Random Matrix Theory and Its Applications in ML: Part 2 Deep Learning Applications Short Course @ Jiangsu Normal University

Zhenyu Liao

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

January 12, 2024



Outline

In Introduction Deep Learning for Statisticians

2 Important Theoretical Questions for DL

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

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

④ Conclusion

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 as 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

Z. Liao (EIC, HUST)

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



Figure: Labeled data points $x \in \mathbb{R}^2$. Circles denote points in class C_1 . Crosses denote points in class C_2 .

- ▶ build a model/function *f* (from above historical data) that takes any points x ∈ ℝ² and returns C₁ or C₂
- ► logistic regression: $f(\mathbf{x}) = \sigma(\mathbf{w}^{\mathsf{T}}\mathbf{x} + b)$ 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 w, b from data/samples, by minimizing some cost function (e.g., negative likelihood, MSE)
- ▶ predict $\mathbf{x} \in C_1$ if $f(\mathbf{x}) < 1/2$ and $\mathbf{x} \in C_2$ otherwise.

Z. Liao (EIC, HUST)

RMT4ML

Neural networks are nothing but "cascaded" logistic regressors

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

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

•
$$\sigma(\mathbf{W}_2 x + \mathbf{b}_2) \in \mathbb{R}^2$$
, $\sigma(\mathbf{W}_3 \sigma(\mathbf{W}_2 x + \mathbf{b}_2) + \mathbf{b}_3) \in \mathbb{R}^3$
• $f_{4L-NN}(\mathbf{x}) = \sigma(\mathbf{W}_4 \sigma(\mathbf{W}_3 \sigma(\mathbf{W}_2 \mathbf{x} + \mathbf{b}_2) + \mathbf{b}_3) + \mathbf{b}_4) \in \mathbb{R}^2$
Define the label/target output as

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.

(2)



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:

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

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*:

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

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

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

Matlab code to train a simple NN

XXXXXXXX DATA XXXXXXXXXXXXXXX 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); zeros(1,5)]; zeros(1,5)]; zeros(1,5) = zeros(1,5) = zeros(1,5) = zeros(1,5); zeros(1,5) = zeros(1,5)% Initialize weights and biases $W^2 = 0.5*randn(2,2)$; $W^3 = 0.5*randn(3,2)$; $W^4 = 0.5*randn(2,3)$; $b^2 = 0.5*randn(2,1)$; $b^3 = 0.5*randn(3,1)$; $b^4 = 0.5*randn(2,1)$; % Forward and Back propagate eta = 0.05;% learning rate % number of SG iterations Niter = 1e6: 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 $a^2 = activate(x, W^2, b^2); a^3 = activate(a^2, W^3, b^3); a^4 = activate(a^3, W^4, b^4);$ % Backward pass delta4 = a4.*(1-a4).*(a4-v(:,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 semilogv([1:1e4:Niter].savecost(1:1e4:Niter)) function costval = cost(W2,W3,W4,b2,b3,b4)costvec = zeros(10.1);for i = 1:10x = [x1(i); x2(i)];a2 = activate(x, W2, b2); a3 = activate(a2, W3, b3); a4 = activate(a3, W4, b4);costvec(i) = norm(v(:,i) - a4,2);end $costval = norm(costvec.2)^2$: and % of negted function

Matlab code to train a simple NN

```
function y = activate(x,W,b)
XACTIVATE Evaluates sigmoid function.
X
X is the input vector, y is the output vector
X W contains the weights, b contains the shifts
X
The ith component of y is activate((Wx+b)_i)
X where activate(z) = 1/(1+exp(-z))
```

```
y = 1./(1+exp(-(W*x+b)));
```



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

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

different choice of activation and/or cost function:

- rectified linear unit, or **ReLU**, activation: $\sigma(t) = \max(t, 0)$
- **cross-entropy** cost function:

$$-\sum_{i=1}^{N} \log \left(\frac{e^{v_{l_i}^{(i)}}}{\sum_{j=1}^{K} e^{v_j^{(i)}}} \right).$$
(6)

- dropout, batch normalization, and other types of normalization, etc.
- use of **tensors** instead of vectors or matrices for input data or intermediate representations

What does deep learning theory care about and why?

- **b** 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
- 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., $\mathbf{x} = [x_1, ..., x_p] \subset [0, 1]^p$, 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)?
 - nice property of the (over-parameterized) DL model: Neural Tangent Kernel [JGH18]
 - Inductive bias due to algorithm: Double Descent or Benign Overfitting [BMR21]



A Good DL theory should cover both optimization and generalization!

Z. Liao (EIC, HUST)

January 12, 2024

²Arthur Jacot, Franck Gabriel, and Clément Hongler. "Neural tangent kernel: Convergence and generalization in neural networks". In: Advances in neural information processing systems. 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 $K(\cdot, \cdot)$: $\mathbb{R}^p \times \mathbb{R}^p \to \mathbb{R}$, similarity measure between input data points in \mathbb{R}^p
- examples include:
 - linear kernel $K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^{\mathsf{T}} \mathbf{y}$, cosine kernel $= \frac{\mathbf{x}^{\mathsf{T}} \mathbf{y}}{\|\mathbf{x}\| \cdot \|\mathbf{y}\|}$, Gaussian (RBF) kernel $= \exp(\|\mathbf{x} \mathbf{y}\|^2 / \gamma^2)$
 - kernel induced by NN: $K(\mathbf{x}, \mathbf{y}) = \sigma(\mathbf{W}\mathbf{x})^{\mathsf{T}} \sigma(\mathbf{W}\mathbf{y})$, 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

hidden-layer of N neurons



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

$$f(\mathbf{x};\boldsymbol{\theta}) = \boldsymbol{\beta}^{\mathsf{T}} \sigma(\mathbf{W}\mathbf{x}) = \sum_{\ell=1}^{n} \beta_{\ell} \sigma(\mathbf{w}_{\ell}^{\mathsf{T}}\mathbf{x}), \quad \boldsymbol{\theta} = [\beta_{1}, \ldots, \beta_{N}; \mathbf{w}_{1}, \ldots, \mathbf{w}_{N}].$$
(7)

linearization of the network at initialization, by Taylor expansion

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

and

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

Z. Liao (EIC, HUST)

January 12, 2024

17/44

The big picture of NTK

• around initialization $\theta \approx \theta_0$, linearized network output

$$f(\mathbf{x};\boldsymbol{\theta}) \approx f_{\text{lin}}(\mathbf{x};\boldsymbol{\theta}_0 + \boldsymbol{\delta}) = f(\mathbf{x};\boldsymbol{\theta}_0) + \boldsymbol{\delta}^{\mathsf{T}} \boldsymbol{\phi}_{\text{NTK}}(\mathbf{x}), \quad K_{e-NTK}(\mathbf{x},\mathbf{y}) = \boldsymbol{\phi}_{\text{NTK}}(\mathbf{x})^{\mathsf{T}} \boldsymbol{\phi}_{\text{NTK}}(\mathbf{y}), \tag{10}$$

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

- 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
- **(a)** it suffices to optimize in $B(\theta_0)$ to reach an approx. global min, i.e., $f(\mathbf{x}; \theta_0) \approx f_{\text{lin}}(\mathbf{x}; \theta_0) \approx 0$
- from an optimization viewpoint, optimizing $f(\mathbf{x}; \boldsymbol{\theta}) \approx \text{optimizing } f_{\text{lin}}$ and will not leave $B(\boldsymbol{\theta}_0)$

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) $Cost(f_{\theta}(\mathbf{x}), \mathbf{y}) \approx Cost(f_{lin}(\mathbf{x}), \mathbf{y})$ linear (in the parameter θ) and convex!
- ► for MSE, $\operatorname{Cost}(f_{\operatorname{lin}}(\mathbf{X}), \mathbf{y}) = \frac{1}{n} \sum_{i=1}^{n} (f_{\operatorname{lin}}(\mathbf{x}_i) y_i)^2$, nothing but linear regression of type $\operatorname{Cost} = \|\mathbf{y}' \mathbf{\Phi}_{NTK}(\mathbf{X})^{\mathsf{T}} \boldsymbol{\delta}\|^2$ with $y'_i = f(\mathbf{x}_i; \boldsymbol{\theta}_0) y$

Precise Characterization of Double Descent Curves



- larger model, the better?! Maybe, due to double descent [Has+22] and implicit (norm-based?) bias
- ► case of linear regression model $\text{Cost} = \frac{1}{n} \sum_{i=1}^{n} (\beta^{\mathsf{T}} \mathbf{x}_{i} - y_{i})^{2}$, with β , $\mathbf{x}_{i} \in \mathbb{R}^{p}$, 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



Linear versus nonlinear

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 $\text{Cost} = \frac{1}{n} \sum_{i=1}^{n} (\beta^{\mathsf{T}} \mathbf{x}_i - y_i)^2$, 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 β̂ = (XX^T)⁻¹Xy or β̂ = X(X^TX)⁻¹y 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



- 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 β by minimizing Cost = $\frac{1}{n} \sum_{i=1}^{n} (y_i - \beta^T \sigma(\mathbf{W}\mathbf{x}_i))^2 + \gamma \|\beta\|^2$ for some regularization parameter $\gamma > 0$, then

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

training MSE (on the given training set (X, y)) reads

$$E_{\text{train}} = \frac{1}{n} \| \mathbf{y} - \boldsymbol{\Sigma}^{\mathsf{T}} \boldsymbol{\beta} \|_{F}^{2} = \frac{\gamma^{2}}{n} \mathbf{y} \mathbf{Q}^{2}(\gamma) \mathbf{y}, \quad \mathbf{Q}(\gamma) \equiv \left(\frac{1}{n} \boldsymbol{\Sigma}^{\mathsf{T}} \boldsymbol{\Sigma} + \gamma \mathbf{I}_{n}\right)^{-1}$$
(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} : $E_{\text{test}} = \frac{1}{\hat{n}} \|\hat{\mathbf{y}} - \hat{\boldsymbol{\Sigma}}^{\mathsf{T}} \boldsymbol{\beta}\|_{F}^{2}$, $\hat{\boldsymbol{\Sigma}} = \sigma(\mathbf{W}\hat{\mathbf{X}})$.

Nonlinear resolvent

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

▶ nonlinear $\Sigma^{\mathsf{T}} = \sigma(\mathbf{W}\mathbf{X})^{\mathsf{T}}$ still has i.i.d. columns, but

- ▶ its *i*-th column $\sigma([\mathbf{X}^{\mathsf{T}}\mathbf{W}^{\mathsf{T}}]_{i})$ 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

$$\mathbb{P}\left(\left|\frac{1}{n}\sigma(\mathbf{w}^{\mathsf{T}}\mathbf{X})\mathbf{A}\sigma(\mathbf{X}^{\mathsf{T}}\mathbf{w})-\frac{1}{n}\operatorname{tr}\mathbf{A}\mathbf{K}\right|>t\right)\leq Ce^{-cn\min(t,t^{2})}$$

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

- **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 **Q**, establish (limiting) eigenvalue distribution of $\frac{1}{n}\sigma(\mathbf{WX})^{\mathsf{T}}\sigma(\mathbf{WX})$, etc.

Z. Liao (EIC, HUST)

RMT4ML

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 $n, p, N \to \infty$ at the same pace, for $\mathbf{Q} = (\sigma(\mathbf{X}^{\mathsf{T}}\mathbf{W}^{\mathsf{T}})\sigma(\mathbf{W}\mathbf{X})/n + \gamma \mathbf{I}_n)^{-1}$ with $\gamma > 0$,

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

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

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 $n, p, N \to \infty$, we have $E_{\text{train}} - \bar{E}_{\text{train}} \to 0$ and $E_{\text{test}} - \bar{E}_{\text{test}} \to 0$ with

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

⁵Cosme Louart, Zhenyu Liao, and Romain Couillet. "A Random Matrix Approach to Neural Networks". In: *The Annals of Applied Probability* 28.2 (2018), pp. 1190–1248

Z. Liao (EIC, HUST)

Numerical results





Z. Liao (EIC, HUST)

Numerical results: double descent



Some further RMT investigations on the two-layer model

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

- ▶ [PW17] first guess expression of the eigenvalue behavior
- ► [BP21]: eigenvalue distribution of $\frac{1}{n}\sigma(\mathbf{W}\mathbf{X})^{\mathsf{T}}\sigma(\mathbf{W}\mathbf{X})$ for \mathbf{W}, \mathbf{X} having sub-gaussian entries
- for "centered" $\sigma(\cdot)$ with respect to Gaussian measure: $\mathbb{E}[\sigma(\xi)] = 0$ for $\xi \sim \mathcal{N}(0, 1)$
- **a** 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 $\frac{1}{n}\sigma(\mathbf{W}\mathbf{X})^{\mathsf{T}}\sigma(\mathbf{W}\mathbf{X})$ for sub-gaussian \mathbf{W}, \mathbf{X} and centered $\sigma(\cdot)$
- despite being a white model, spikes may appear!
- if $\mathbb{E}[\xi^2 \sigma(\xi)] = 0$, then no spike
- Otherwise, at most two spikes

Question: what happen if either W or 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. 2017, pp. 2634–2643

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: Advances in Neural Information Processing Systems. Vol. 33. Curran Associates, Inc., 2020, pp. 7710–7721

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 X = [x₁,..., x_n] ∈ ℝ^{p×n} with associated labels/targets y = [y₁,..., y_n] ∈ ℝⁿ, w ∈ ℝ^p is learned via gradient descent by minimizing the (ridge-regularized) squared loss

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

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

• gradient given by $\nabla L(\mathbf{w}) = -\frac{1}{n}\mathbf{X}(\mathbf{y} - \mathbf{X}^{\mathsf{T}}\mathbf{w}) + \gamma \mathbf{w}$ 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} :

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

solution explicitly given by

$$\mathbf{w}(t) = e^{-\alpha t \left(\frac{1}{n} \mathbf{X} \mathbf{X}^{\mathsf{T}} + \gamma \mathbf{I}_{p}\right)} \mathbf{w}_{0} + \left(\mathbf{I}_{p} - e^{-\alpha t \left(\frac{1}{n} \mathbf{X} \mathbf{X}^{\mathsf{T}} + \gamma \mathbf{I}_{p}\right)}\right) \mathbf{w}_{\infty}$$
(15)

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

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

the ridge regression solution with regularization parameter γ .

Z. Liao (EIC, HUST)

RMT4ML

32/44

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

$$C_1: \mathbf{x}_i \sim \mathcal{N}(-\boldsymbol{\mu}, \mathbf{I}_p) \quad C_2: \mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{I}_p)$$

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

study training and test misclassification error rates as

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

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 $\hat{y} = -1$. \blacktriangleright 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 $n, p \to \infty$ with $p/n \to c \in (0, \infty)$, we have

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

almost surely, where

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

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

with $\rho = \lim_{p \to \infty} \|\mu\|^2$, Γ a positive contour surrounding the support of the Marčenko–Pastur law (shifted by $\gamma \ge 0$) and the points $(\gamma, 0)$ and $(\gamma + \lambda_s, 0)$ with $\lambda_s = c + 1 + \rho + c/\rho$, $f_t(z) \equiv \exp(-\alpha tz)$ 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 (isf 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 $\rho = \lim \|\mu\|^2$, $f_t(x) = \exp(-\alpha tx)$, $\mu(x)$ the MP law

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

and

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

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

Z. Liao (EIC, HUST)

January 12, 2024

35 / 44

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. 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.

Two-layer random network after one step training



two-layer NN having *N* neurons, with output f(**x**) = 1/√N β^T σ(**Wx**), for input **x** ∈ ℝ^p, first-layer weight **W** ∈ ℝ^{N×p}, second-layer weight β ∈ ℝ^N, and nonlinear σ
 model trained on {(**x**_i, y_i)}ⁿ_{i=1} of size *n*, by minimizing

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

first-layer gradient explicitly given by

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

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

Two-layer random network after one step GD training

- consider first step gradient update on W as $W_1 = W_0 + \sqrt{N}\eta_0 G_0$, with $\mathbf{G}_{0} = \frac{1}{n} \left(\left(\frac{1}{\sqrt{N}} \boldsymbol{\beta}_{0} \left(\mathbf{y}^{\mathsf{T}} - \frac{1}{\sqrt{N}} \boldsymbol{\beta}_{0}^{\mathsf{T}} \sigma(\mathbf{W}_{0} \mathbf{X}) \right) \right) \odot \sigma'(\mathbf{W}_{0} \mathbf{X}) \mathbf{X}^{\mathsf{T}}$
- **•** key observation made in [Ba+22]: under standard assumption and for Gaussian W_0 , β_0 and X, the first step gradient G_0 is approximately of rank one!

$$\left\|\mathbf{G}_{0} - \frac{\mathbb{E}[\sigma'(\boldsymbol{\xi})]}{n\sqrt{N}}\boldsymbol{\beta}_{0}\mathbf{y}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\right\| \to 0.$$
(21)

 \triangleright result obtained by (kind of conditioned on **X**, **y** and β_0) and playing with the randomness in **W**₀ built upon this, results on generalization can be obtained, etc.



Figure 3: Main: empirical singular values of W_1 (blue) vs. analytic prediction (red) Subfigure: overlap between u, 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:
- small step $\eta = O(1)$ (same order in spectral norm): improve over initial CK, but not as good as optimal linear model
- large step η = O(√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:
- from shallow to deep random NNs
- I from random to non-so-random NNs
- what is a good theory for DNN?
- Model and data/task dependent, can be used to guild 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?