p$ and $y_i \in \mathbb{R}$

$$f(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\beta}^\top \sigma(\mathbf{W}\mathbf{x}) = \sum_{\ell=1}^n \beta_\ell \sigma(\mathbf{w}_\ell^\top \mathbf{x}), \quad \boldsymbol{\theta} = [\beta_1, \dots, \beta_N; \mathbf{w}_1, \dots, \mathbf{w}_N]. \quad (7)$$

- linearization of the network at initialization, by Taylor expansion

$$f(\mathbf{x}; \boldsymbol{\theta}) \approx f_{\text{lin}}(\mathbf{x}; \boldsymbol{\theta}) = f(\mathbf{x}; \boldsymbol{\theta}_0) + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \boxed{\nabla_{\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}_0)}. \quad (8)$$

and

$$f_{\text{lin}}(\mathbf{x}; \boldsymbol{\theta}_0 + \boldsymbol{\delta}) = f(\mathbf{x}; \boldsymbol{\theta}_0) + \boldsymbol{\delta}^\top \phi_{\text{NTK}}(\mathbf{x}), \quad K_{e-\text{NTK}}(\mathbf{x}, \mathbf{y}) = \phi_{\text{NTK}}(\mathbf{x})^\top \phi_{\text{NTK}}(\mathbf{y}). \quad (9)$$

The big picture of NTK

- ▶ around initialization $\theta \approx \theta_0$, linearized network output

$$f(\mathbf{x}; \theta) \approx f_{\text{lin}}(\mathbf{x}; \theta_0 + \delta) = f(\mathbf{x}; \theta_0) + \delta^T \phi_{\text{NTK}}(\mathbf{x}), \quad K_{e-\text{NTK}}(\mathbf{x}, \mathbf{y}) = \phi_{\text{NTK}}(\mathbf{x})^T \phi_{\text{NTK}}(\mathbf{y}), \quad (10)$$

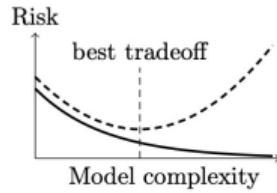
Now, if there exists a neighborhood $B(\theta_0)$ of θ_0 such that

- ① for any $\theta \in B(\theta_0)$, we have $f(\mathbf{x}; \theta) \approx f_{\text{lin}}(\mathbf{x}; \theta)$, and closeness in cost function
- ② it suffices to optimize in $B(\theta_0)$ to reach an approx. global min, i.e., $f(\mathbf{x}; \theta_0) \approx f_{\text{lin}}(\mathbf{x}; \theta_0) \approx 0$
- ③ from an optimization viewpoint, optimizing $f(\mathbf{x}; \theta) \approx f_{\text{lin}}$ and will not leave $B(\theta_0)$

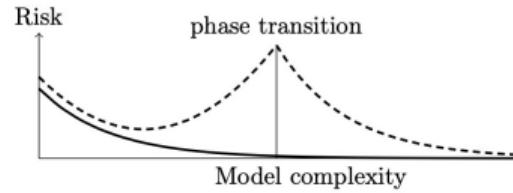
Till now, the major way to reach the above is **over-parameterization** and/or **proper random initialization**, with **small** stochasticity (e.g., small learning rate or full batch GD)

- ▶ cost function (e.g., MSE) $\text{Cost}(f_\theta(\mathbf{x}), \mathbf{y}) \approx \text{Cost}(f_{\text{lin}}(\mathbf{x}), \mathbf{y})$ linear (in the parameter θ) and convex!
- ▶ for MSE, $\text{Cost}(f_{\text{lin}}(\mathbf{X}), \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n (f_{\text{lin}}(\mathbf{x}_i) - y_i)^2$, nothing but linear regression of type
 $\text{Cost} = \|\mathbf{y}' - \Phi_{\text{NTK}}(\mathbf{X})^T \delta\|^2$ with $y'_i = f(\mathbf{x}_i; \theta_0) - y_i$

Precise Characterization of Double Descent Curves

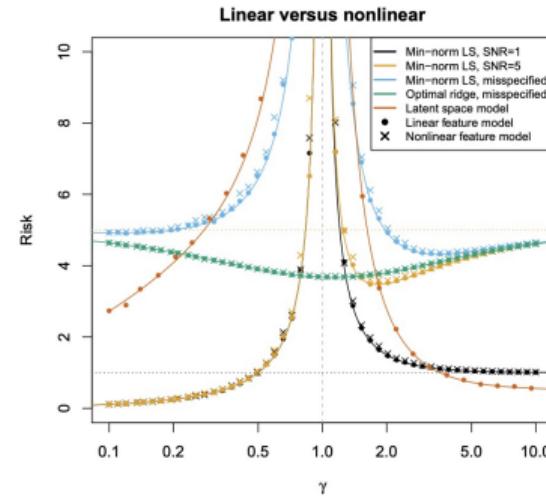


(a) Classical U-shaped curve



(b) Modern double descent “UL”-shaped test curve

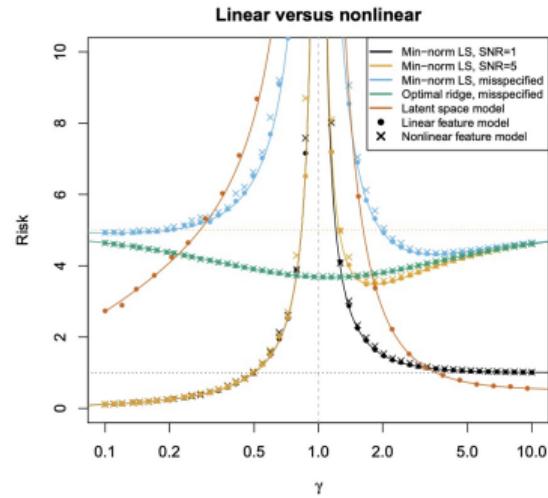
- ▶ larger model, the better?! Maybe, due to double descent [Has+22] and implicit (norm-based?) bias
- ▶ case of linear regression model
 $\text{Cost} = \frac{1}{n} \sum_{i=1}^n (\beta^\top x_i - y_i)^2$, with $\beta, x_i \in \mathbb{R}^p$,
depend on the sign $n - p$, either in the
over-parameterized or under-parameterized
(with min-norm solution) regime
- ▶ **generalization** risk shows a double descent curve



⁴Trevor Hastie et al. “Surprises in High-Dimensional Ridgeless Least Squares Interpolation”. In: *The Annals of Statistics* 50.2 (Apr. 2022), pp. 949–986

Precise Characterization of Double Descent Curves

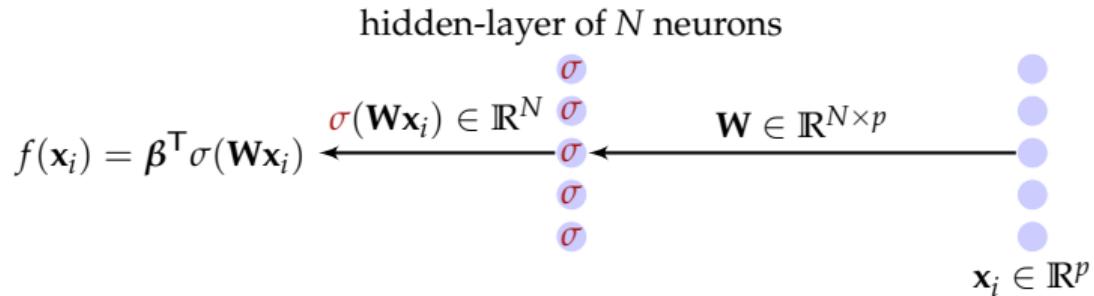
- ▶ case of linear regression model
 $\text{Cost} = \frac{1}{n} \sum_{i=1}^n (\beta^\top \mathbf{x}_i - y_i)^2$, with $\beta, \mathbf{x}_i \in \mathbb{R}^p$,
depend on the sign $n - p$, either in the
over-parameterized or **under-parameterized**
(with min-norm solution) regime
- ▶ **generalization** risk shows a double descent
curve [Has+22]
- ▶ very **understandable** for RMT experts:
- ▶ ridgeless least squares $\hat{\beta} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}$ or
 $\hat{\beta} = \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{y}$ and there is a **singular behavior**
in the spectrum at $p = n$
- ▶ tons of extensions: relaxing assumption, (slightly)
more involved models, etc., less progress in the
sense of **deep** though



Technical challenges and opportunities for RMT in DL theory

- ▶ entry-wise non-linearity and depth: some **successful** efforts
- ▶ gradient descent leads to involved correlation structure: even a **single** step makes things complicated
- ▶ statistical assumption to work with: largely **open!**

Two-layer network with random first layer



- ▶ for random (first-layer) weights $\mathbf{W} \in \mathbb{R}^{N \times p}$ having say i.i.d. standard Gaussian entries
- ▶ get second-layer $\boldsymbol{\beta}$ by minimizing $\text{Cost} = \frac{1}{n} \sum_{i=1}^n (y_i - \boldsymbol{\beta}^\top \sigma(\mathbf{W}\mathbf{x}_i))^2 + \gamma \|\boldsymbol{\beta}\|^2$ for some regularization parameter $\gamma > 0$, then

$$\boldsymbol{\beta} \equiv \frac{1}{n} \boldsymbol{\Sigma} \left(\frac{1}{n} \boldsymbol{\Sigma}^\top \boldsymbol{\Sigma} + \gamma \mathbf{I}_n \right)^{-1} \mathbf{y}, \quad (11)$$

- ▶ training MSE (on the given training set (\mathbf{X}, \mathbf{y})) reads

$$E_{\text{train}} = \frac{1}{n} \|\mathbf{y} - \boldsymbol{\Sigma}^\top \boldsymbol{\beta}\|_F^2 = \frac{\gamma^2}{n} \mathbf{y} \mathbf{Q}^2(\gamma) \mathbf{y}, \quad \boxed{\mathbf{Q}(\gamma) \equiv \left(\frac{1}{n} \boldsymbol{\Sigma}^\top \boldsymbol{\Sigma} + \gamma \mathbf{I}_n \right)^{-1}} \quad (12)$$

- ▶ Similarly, the test MSE on a test set $(\hat{\mathbf{X}}, \hat{\mathbf{y}}) \in \mathbb{R}^{p \times \hat{n}} \times \mathbb{R}^{d \times \hat{n}}$ of size \hat{n} : $E_{\text{test}} = \frac{1}{\hat{n}} \|\hat{\mathbf{y}} - \hat{\boldsymbol{\Sigma}}^\top \boldsymbol{\beta}\|_F^2$, $\hat{\boldsymbol{\Sigma}} = \sigma(\mathbf{W}\hat{\mathbf{X}})$.

$$\mathbf{Q}(\gamma) = \left(\frac{1}{n} \sigma(\mathbf{W}\mathbf{X})^\top \sigma(\mathbf{W}\mathbf{X}) + \gamma \mathbf{I}_n \right)^{-1} \quad (13)$$

- ▶ nonlinear $\Sigma^\top = \sigma(\mathbf{W}\mathbf{X})^\top$ still has i.i.d. columns, but
- ▶ its i -th column $\sigma([\mathbf{X}^\top \mathbf{W}^\top]_{\cdot i})$ no longer has i.i.d. or linearly dependent entries
- ▶ **trace lemma** does not apply

Lemma (Concentration of nonlinear quadratic form, [LLC18, Lemma 1])

For $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$, 1-Lipschitz $\sigma(\cdot)$, and $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{X} \in \mathbb{R}^{p \times n}$ such that $\|\mathbf{A}\|, \|\mathbf{X}\|$ bounded, then

$$\mathbb{P} \left(\left| \frac{1}{n} \sigma(\mathbf{w}^\top \mathbf{X}) \mathbf{A} \sigma(\mathbf{X}^\top \mathbf{w}) - \frac{1}{n} \text{tr } \mathbf{AK} \right| > t \right) \leq C e^{-cn \min(t, t^2)}$$

for some $C, c > 0$, $p/n \in (0, \infty)$ with $\mathbf{K} \equiv \mathbf{K}_{\mathbf{XX}} \equiv \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)} [\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})] \in \mathbb{R}^{n \times n}$.

- ▶ \mathbf{K} is in fact the **conjugate kernel (CK)** matrix
- ▶ for well-behaved (e.g., Lipschitz) non-linearity, trace lemma holds in this **nonlinear** case
- ▶ get deterministic equivalent for \mathbf{Q} , establish (limiting) eigenvalue distribution of $\frac{1}{n} \sigma(\mathbf{W}\mathbf{X})^\top \sigma(\mathbf{W}\mathbf{X})$, etc.

Theorem (Resolvent for nonlinear Gram matrix, [LLC18])

Let $\mathbf{W} \in \mathbb{R}^{N \times p}$ be a random matrix with i.i.d. standard Gaussian entries, $\sigma(\cdot)$ be 1-Lipschitz, and $\mathbf{X} \in \mathbb{R}^{p \times n}$ be of bounded operator norm. Then, as $n, p, N \rightarrow \infty$ at the same pace, for $\mathbf{Q} = (\sigma(\mathbf{X}^\top \mathbf{W}^\top) \sigma(\mathbf{W} \mathbf{X}) / n + \gamma \mathbf{I}_n)^{-1}$ with $\gamma > 0$,

$$\|\mathbb{E}[\mathbf{Q}] - \bar{\mathbf{Q}}\| \rightarrow 0, \quad \bar{\mathbf{Q}} \equiv \left(\frac{N}{n} \frac{\mathbf{K}}{1+\delta} + \gamma \mathbf{I}_n \right)^{-1}$$

for δ the unique positive solution to $\delta = \frac{1}{n} \operatorname{tr} \bar{\mathbf{Q}} \mathbf{K}$ and $\mathbf{K} = \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)} [\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})] \in \mathbb{R}^{n \times n}$.

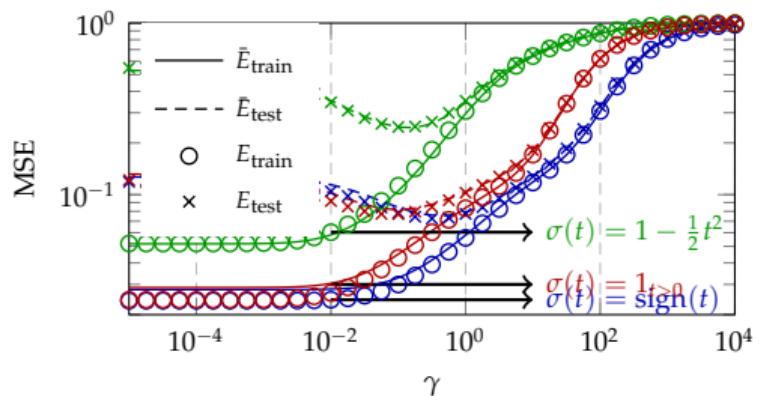
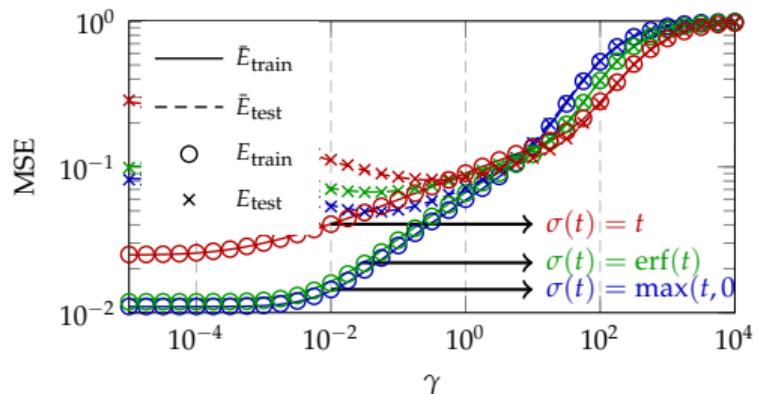
Corollary (Asymptotic training and test MSEs)

Under the setting and notations of Theorem 2, for bounded $\mathbf{X}, \hat{\mathbf{X}}, \mathbf{y}, \hat{\mathbf{y}}$, then the training and test MSES, satisfy, as $n, p, N \rightarrow \infty$, we have $E_{\text{train}} - \bar{E}_{\text{train}} \rightarrow 0$ and $E_{\text{test}} - \bar{E}_{\text{test}} \rightarrow 0$ with

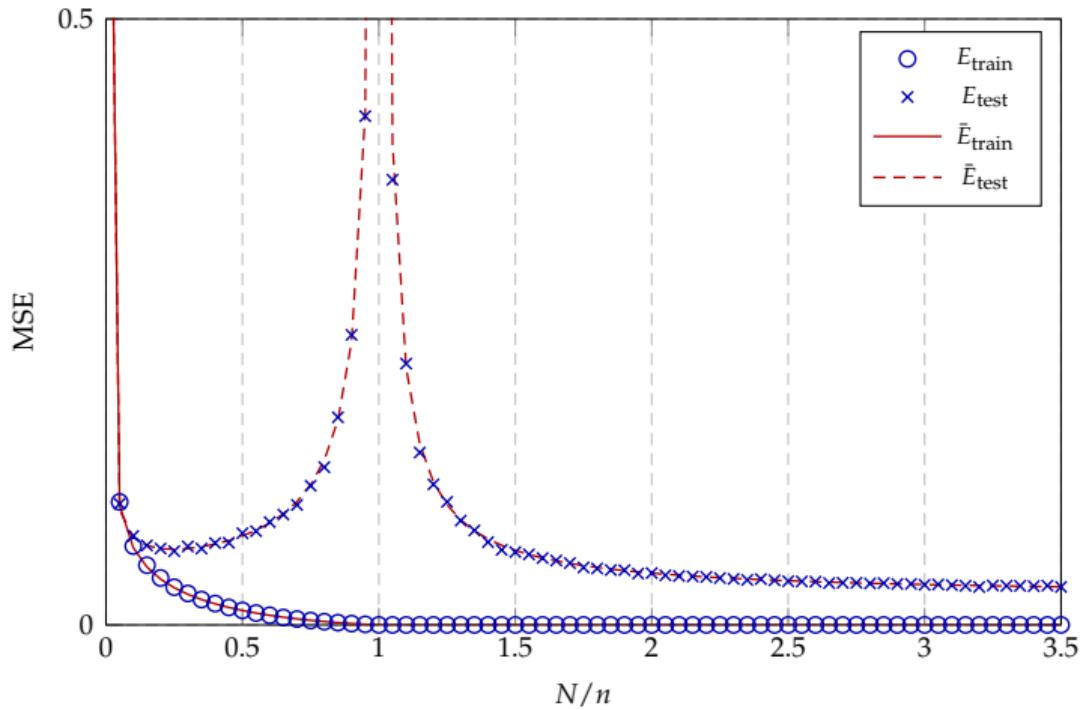
$$\begin{aligned}\bar{E}_{\text{train}} &= \frac{\gamma^2}{n} \mathbf{y}^\top \bar{\mathbf{Q}} \left(\frac{\frac{1}{N} \operatorname{tr} \bar{\mathbf{Q}} \bar{\mathbf{K}} \bar{\mathbf{Q}}}{1 - \frac{1}{N} \operatorname{tr} \bar{\mathbf{K}} \bar{\mathbf{Q}} \bar{\mathbf{K}} \bar{\mathbf{Q}}} \bar{\mathbf{K}} + \mathbf{I}_n \right) \bar{\mathbf{Q}} \mathbf{y} \\ \bar{E}_{\text{test}} &= \frac{1}{\hat{n}} \|\hat{\mathbf{y}} - \bar{\mathbf{K}}_{\mathbf{X} \hat{\mathbf{X}}}^\top \bar{\mathbf{Q}} \mathbf{y}\|_F^2 + \frac{\frac{1}{N} \mathbf{y}^\top \bar{\mathbf{Q}} \bar{\mathbf{K}} \bar{\mathbf{Q}} \mathbf{y}}{1 - \frac{1}{N} \operatorname{tr} \bar{\mathbf{K}} \bar{\mathbf{Q}} \bar{\mathbf{K}} \bar{\mathbf{Q}}} \left(\frac{1}{\hat{n}} \operatorname{tr} \bar{\mathbf{K}}_{\hat{\mathbf{X}} \hat{\mathbf{X}}} - \frac{1}{\hat{n}} \operatorname{tr} (\mathbf{I}_n + \gamma \bar{\mathbf{Q}}) (\bar{\mathbf{K}}_{\mathbf{X} \hat{\mathbf{X}}} \bar{\mathbf{K}}_{\mathbf{X} \hat{\mathbf{X}}}^\top \bar{\mathbf{Q}}) \right)\end{aligned}$$

⁵Cosme Louart, Zhenyu Liao, and Romain Couillet. "A random matrix approach to neural networks". In: *Annals of Applied Probability* 28.2 (2018), pp. 1190–1248

Numerical results



Numerical results: double descent



Some further RMT investigations on the two-layer model

Eigenspectra of $\frac{1}{n}\sigma(\mathbf{W}\mathbf{X})^T\sigma(\mathbf{W}\mathbf{X})$:

- ▶ [PW17] first guess expression of the eigenvalue behavior
- ▶ [BP21]: eigenvalue distribution of $\frac{1}{n}\sigma(\mathbf{W}\mathbf{X})^T\sigma(\mathbf{W}\mathbf{X})$ for \mathbf{W}, \mathbf{X} having sub-gaussian entries
 - ① for “centered” $\sigma(\cdot)$ with respect to Gaussian measure: $\mathbb{E}[\sigma(\xi)] = 0$ for $\xi \sim \mathcal{N}(0, 1)$
 - ② take a rather **explicit** form (3rd order poly ST equation) and depends on σ only via $\mathbb{E}[\sigma^2(\xi)]$ and $\mathbb{E}[\sigma(\xi)\xi]$.
- ▶ [BP22]: behavior of largest eigenvalue of $\frac{1}{n}\sigma(\mathbf{W}\mathbf{X})^T\sigma(\mathbf{W}\mathbf{X})$ for sub-gaussian \mathbf{W}, \mathbf{X} and centered $\sigma(\cdot)$
- ▶ despite being a **white model**, spikes may appear!
 - ① if $\mathbb{E}[\xi^2\sigma(\xi)] = 0$, then **no** spike
 - ② otherwise, at most **two** spikes

Question: what happen if either \mathbf{W} or \mathbf{X} has some structure? Any different **phase transition** behavior?

⁶Jeffrey Pennington and Pratik Worah. “Nonlinear random matrix theory for deep learning”. In: *Advances in Neural Information Processing Systems*. Vol. 30. NIPS’17. Curran Associates, Inc., 2017, pp. 2637–2646

⁷Lucas Benigni and Sandrine Péché. “Eigenvalue Distribution of Some Nonlinear Models of Random Matrices”. In: *Electronic Journal of Probability* 26.none (Jan. 2021), pp. 1–37

⁸Lucas Benigni and Sandrine Péché. *Largest Eigenvalues of the Conjugate Kernel of Single-Layered Neural Networks*. Jan. 2022. arXiv: 2201.04753 [cs, math]

Some further RMT investigations on random DNNs

- ▶ design of DNN to achieve **dynamical isometry**, **accelerate** training at the **beginning** stage of training
- ▶ Jeffrey Pennington, Samuel Schoenholz, and Surya Ganguli. "Resurrecting the sigmoid in deep learning through dynamical isometry: theory and practice". In: *Advances in Neural Information Processing Systems*. Vol. 30. NIPS'17. Curran Associates, Inc., 2017, pp. 4785–4795
- ▶ Minmin Chen, Jeffrey Pennington, and Samuel Schoenholz. "Dynamical Isometry and a Mean Field Theory of RNNs: Gating Enables Signal Propagation in Recurrent Neural Networks". In: *Proceedings of the 35th International Conference on Machine Learning*. Vol. 80. Proceedings of Machine Learning Research. Stockholmsmässan, Stockholm Sweden: PMLR, 2018, pp. 873–882
- ▶ Lechao Xiao et al. "Dynamical Isometry and a Mean Field Theory of CNNs: How to Train 10,000-Layer Vanilla Convolutional Neural Networks". In: *Proceedings of the 35th International Conference on Machine Learning*. Vol. 80. Proceedings of Machine Learning Research. Stockholmsmässan, Stockholm Sweden: PMLR, 2018, pp. 5393–5402
- ▶ Dar Gilboa et al. "Dynamical Isometry and a Mean Field Theory of LSTMs and GRUs". In: *arXiv* (2019). eprint: 1901.08987
- ▶ understand how weight distribution **interact** with activation in DNNs
- ▶ Leonid Pastur. "On Random Matrices Arising in Deep Neural Networks. Gaussian Case". In: *arXiv* (2020). eprint: 2001.06188
- ▶ Leonid Pastur and Victor Slavin. "On Random Matrices Arising in Deep Neural Networks: General I.I.D. Case". In: *Random Matrices: Theory and Applications* 12.01 (Jan. 2023), p. 2250046
- ▶ Leonid Pastur. "Eigenvalue Distribution of Large Random Matrices Arising in Deep Neural Networks: Orthogonal Case". In: *Journal of Mathematical Physics* 63.6 (2022), p. 063505
- ▶ Zhou Fan and Zhichao Wang. "Spectra of the Conjugate Kernel and Neural Tangent Kernel for linear-width neural networks". In: *arXiv* (2020). eprint: 2005.11879

Gradient descent dynamics on linear regression model

- ▶ **gradient descent dynamics (GDDs)** of ridge regression learning (i.e., of a single-layer linear network)
- ▶ given training data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$ with associated labels/targets $\mathbf{y} = [y_1, \dots, y_n] \in \mathbb{R}^n$, $\mathbf{w} \in \mathbb{R}^p$ is learned via gradient descent by minimizing the (ridge-regularized) squared loss

$$L(\mathbf{w}) = \frac{1}{2n} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2 + \frac{\gamma}{2} \|\mathbf{w}\|^2 \quad (14)$$

for some regularization penalty $\gamma \geq 0$.

- ▶ gradient given by $\nabla L(\mathbf{w}) = -\frac{1}{n} \mathbf{X}(\mathbf{y} - \mathbf{X}^\top \mathbf{w}) + \gamma \mathbf{w}$ so that, for small gradient descent steps (or learning rate) α , **continuous-time approximation** (in fact, **gradient flow**) of the time evolution $\mathbf{w}(t)$ of \mathbf{w} :

$$\frac{\partial \mathbf{w}(t)}{\partial t} = -\alpha \nabla L(\mathbf{w}) = \frac{\alpha}{n} \mathbf{X} \mathbf{y} - \alpha \left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top + \gamma \mathbf{I}_p \right) \mathbf{w}$$

solution explicitly given by

$$\boxed{\mathbf{w}(t) = e^{-\alpha t \left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top + \gamma \mathbf{I}_p \right)} \mathbf{w}_0 + \left(\mathbf{I}_p - e^{-\alpha t \left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top + \gamma \mathbf{I}_p \right)} \right) \mathbf{w}_\infty} \quad (15)$$

with $\mathbf{w}_0 = \mathbf{w}(t=0)$ (the initialization of gradient descent) and

$$\mathbf{w}_\infty = \left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top + \gamma \mathbf{I}_p \right)^{-1} \frac{1}{n} \mathbf{X} \mathbf{y} \quad (16)$$

the ridge regression solution with regularization parameter γ .

Some RMT results on GDD in classification

- ▶ to study statistical evolution of $\mathbf{w}(t)$, consider binary Gaussian mixture model for input data

$$\mathcal{C}_1 : \mathbf{x}_i \sim \mathcal{N}(-\boldsymbol{\mu}, \mathbf{I}_p) \quad \mathcal{C}_2 : \mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{I}_p)$$

with associated labels $y_i = -1$ and $y_i = 1$, respectively.

- ▶ study training and test **misclassification error rates** as

$$\mathbb{P}(\mathbf{x}_i^T \mathbf{w}(t) > 0 \mid y_i = -1), \quad \text{and} \quad \mathbb{P}(\hat{\mathbf{x}}^T \mathbf{w}(t) > 0 \mid \hat{y} = -1),$$

for $\hat{\mathbf{x}} \sim \mathcal{N}(-\boldsymbol{\mu}, \mathbf{I}_p)$ a new test datum (independent of the training set (\mathbf{X}, \mathbf{y})) of genuine label $\hat{y} = -1$.

- ▶ we can of course consider different statistical **model** and/or different **task** (e.g., regression)

Some RMT results on GDDs

Theorem (Training and test performance of GDD, [LC18])

For a random initialization $\mathbf{w}_0 \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_p / p)$ independent of \mathbf{X} , \mathbf{x} a column of \mathbf{X} of mean $\boldsymbol{\mu}$ and $\hat{\mathbf{x}}$ an independent copy of \mathbf{x} , as $n, p \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$, we have

$$\mathbb{P}(\hat{\mathbf{x}}^\top \mathbf{w}(t) > 0 \mid \hat{y} = -1) - Q\left(\frac{E_{\text{test}}}{\sqrt{V_{\text{test}}}}\right) \rightarrow 0, \quad \mathbb{P}(\mathbf{x}^\top \mathbf{w}(t) > 0 \mid y = -1) - Q\left(\frac{E_{\text{train}}}{\sqrt{V_{\text{train}}}}\right) \rightarrow 0,$$

almost surely, where

$$E_{\text{test}} = -\frac{1}{2\pi i} \oint_{\Gamma} \frac{1-f_t(z)}{z} \frac{\rho m(z) dz}{(\rho + c) m(z) + 1}, \quad V_{\text{test}} = \frac{1}{2\pi i} \oint_{\Gamma} \left[\frac{\frac{1}{z^2} (1-f_t(z))^2}{(\rho + c) m(z) + 1} - \sigma^2 f_t^2(z) m(z) \right] dz$$
$$E_{\text{train}} = -\frac{1}{2\pi i} \oint_{\Gamma} \frac{1-f_t(z)}{z} \frac{dz}{(\rho + c) m(z) + 1}, \quad V_{\text{train}} = \frac{1}{2\pi i} \oint_{\Gamma} \left[\frac{\frac{1}{z} (1-f_t(z))^2}{(\rho + c) m(z) + 1} - \sigma^2 f_t^2(z) z m(z) \right] dz - E_{\text{train}}^2$$

with $\rho = \lim_{p \rightarrow \infty} \|\boldsymbol{\mu}\|^2$, Γ a positive contour surrounding the support of the Marčenko–Pastur law (shifted by $\gamma \geq 0$) and the points $(\gamma, 0)$ and $(\gamma + \lambda_s, 0)$ with $\lambda_s = c + 1 + \rho + c/\rho$, $f_t(z) \equiv \exp(-\alpha t z)$ and $m(z)$ unique ST solution to $c(z - \gamma)m^2(z) - (1 - c - z + \gamma)m(z) + 1 = 0$.

Some further simplifications

- ▶ choose the contour Γ as, e.g., rectangle circling around both **main bulk** and **isolated eigenvalue** (if any)

This leads to

$$E_{\text{test}} = \int \frac{1 - f_t(x + \gamma)}{x + \gamma} \omega(dx) \quad V_{\text{test}} = \frac{\rho + c}{\rho} \int \frac{(1 - f_t(x + \gamma))^2 \omega(dx)}{(x + \gamma)^2} + \sigma^2 \int f_t^2(x + \gamma) \mu(dx)$$

$$E_{\text{train}} = \frac{\rho + c}{\rho} \int \frac{1 - f_t(x + \gamma)}{x + \gamma} \omega(dx), \quad V_{\text{train}} = \frac{\rho + c}{\rho} \int \frac{x(1 - f_t(x + \gamma))^2 \omega(dx)}{(x + \gamma)^2} + \sigma^2 \int x f_t^2(x + \gamma) \mu(dx) - E_{\text{train}}^2$$

where we recall $\rho = \lim \|\mu\|^2$, $f_t(x) = \exp(-\alpha t x)$, $\mu(x)$ the MP law

$$\mu(dx) = \frac{\sqrt{(x - \lambda_-)^+ (\lambda_+ - x)^+}}{2\pi cx} dx + (1 - c^{-1})^+ \delta(x), \quad (17)$$

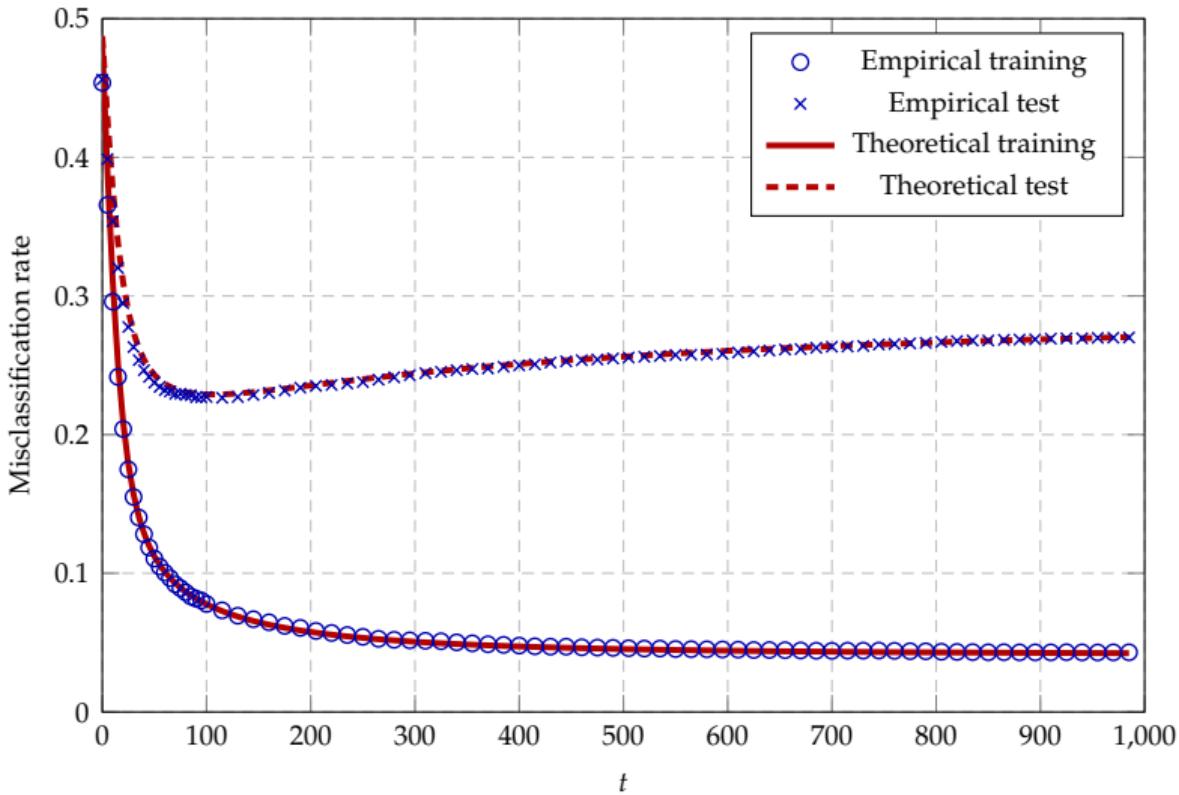
and

$$\omega(dx) \equiv \frac{\sqrt{(x - \lambda_-)^+ (\lambda_+ - x)^+}}{2\pi(\lambda_s - x)} dx + \frac{(\rho^2 - c)^+}{\rho} \delta_{\lambda_s}(x) \quad (18)$$

for $\lambda_s = c + 1 + \rho + c/\rho$ the (possible) spike location.

⁹Zhenyu Liao and Romain Couillet. "The Dynamics of Learning: A Random Matrix Approach". In: *Proceedings of the 35th International Conference on Machine Learning*. Vol. 80. Proceedings of Machine Learning Research. Stockholm, Sweden: PMLR, 2018, pp. 3072–3081

Numerical results



Some further RMT efforts on high-dimensional dynamics

From the statistical physics community: reduces to low-dimensional ODE or SDE

- ▶ Sebastian Goldt et al. "Dynamics of Stochastic Gradient Descent for Two-Layer Neural Networks in the Teacher-Student Setup". In: *Advances in Neural Information Processing Systems*. Vol. 32. Curran Associates, Inc., 2019
- ▶ Francesca Mignacco et al. "Dynamical Mean-Field Theory for Stochastic Gradient Descent in Gaussian Mixture Classification". In: *Advances in Neural Information Processing Systems*. Vol. 33. Curran Associates, Inc., 2020, pp. 9540–9550
- ▶ Rodrigo Veiga et al. "Phase Diagram of Stochastic Gradient Descent in High-Dimensional Two-Layer Neural Networks". In: *Advances in Neural Information Processing Systems* 35 (Dec. 2022), pp. 23244–23255

From the optimization community: how RMT results apply to characterize average-case behavior in optimization

- ▶ Courtney Paquette et al. "SGD in the Large: Average-case Analysis, Asymptotics, and Stepsize Criticality". In: *Proceedings of Thirty Fourth Conference on Learning Theory*. PMLR, July 2021, pp. 3548–3626
- ▶ Courtney Paquette et al. "Halting Time Is Predictable for Large Models: A Universality Property and Average-Case Analysis". In: *Foundations of Computational Mathematics* 23.2 (Apr. 2023), pp. 597–673

And from the RMT community as well

- ▶ Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. "Online Stochastic Gradient Descent on Non-Convex Losses from High-Dimensional Inference". In: *Journal of Machine Learning Research* 22.106 (2021), pp. 1–51
- ▶ Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. "High-Dimensional Limit Theorems for SGD: Effective Dynamics and Critical Scaling". In: *Advances in Neural Information Processing Systems* 35 (Dec. 2022), pp. 25349–25362
- ▶ Gerard Ben Arous et al. *High-Dimensional SGD Aligns with Emerging Outlier Eigenspaces*. Oct. 2023. arXiv: 2310.03010 [cs, math, stat]

One step gradient beyond random network

- ▶ extends to **wide** DNN model via NTK, see, e.g., Y. Du, Z. Ling, R. C. Qiu, Z. Liao, "High-dimensional Learning Dynamics of Deep Neural Nets in the Neural Tangent Regime", High-dimensional Learning Dynamics Workshop, The Fortieth International Conference on Machine Learning (ICML'2023), 2023
- ▶ however, limited in the NTK and **linearized** regime
- ▶ what about **nonlinear feature learning** during gradient descent (**different** from initialization)?
- ▶ **empirical observation:** spikes appear in the NTK spectra during gradient descent training [FW20]

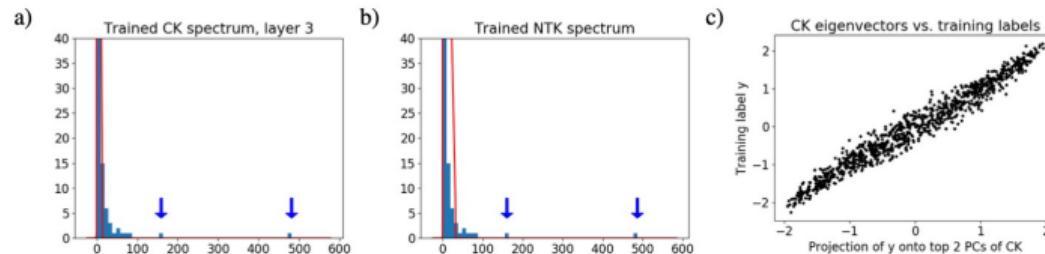
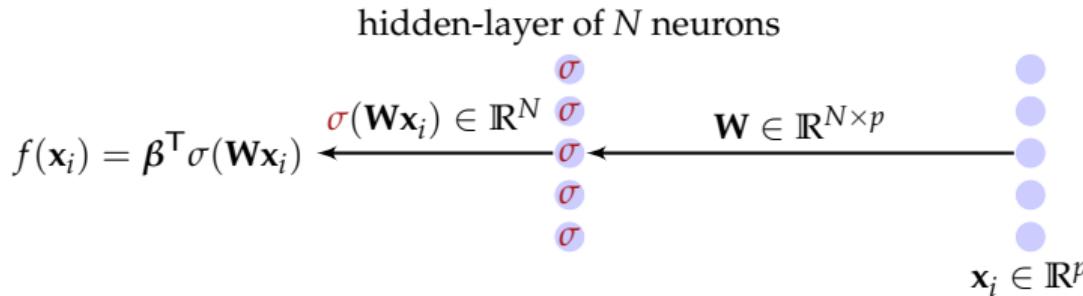


Figure 3: Eigenvalues of (a) K^{CK} and (b) K^{NTK} in a *trained* network, for training labels $y_\alpha = \sigma(\mathbf{x}_\alpha^\top \mathbf{v})$. The limit spectra at random initialization of weights are shown in red. Large outlier eigenvalues, indicated by blue arrows, emerge over training. (c) The projection of training labels onto the first 2 eigenvectors of the trained matrix K^{CK} accounts for 96% of the training label variance.

Interesting connection between optimization and RMT

Two-layer random network after one step training



- ▶ two-layer NN having N neurons, with output $f(\mathbf{x}) = \frac{1}{\sqrt{N}}\beta^T \sigma(\mathbf{W}\mathbf{x})$, for input $\mathbf{x} \in \mathbb{R}^p$, first-layer weight $\mathbf{W} \in \mathbb{R}^{N \times p}$, second-layer weight $\beta \in \mathbb{R}^N$, and nonlinear σ
- ▶ model trained on $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ of size n , by minimizing

$$\text{Cost} = \frac{1}{2n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2. \quad (19)$$

- ▶ first-layer gradient explicitly given by

$$\frac{\partial \text{Cost}}{\partial \mathbf{W}} = -\frac{1}{n} \left(\left(\frac{1}{\sqrt{N}} \beta \left(\mathbf{y}^T - \frac{1}{\sqrt{N}} \beta^T \sigma'(\mathbf{W}\mathbf{X}) \right) \right) \odot \sigma'(\mathbf{W}\mathbf{X}) \right) \mathbf{X}^T \in \mathbb{R}^{N \times p}, \quad (20)$$

with $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$, and $\mathbf{y} = [y_1, \dots, y_n]^T \in \mathbb{R}^n$.

Two-layer random network after one step GD training

- ▶ consider first step gradient update on \mathbf{W} as $\mathbf{W}_1 = \mathbf{W}_0 + \sqrt{N}\eta_0 \mathbf{G}_0$, with
$$\mathbf{G}_0 = \frac{1}{n} \left(\left(\frac{1}{\sqrt{N}} \beta_0 \left(\mathbf{y}^\top - \frac{1}{\sqrt{N}} \beta_0^\top \sigma(\mathbf{W}_0 \mathbf{X}) \right) \right) \odot \sigma'(\mathbf{W}_0 \mathbf{X}) \right) \mathbf{X}^\top$$
- ▶ **key observation** made in [Ba+22]: under standard assumption and for Gaussian \mathbf{W}_0, β_0 and \mathbf{X} , the first step gradient \mathbf{G}_0 is **approximately of rank one!**

$$\left\| \mathbf{G}_0 - \frac{\mathbb{E}[\sigma'(\xi)]}{n\sqrt{N}} \beta_0 \mathbf{y}^\top \mathbf{X}^\top \right\| \rightarrow 0. \quad (21)$$

- ▶ result obtained by (kind of conditioned on \mathbf{X}, \mathbf{y} and β_0) and playing with the randomness in \mathbf{W}_0
- ▶ built upon this, results on **generalization** can be obtained, etc.

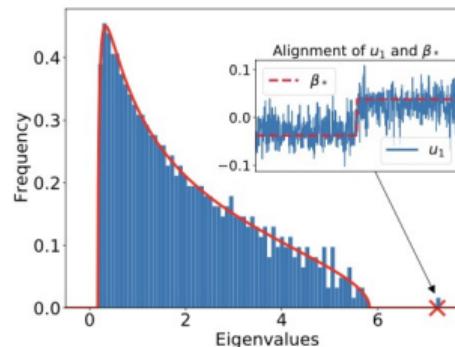
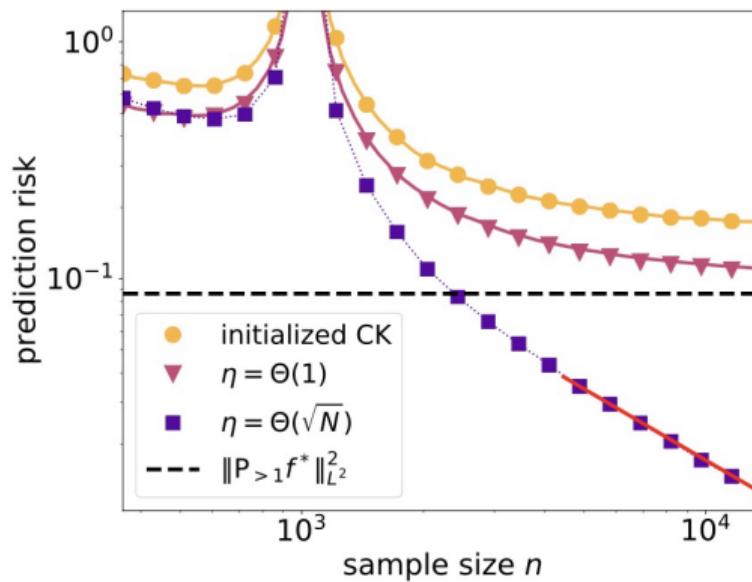


Figure 3: Main: empirical singular values of \mathbf{W}_1 (blue) vs. analytic prediction (red). Subfigure: overlap between u_1 and RMT4ML

Discussion on the step size and its impact

- ▶ since $\|\mathbf{W}_0\| = O(1)$, $\|\mathbf{W}_0\|_F = \sqrt{N}$, and $\sqrt{N}\|\mathbf{G}_0\| = O(1)$, $\sqrt{N}\|\mathbf{G}_0\|_F = O(1)$, may consider:
- ① small step $\eta = O(1)$ (same order in spectral norm): improve over initial CK, but **not as good** as optimal linear model
- ② large step $\eta = O(\sqrt{N})$ (same order in Frobenius norm): improve over a class of **nonlinear** model, match **neural scaling law** in some cases



Conclusion and take-away message

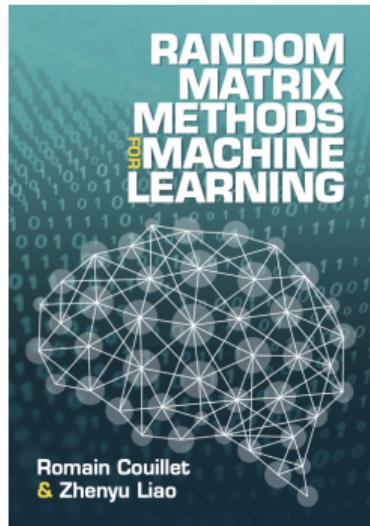
Take-away message:

- ▶ basics in ML and DL
- ▶ DL theory: **optimization+generalization**
- ▶ some **not so fantastic** story on neural tangent kernel and double descent
- ▶ opportunities in RMT for DL:
 - ① from shallow to deep random NNs
 - ② from random to non-so-random NNs
- ▶ what is a good theory for DNN?
- ▶ **Model and data/task dependent**, can be used to **guide DNN model design.**

RMT for machine learning: from theory to practice!

Random matrix theory (RMT) for machine learning:

- ▶ **change of intuition** from small to large dimensional learning paradigm!
- ▶ **better understanding** of existing methods: why they work if they do, and what the issue is if they do not
- ▶ **improved novel methods** with performance guarantee!



- ▶ book "*Random Matrix Methods for Machine Learning*"
- ▶ by Romain Couillet and **Zhenyu Liao**
- ▶ Cambridge University Press, 2022
- ▶ a pre-production version of the book and exercise solutions at <https://zhenyu-liao.github.io/book/>
- ▶ MATLAB and Python codes to reproduce all figures at <https://github.com/Zhenyu-LIAO/RMT4ML>

Thank you! Q & A?

Study of CK in the infinite-neuron regime

- ▶ **Key object:** empirical CK $\frac{1}{N}\Sigma^T\Sigma$, correlation in the feature space, for random initialization:
 $\mathbf{W}_{ij} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$, relates to linearized model f_{lin}
- ▶ $\frac{1}{N}\Sigma^T\Sigma = \frac{1}{N} \sum_{i=1}^N \sigma(\mathbf{X}^T \mathbf{w}_i) \sigma(\mathbf{w}_i^T \mathbf{X})$ for independent $\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.
- ▶ In the **infinite-neuron limit** ($N \rightarrow \infty$), convergence to the **limiting** CK matrix

$$\boxed{\frac{1}{N}\Sigma^T\Sigma \rightarrow \mathbf{K}_{\text{CK}}(\mathbf{X}) \equiv \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)} [\sigma(\mathbf{X}^T \mathbf{w}) \sigma(\mathbf{w}^T \mathbf{X})] \in \mathbb{R}^{n \times n}}$$

- ▶ theoretical **understanding** of NN model: generalization? optimization?
- ▶ **Application:** compress NN by carefully choosing **weights W** and/or **activation? σ** , e.g., **without** changing \mathbf{K}_{CK} ?

Problem settings

Data: K-class Gaussian mixture model (GMM)

Let $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ be independently drawn (non-necessarily uniformly) from one of the K classes:

$$\mathcal{C}_a : \sqrt{p}\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}_a, \mathbf{C}_a), \quad a \in \{1, \dots, K\} \quad (22)$$

Large dimensional asymptotics

As $n, p \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$ and some additional growth-rate assumptions on the difference $\|\boldsymbol{\mu}_a - \boldsymbol{\mu}_b\|$ and $\|\mathbf{C}_a - \mathbf{C}_b\|$, $a, b \in \{1, \dots, K\}$, as $n, p \rightarrow \infty$.

Theorem (Asymptotic equivalent for \mathbf{K} , [AZC22])

For CK matrix $\mathbf{K}_{CK} = \{\mathbb{E}[\sigma(\mathbf{x}_i^\top \mathbf{w})\sigma(\mathbf{w}^\top \mathbf{x}_j)]\}_{i,j=1}^n$ defined above, one has, as $n, p \rightarrow \infty$ that $\|\mathbf{K}_{CK} - \tilde{\mathbf{K}}_{CK}\| \rightarrow 0$, for some random matrix $\tilde{\mathbf{K}}_{CK}$ dependent of data \mathbf{X} , of activation σ but **only** via the following scalars

$$d_0 = \mathbb{E}[\sigma^2(\sqrt{\tau}z)] - \mathbb{E}[\sigma(\sqrt{\tau}z)]^2 - \tau\mathbb{E}[\sigma'(\sqrt{\tau}z)]^2, \quad d_1 = \mathbb{E}[\sigma'(\sqrt{\tau}z)]^2, \quad d_2 = \frac{1}{4}\mathbb{E}[\sigma''(\sqrt{\tau}z)]^2$$

and **independent** of the distribution of \mathbf{W} , as long as of normalized to have zero mean and unit variance.

Main result and the proof

Theorem (Asymptotic equivalent for \mathbf{K} , [AZC22])

For CK matrix $\mathbf{K}_{CK} = \{\mathbb{E}[\sigma(\mathbf{x}_i^\top \mathbf{w})\sigma(\mathbf{w}^\top \mathbf{x}_j)]\}_{i,j=1}^n$ defined above, one has, as $n, p \rightarrow \infty$ that $\|\mathbf{K}_{CK} - \tilde{\mathbf{K}}_{CK}\| \rightarrow 0$, for some random matrix $\tilde{\mathbf{K}}_{CK}$ dependent of data \mathbf{X} , of activation σ but **only** via the following scalars

$$d_0 = \mathbb{E}[\sigma^2(\sqrt{\tau}z)] - \mathbb{E}[\sigma(\sqrt{\tau}z)]^2 - \tau \mathbb{E}[\sigma'(\sqrt{\tau}z)]^2, \quad d_1 = \mathbb{E}[\sigma'(\sqrt{\tau}z)]^2, \quad d_2 = \frac{1}{4} \mathbb{E}[\sigma''(\sqrt{\tau}z)]^2$$

and **independent** of the distribution of \mathbf{W} , as long as of normalized to have zero mean and unit variance.

Proof sketch:

- ▶ We are interested in the kernel matrix \mathbf{K} , the (i, j) entry of which $\mathbf{K}_{ij} = \mathbb{E}_{\mathbf{w}}[\sigma(\mathbf{x}_i^\top \mathbf{w})\sigma(\mathbf{w}^\top \mathbf{x}_j)]$.
- ▶ Conditioned on $\mathbf{x}_i, \mathbf{x}_j, \mathbf{w}^\top \mathbf{x}_i \equiv \|\mathbf{x}_i\| \cdot \xi_i$ and $\mathbf{w}^\top \mathbf{x}_j$ are asymptotically **Gaussian**, but **correlated**!
- ▶ Gram-Schmidt to **de-correlate** $\mathbf{w}^\top \mathbf{x}_j = \frac{\mathbf{x}_i^\top \mathbf{x}_j}{\|\mathbf{x}_i\|} \xi_i + \sqrt{\|\mathbf{x}_j\|^2 - \frac{(\mathbf{x}_i^\top \mathbf{x}_j)^2}{\|\mathbf{x}_i\|^2}} \xi_j$, for Gaussian ξ_j now **independent** of ξ_i
- ▶ Use the fact $\mathbf{x}_i^\top \mathbf{x}_j = O(p^{-1/2})$ and $\|\mathbf{x}_i\|^2 \approx \tau/2 = O(1)$, Taylor-expand to “**linearize**” $\sigma(\cdot)$ to order $o(n^{-1})$
- ▶ Since $\|\mathbf{A}\|_2 \leq n \|\mathbf{A}\|_\infty$, with $\|\mathbf{A}\|_\infty = \max_{ij} |\mathbf{A}_{ij}|$, obtain **spectral approximation** $\tilde{\mathbf{K}}$.

¹⁰Hafiz Tiomoko Ali, Zhenyu Liao, and Romain Couillet. “Random matrices in service of ML footprint: ternary random features with no performance loss”. In: International Conference on Learning Representations (ICLR 2022). 2022

Practical consequence of the theory

According to theorem, allowed to choose **arbitrary** weights \mathbf{W} and activation σ , without affecting \mathbf{K} asymptotically, under the following conditions:

- ▶ weights \mathbf{W} have **independent** entries with zero mean and unit variance
- ▶ activation σ has the **same** few parameters as the original net

$$d_0 = \mathbb{E}[\sigma^2(\sqrt{\tau}z)] - \mathbb{E}[\sigma(\sqrt{\tau}z)]^2 - \tau\mathbb{E}[\sigma'(\sqrt{\tau}z)]^2, \quad d_1 = \mathbb{E}[\sigma'(\sqrt{\tau}z)]^2, \quad d_2 = \frac{1}{4}\mathbb{E}[\sigma''(\sqrt{\tau}z)]^2, \quad (23)$$

In particular,

- ▶ **sparse and binarized** (e.g., Bernoulli distributed) weights \mathbf{W} instead of dense Gaussian weights

$$[\mathbf{W}]_{ij} = 0 \text{ with proba } \varepsilon \in [0, 1), \quad [\mathbf{W}]_{ij} = \pm(1 - \varepsilon)^{-1/2} \text{ each with proba } 1/2 - \varepsilon/2, \quad (24)$$

- ▶ **sparse quantized** (e.g., binarized) activation σ shares the same d_0, d_1 , and d_2

Numerical results

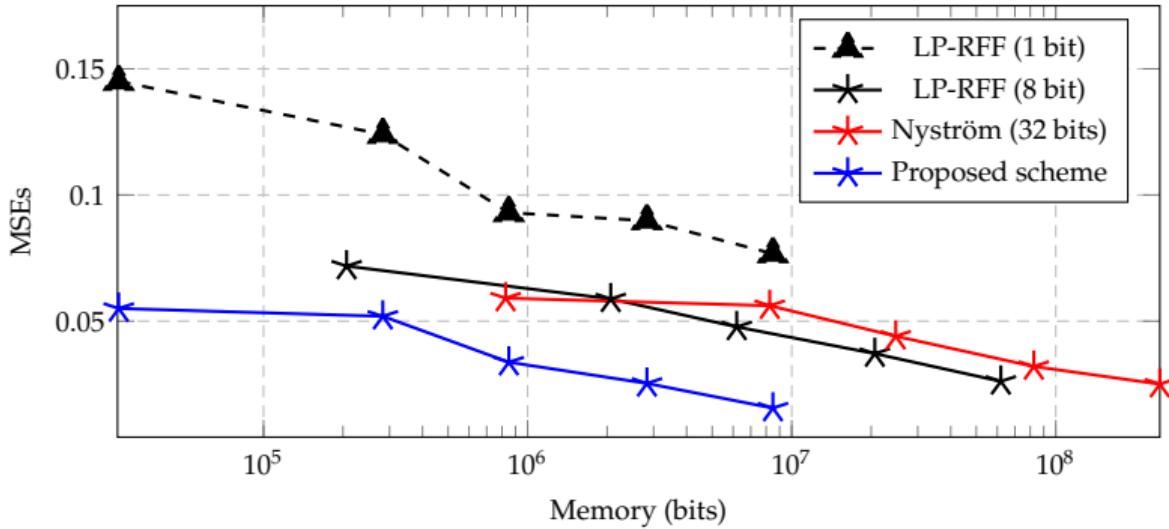


Figure: Test mean square errors of ridge regression on quantized single-hidden-layer random nets for different numbers of features $N \in \{5.10^2, 10^3, 5.10^3, 10^4, 5.10^4\}$, using LP-RFF, Nyström approximation, versus the proposed approach, on the Census dataset, with $n = 16\,000$ training samples, $n_{\text{test}} = 2\,000$ test samples, and data dimension $p = 119$.

CK of fully-connected deep neural networks

- ▶ everyone cares more about **deep** neural networks
- ▶ with some additional efforts, theory extends to fully-connected **deep** neural networks of depth L ,

$$f(\mathbf{x}) = \frac{1}{\sqrt{d_L}} \mathbf{w}^\top \sigma_L \left(\frac{1}{\sqrt{d_{L-1}}} \mathbf{W}_L \sigma_{L-1} \left(\dots \frac{1}{\sqrt{d_2}} \sigma_2 \left(\frac{1}{\sqrt{d_1}} \mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}) \right) \right) \right), \quad (25)$$

again for random $\mathbf{W}_1, \dots, \mathbf{W}_L$ and activations $\sigma_1(\cdot), \dots, \sigma_L(\cdot)$.

Theorem (Asymptotic equivalents for conjugate kernels, informal)

Under the same condition, define output features of layer $\ell \in \{1, \dots, L\}$, as

$$\Sigma_\ell = \frac{1}{\sqrt{d_\ell}} \sigma_\ell \left(\frac{1}{\sqrt{d_{\ell-1}}} \mathbf{W}_\ell \sigma_{\ell-1} \left(\dots \frac{1}{\sqrt{d_2}} \sigma_2 \left(\frac{1}{\sqrt{d_1}} \mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}) \right) \right) \right). \quad (26)$$

we have for the Conjugate Kernel $\mathbf{K}_{\text{CK},\ell}$ at layer ℓ defined as

$$\mathbf{K}_{\text{CK},\ell} = \mathbb{E}[\Sigma_\ell^\top \Sigma_\ell] \in \mathbb{R}^{n \times n}, \quad (27)$$

that $\|\mathbf{K}_{\text{CK},\ell} - \tilde{\mathbf{K}}_{\text{CK},\ell}\| \rightarrow 0$, some random matrix $\tilde{\mathbf{K}}_{\text{CK},\ell}$ dependent of data, of activation σ_ℓ but **only** via a few parameters, and **independent** of the distribution of \mathbf{W} , as long as of normalized to have zero mean and unit variance.

Theorem (Asymptotic equivalents for CK matrices, formal, [Gu+22])

Let $\tau_0, \tau_1, \dots, \tau_L \geq 0$ be a sequence of non-negative numbers satisfying the following recursion:

$$\tau_\ell = \sqrt{\mathbb{E}[\sigma_\ell^2(\tau_{\ell-1}\xi)]}, \quad \xi \sim \mathcal{N}(0, 1), \quad \ell \in \{1, \dots, L\}. \quad (28)$$

Further assume that the activation functions $\sigma_\ell(\cdot)$ s are “centered,” such that $\mathbb{E}[\sigma_\ell(\tau_{\ell-1}\xi)] = 0$. Then, for the CK matrix $\mathbf{K}_{\text{CK},\ell}$ of layer $\ell \in \{1, \dots, L\}$ defined in (27), as $n, p \rightarrow \infty$, one has that:

$$\|\mathbf{K}_{\text{CK},\ell} - \tilde{\mathbf{K}}_{\text{CK},\ell}\| \rightarrow 0, \quad \tilde{\mathbf{K}}_{\text{CK},\ell} \equiv \alpha_{\ell,1} \mathbf{X}^\top \mathbf{X} + \mathbf{V} \mathbf{A}_\ell \mathbf{V}^\top + (\tau_\ell^2 - \tau_0^2 \alpha_{\ell,1}) \mathbf{I}_n, \quad (29)$$

almost surely, with $\mathbf{V} = [\mathbf{J}/\sqrt{p}, \boldsymbol{\psi}] \in \mathbb{R}^{n \times (K+1)}$, $\mathbf{A}_\ell = \begin{bmatrix} \alpha_{\ell,2} \mathbf{t} \mathbf{t}^\top + \alpha_{\ell,3} \mathbf{T} & \alpha_{\ell,2} \mathbf{t} \\ \alpha_{\ell,2} \mathbf{t}^\top & \alpha_{\ell,2} \end{bmatrix} \in \mathbb{R}^{(K+1) \times (K+1)}$, for class label vectors $\mathbf{J} = [\mathbf{j}_1, \dots, \mathbf{j}_K] \in \mathbb{R}^{n \times K}$, “second-order” data fluctuation vector $\boldsymbol{\psi} \in \mathbb{R}^n$, second-order data statistics $\mathbf{t} = \{\text{tr } \mathbf{C}_a^\circ / \sqrt{p}\}_{a=1}^K \in \mathbb{R}^K$ and $\mathbf{T} = \{\text{tr } \mathbf{C}_a \mathbf{C}_b / p\}_{a,b=1}^K \in \mathbb{R}^{K \times K}$, as well as non-negative $\alpha_{\ell,1}, \alpha_{\ell,2}, \alpha_{\ell,3}$ satisfying

$$\alpha_{\ell,1} = \mathbb{E}[\sigma'_\ell(\tau_{\ell-1}\xi)]^2 \alpha_{\ell-1,1}, \quad \alpha_{\ell,2} = \mathbb{E}[\sigma'_\ell(\tau_{\ell-1}\xi)]^2 \alpha_{\ell-1,2} + \frac{1}{4} \mathbb{E}[\sigma''_\ell(\tau_{\ell-1}\xi)]^2 \alpha_{\ell-1,4}^2, \quad (30)$$

$$\alpha_{\ell,3} = \mathbb{E}[\sigma'_\ell(\tau_{\ell-1}\xi)]^2 \alpha_{\ell-1,3} + \frac{1}{2} \mathbb{E}[\sigma''_\ell(\tau_{\ell-1}\xi)]^2 \alpha_{\ell-1,1}^2. \quad (31)$$

with $\alpha_{\ell,4} = \mathbb{E}[(\sigma'_\ell(\tau_{\ell-1}\xi))^2 + \sigma_\ell(\tau_{\ell-1}\xi)\sigma''_\ell(\tau_{\ell-1}\xi)]$ $\alpha_{\ell-1,4}$ for $\xi \sim \mathcal{N}(0, 1)$.

Fully-connected deep nets: CK, NTK, and beyond

- ▶ happy with the study of (limiting) CK for DNN models
- ▶ extension to NTK via intrinsic **connection** between CK and NTK [JGH18]

$$\mathbf{K}_{\text{NTK},\ell}(\mathbf{X}) = \mathbf{K}_{\text{CK},\ell}(\mathbf{X}) + \mathbf{K}_{\text{NTK},\ell-1}(\mathbf{X}) \circ \mathbf{K}'_{\text{CK},\ell}(\mathbf{X}), \quad \mathbf{K}_{\text{NTK},0}(\mathbf{X}) = \mathbf{K}_{\text{CK},0}(\mathbf{X}) = \mathbf{X}^T \mathbf{X}, \quad (32)$$

and some additional efforts

- ▶ **convergence** and **generalization** theory via NTK [JGH18]: for **(i)** sufficiently wide nets **(ii)** trained with gradient descent of sufficiently small step size
- ▶ NTK is **determined** at random initialization and remains **unchanged** during training, and applies to **explicitly** characterize DNN convergence and generalization properties
- ▶ we can use the theory for DNN compression!

²Arthur Jacot, Franck Gabriel, and Clément Hongler. "Neural Tangent Kernel: Convergence and Generalization in Neural Networks". In: *Advances in Neural Information Processing Systems*. Vol. 31. NIPS'18. Curran Associates, Inc., 2018, pp. 8571–8580

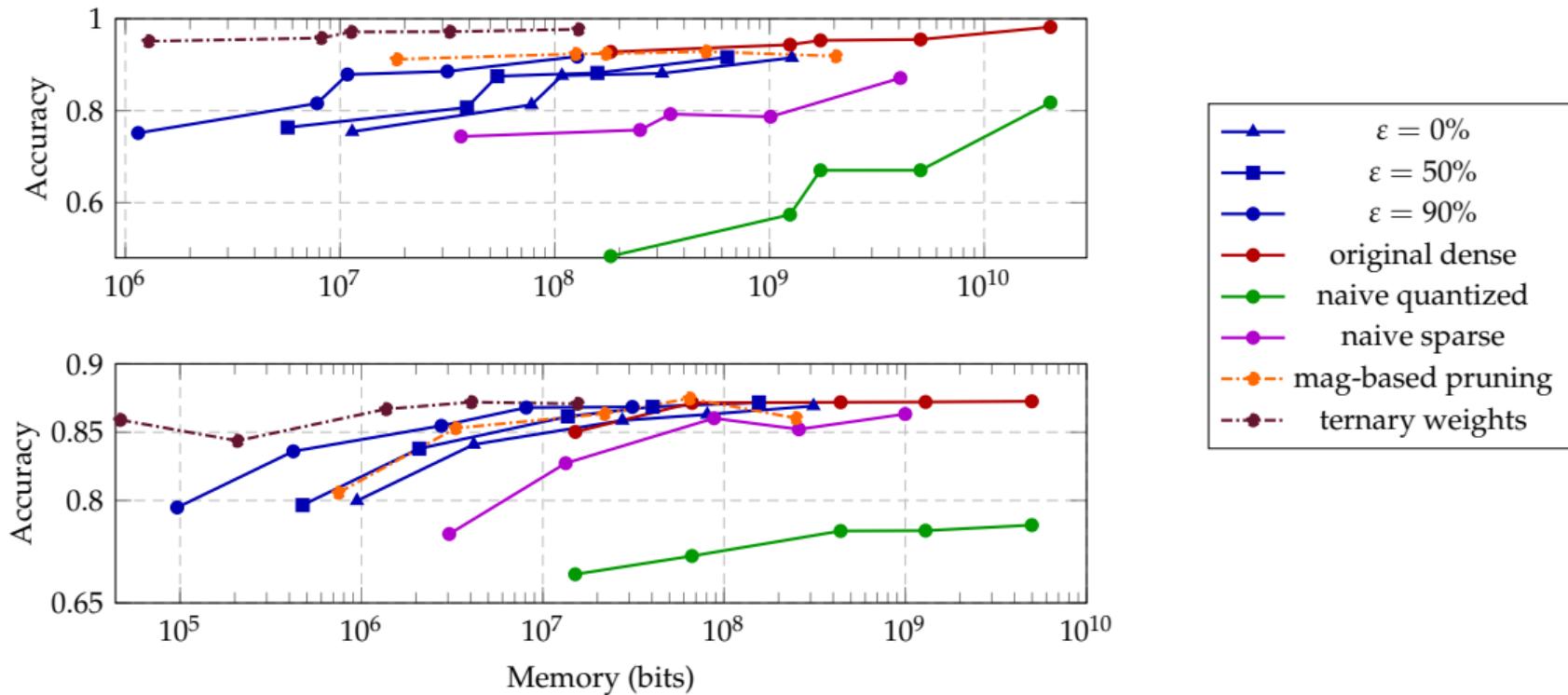


Figure: Test accuracy of classification on MNIST (**top**) and CIFAR10 (**bottom**) datasets. **Blue:** proposed NTK-LC approach with different levels of sparsity $\epsilon \in \{0\%, 50\%, 90\%\}$, **purple**: heuristic sparsification approach by uniformly zeroing out 80% of the weights, **green**: heuristic quantization approach with binary activation $\sigma(t) = 1_{t < -1} + 1_{t > 1}$, **red**: original network, **orange**: NTK-LC without activation quantization, and **brown**: magnitude-based pruning with same sparsity level as **orange**. Memory varies due to the **change of layer width** of the network.

Connection between Implicit and Explicit NNs

Deep equilibrium model, DEQ, [BKK19]

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

$$f(\mathbf{x}_i) = \boldsymbol{\beta}^T \mathbf{z}_i^*, \quad (33)$$

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

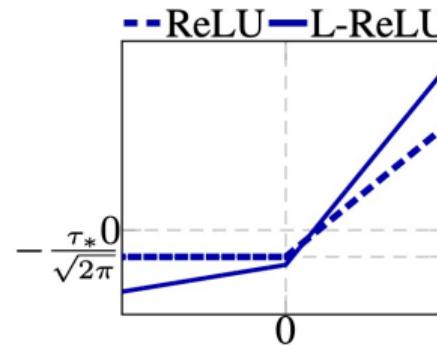
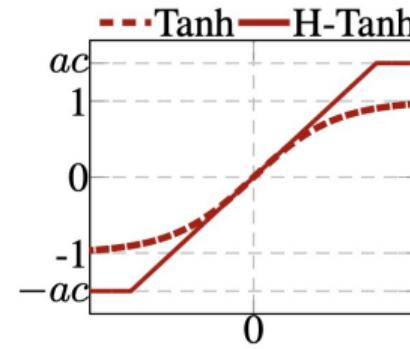
$$\mathbf{z}_i^{(l)} = \frac{1}{\sqrt{m}} \phi \left(\sigma_a \mathbf{A} \mathbf{z}_i^{(l-1)} + \sigma_b \mathbf{B} \mathbf{x}_i \right) \in \mathbb{R}^m, \text{ for } l \geq 1, \quad (34)$$

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

$$\mathbf{z}_i^* = \frac{1}{\sqrt{m}} \phi \left(\sigma_a \mathbf{A} \mathbf{z}_i^* + \sigma_b \mathbf{B} \mathbf{x}_i \right). \quad (35)$$

Connection between Implicit and Explicit NNs

- ▶ similar analysis can be performed for such **Implicit-NN** models as well
- ▶ leading to high-dimensional “**equivalence**” (in the sense of CK and/or NTK) between **Implicit** and **Explicit** NNs
- ▶ see for detailed results:
 - Zenan Ling et al. *Deep Equilibrium Models Are Almost Equivalent to Not-so-deep Explicit Models for High-dimensional Gaussian Mixtures*. May 2024. arXiv: 2402.02697 [cs, stat], to be presented at ICML 2024.



Numerical results

