## Practice, Theory, and Theorems for Random Matrix Theory in Modern Machine Learning

#### Michael. W. Mahoney

#### ICSI and Department of Statistics University of California, Berkeley, USA

(Joint work with Charles Martin and Zhenyu Liao)

June 24, 2022

## Outline



## 2 Theory

## More practice

## Theorems

- Introduction
- Sample covariance matrix for large dimensional data: from concentration to RMT
- A random matrix perspective of the "curse of dimensionality"
- Kernel spectral clustering for large dimensional data
- A random matrix approach to large neural networks and random features

## Table of Contents

## Practice

## 2 Theory

## 3 More practice

## Interval 1 Theorems

- Introduction
- Sample covariance matrix for large dimensional data: from concentration to RMT
- A random matrix perspective of the "curse of dimensionality"
- Kernel spectral clustering for large dimensional data
- A random matrix approach to large neural networks and random features

# *Lots* of DNNs analyzed: Look at nearly every publicly-available SOTA model in CV and NLP

- Don't evaluate your method on one/two/three NNs, evaluate it on:
  - dozens (2017)
  - hundreds (2019)
  - thousands (2021)
- Don't use bad/toy models, use SOTA models.
  - If you do, don't be surprised if low-quality/toy models are different than high-quality/SOTA models.
- Don't train models, instead validate pre-trained models.
  - Validating models is harder than training models.

# Results: LeNet5 (an old/small NN example)



Figure: Full and zoomed-in ESD for LeNet5, Layer FC1.

Older and/or smaller and/or less well-trained models look like bulk+spike.

Mahoney (	(UC Berkeley)
-----------	---------------

# Results: AlexNet (a typical modern/large DNN example)



Figure: Zoomed-in ESD for Layer FC1 and FC3 of AlexNet.

Newer SOTA models have heavy-tail structure in their weight matrix correlations (i.e., not elements but eigenvalues).

Mahoney (UC Berkeley)

WeightWatcher

## Table of Contents

Practice

## 2 Theory

## 3 More practice

## Intervention Theorems

Introduction

- Sample covariance matrix for large dimensional data: from concentration to RMT
- A random matrix perspective of the "curse of dimensionality"
- Kernel spectral clustering for large dimensional data
- A random matrix approach to large neural networks and random features

Random Matrix Theory 101: Wigner and Tracy-Widom

- Wigner: global bulk statistics approach universal semi-circular form
- Tracy-Widom: local edge statistics fluctuate in universal way



Problems with Wigner and Tracy-Widom:

- Weight matrices usually not square
- Typically do only a single training run

# Random Matrix Theory 102': Marchenko-Pastur



(c) Vary aspect ratios



< □ > < □ > < □ > < □ >

Figure: Marchenko-Pastur (MP) distributions.

Important points:

- Global bulk stats: The overall shape is deterministic, fixed by Q and  $\sigma$ .
- Local edge stats: The edge  $\lambda^+$  is very crisp, i.e.,  $\Delta \lambda_M = |\lambda_{max} - \lambda^+| \sim O(M^{-2/3})$ , plus Tracy-Widom fluctuations.

We use both global bulk statistics as well as local edge statistics in our theory.

# Random Matrix Theory 103: Heavy-tailed RMT

Go beyond the (relatively easy) Gaussian Universality class:

• model strongly-correlated systems ("signal") with heavy-tailed random matrices.

	Generative Model	Finite-N	Limiting	Bulk edge	(far) Tail
	w/ elements from	Global shape	Global shape	Local stats	Local stats
	Universality class	$\rho_N(\lambda)$	$\rho(\lambda), N \to \infty$	$\lambda \approx \lambda^+$	$\lambda \approx \lambda_{max}$
Basic MP	Gaussian	MP distribution	MP	TW	No tail.
Spiked- Covariance	Gaussian, + low-rank perturbations	MP + Gaussian spikes	MP	TW	Gaussian
Heavy tail, $4 < \mu$	(Weakly) Heavy-Tailed	MP + PL tail	MP	Heavy-Tailed*	Heavy-Tailed*
Heavy tail, $2 < \mu < 4$	(Moderately) Heavy-Tailed (or "fat tailed")	$\sim \lambda^{-(a\mu+b)}$	$\sim \lambda^{-(\frac{1}{2}\mu+1)}$	No edge.	Frechet
Heavy tail, $0 < \mu < 2$	(Very) Heavy-Tailed	$\sim \lambda^{-(\frac{1}{2}\mu+1)}$	$\sim \lambda^{-(\frac{1}{2}\mu+1)}$	No edge.	Frechet

Basic MP theory, and the spiked and Heavy-Tailed extensions we use, including known, empirically-observed, and conjectured relations between them. Boxes marked "\*" are best described as following "TW with large finite size corrections" that are likely Heavy-Tailed, leading to bulk edge statistics and far tail statistics that are indistinguishable. Boxes marked "\*" are phenomenological fits, describing large ( $2 < \mu < 4$ ) or small ( $0 < \mu < 2$ ) finite-size corrections on  $N \to \infty$  behavior.

# RMT-based 5+1 Phases of Training (in pictures)



Figure: The 5+1 phases of learning we identified in DNN training.

Mahoney (UC Berkeley)

▲ ■ ► ■ つへの April 2022 21/50

# Bulk+Spikes: Small Models $\sim$ Tikhonov regularization



Perturbative correction

$$egin{aligned} \lambda_{max} &= & \sigma^2 \left( rac{1}{Q} + rac{|\Delta|^2}{N} 
ight) \left( 1 + rac{N}{|\Delta|^2} 
ight) \ & & |\Delta| > (Q)^{-rac{1}{4}} \end{aligned}$$

simple scale threshold

$$\mathbf{x} = \left(\hat{\mathbf{X}} + lpha \mathbf{I}
ight)^{-1} \hat{\mathbf{W}}^{T} \mathbf{y}$$

eigenvalues  $> \alpha$  (Spikes) carry most of the signal/information

Smaller, older models like LeNet5 exhibit traditional regularization and can be described perturbatively with Gaussian RMT

Mahoney (UC Berkeley)

Low-rank perturbation

 $\mathbf{W}_{l} \simeq \mathbf{W}_{l}^{rand} + \Delta^{large}$ 

WeightWatcher

April 2022 22 / 50

# Heavy-tailed Self-regularization

 $\boldsymbol{\mathsf{W}}$  is strongly-correlated and highly non-random

- We model strongly-correlated systems by heavy-tailed random matrices
- We model signal (not noise) by heavy-tailed random matrices

Then RMT/MP ESD will also have heavy tails.

• The eigenvalues are heavy-tailed; the weights are NOT.



"All" larger, modern DNNs exhibit novel Heavy-tailed self-regularization

Mahoney (	UC Berkeley
-----------	-------------

## Table of Contents

Practice

## 2 Theory

## 3 More practice

### 4 Theorems

Introduction

- Sample covariance matrix for large dimensional data: from concentration to RMT
- A random matrix perspective of the "curse of dimensionality"
- Kernel spectral clustering for large dimensional data
- A random matrix approach to large neural networks and random features

# Watching weights with WeightWatcher

https://github.com/CalculatedContent/WeightWatcher

## Analyzing DNN Weight matrices with WeightWatcher



Compare multiple layers of pre-trained model

Monitor NN properties as you train your own model

#### "pip install weightwatcher"

# Using the theory

Different ways one could *use* a theory.

- Perform diagnostics for model validation, to develop hypotheses, etc.\*
- Make predictions about model quality, generalization, transferability, etc.\*
- Did post-training modifications damage my model?\*
- Will buying more data help?\*
- Will training longer help?\*
- Will quantizing or distilling help?\*
- Construct a regularizer to do model training.\*\*

\*Ideally, by peeking at very little or no data.

\*\*If you have lots of data, lots of GPUs, etc.

# Predicting test accuracies ... lots of metrics ...

• Average log norm (a VC-like data-dependent capacity metric):

$$\langle \log \| \mathbf{W} \| 
angle = rac{1}{N} \sum_{l,i} \log \| \mathbf{W}_{l,i} \| = rac{1}{N} \sum_{l,i} \log(\lambda_{l,i}^{max})$$

• Average alpha (also data-dependent, from HT-SR theory):

$$\alpha = \frac{1}{N} \sum_{I,i} \alpha_{I,i}$$

• Combine the two into a weighted average (weighted to compensate for different size and scale of feature maps):

$$\hat{\alpha} = \frac{1}{N} \sum_{l,i} \log(\lambda_{l,i}^{max}) \alpha_{l,i}$$

• In a special case ( $\alpha \approx 2$ ), for each layer:

**PL–Norm Relation:**  $\alpha \log \lambda^{max} \approx \log \|\mathbf{W}\|_{F}^{2}$ .

#### "pip install weightwatcher"

# (The first) large-scale study (meta-analysis) of hundreds of SOTA pretrained models $^\ddagger$



Series	#	Metric	$(\log   \mathbf{W}  _F^2)$	$(\log \ \mathbf{W}\ _{\infty}^2)$	â	$(\log \ \mathbf{X}\ _{a}^{\alpha})$
		RMSE	0.56	0.23	0.48	0.34
VGG	6	$R^2$	0.88	0.98	0.92	0.96
		Kendall- $\tau$	-0.79	-0.93	-0.93	-0.93
		RMSE	0.9	0.97	0.61	0.66
ResNet	5	$R^2$	0.92	0.9	0.96	0.9
		Kendall- $\tau$	-1.0	-1.0	-1.0	-1.0
DerWet		RMSE	2.4	2.8	1.8	1.9
nesivet-	19	$R^2$	0.81	0.74	0.89	0.88
IK		Kendall- $\tau$	-0.79	-0.79	-0.89	-0.88
		RMSE	0.3	0.11	0.16	0.21
DenseNet	4	$R^2$	0.93	0.99	0.98	0.97
		Kendall- $\tau$	-1.0	-1.0	-1.0	-1.0

Table 1: Quality metrics (for RMSE, smaller is better; for R<sup>2</sup>, larger is better; and for Kendallrank correlation, larger magnitude is better) for reported Top1 test error for pretrained models in each architecture series. Column # refers to number of models. VGG, ResNet, and DenseNet were pretrained on ImageNet. ResNet.1K was protrained on ImageNet. ResNet.3

#### Summary statistics: VGG; ResNet; DenseNet.

	$\log  \cdot _F^2$	$\log  \cdot _{\infty}^2$	â	$\log \  \cdot \ _{\alpha}^{\alpha}$
RMSE (mean)	4.84	5.57	4.58	4.55
RMSE (std)	9.14	9.16	9.16	9.17
R2 (mean)	3.9	3.85	3.89	3.89
R2 (std)	9.34	9.36	9.34	9.34
Kendal-tau (mean)	3.84	3.77	3.86	3.85
Kendal-tau (std)	9.37	9.4	9.36	9.36

Table 3: Comparison of linear regression fits for different average Log Norm and Weighted Alpha metrics across 5 CV datasets, 17 architectures, covering 108 (out of over 400) different pretrained

Figure 2: Comparison of Average Log Norm and Weighted Alpha quality metrics versus re metrics across 5 CV datasets, 17 architectures, covering 108 (out of over 400) different pretrained test accuracy for pretrained VGG models: VGG11, VGG13, VGG16, and VGG19, with and

#### Different metrics on pre-trained VGG.

Summary statistic	s: hundreds	of	models.
-------------------	-------------	----	---------

A B A B A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 A
 A
 A
 A

Lots more plots to prove we can "predict trends ... without access ....'

<sup>1</sup> "Predicting trends in the quality of state-of-the-art neural networks without access to training or testing data," Martin,

Peng, and Mahoney, arXiv:2002.06716, Nature Communications, 2021.

Mahoney (UC Berkeley)

WeightWatcher

April 2022 29 / 50

# Using a theory: on SOTA models

# Analyzing pre-trained models: properties of VGG vs ResNet vs DenseNet leads to the idea of *correlation flow*.



Figure 4: PL exponent ( $\alpha$ ) versus layer id, for the least and the most accurate models in VGG (a), ResNet (b), and DenseNet (c) series. (VGG is without BN; and note that the Y axes on

Alpha versus depth: VGG, ResNet, DenseNet.

WeightWatcher

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

## Using a theory: on SOTA models

Analyzing pre-trained models: properties of GPTx series leads to the idea of *scale collapse*.



Figure 6: Histogram of PL exponents and Log Spectral Norms for weight matrices from the OpenAI GPT and GPT2-small pretrained models.



Figure 7: Log Spectral Norms (in (a)) and PL exponents (in (b)) for weight matrices from the OpenAl GPT and GPT2-small pretrained models. (Note that the quantities shown on each Y axis are different.) In the text, this is interpreted in terms of Scale Collapse and Correlation Flow.

Histogram and depth plots of  $\alpha_{l,i}$  and  $\lambda_{l,i}^{max}$ .

Mahoney (UC Berkeley)

< □ > < □ > < □ > < □ > < □ > < □ >

## Using a theory: easy to break popular SLT metrics

Easy to "break" popular SLT metrics:

- they are not validated counterfactually
- (but they drive the development of models)



Figure 5: ResNet20, distilled with Group Regularization, as implemented in the distiller (4D.regularized.5Lremoved) pretrained models. Log Spectral Norm  $(\log \lambda_{max})$  and PL exponent ( $\alpha$ ) for individual layers, versus layer id, for both baseline (before distillation, green) and finetuned (after distillation, red) pretrained models.

#### Intel's distillation "broke" their models.

Series	#	$\langle \log \  \mathbf{W} \ _F \rangle$	$(\log    W   _{\infty})$	â	$(\log \ \mathbf{X}\ _{\alpha}^{\alpha})$
GPT	49	1.64	1.72	7.01	7.28
GPT2-small	49	2.04	2.54	9.62	9.87
GPT2-medium	98	2.08	2.58	9.74	10.01
GPT2-large	146	1.85	1.99	7.67	7.94
GPT2-xl	194	1.86	1.92	7.17	7.51

Table 2: Average value for the average Log Norm and Weighted Alpha metrics for pretrained OpenAI GPT and GPT2 models. Column # refers to number of layers treated. Averages do

GPTx series: how does a model trained to "bad" data differ from one trained to "good" data?

< ロ > < 同 > < 回 > < 回 >

# Using a theory: leads to predictions

Based on analyzing hundreds of pre-trained SOTA models:

## • "Correlation flow":

 "Shape" of ESD of adjacent layers, as well as overlap between eigenvectors of adjecent layers, should be well-aligned.

## • "Scale collapse":

 "Size" of ESD of one or more layers changes dramatically, while the size of other layers changes very little, as a function of some perturbation of a model, during training (or post-training modification).

## • "Correlation traps":

 Spuriously large eigenvalues<sup>§</sup> may appear, and they may even be important for model convergence.

We can measure these quantities with Weightwatcher—so can you!

<sup>&</sup>lt;sup>§</sup>Eigenvalues not due to signal in the data—we have theorems-style theory for Hessians ("Hessian Eigenspectra of More Realistic Nonlinear Models." Liao and Mahonev. https://arxiv.org/abs/2103.01519). but it's still open for Weights

## Table of Contents

Practice

## 2 Theory



## Theorems

- Introduction
- Sample covariance matrix for large dimensional data: from concentration to RMT
- A random matrix perspective of the "curse of dimensionality"
- Kernel spectral clustering for large dimensional data
- A random matrix approach to large neural networks and random features

## Understanding the mechanism of large dimensional machine learning



- ▶ Big Data era: exploit large *n*, *p*, *N*
- counterintuitive phenomena when n ≫ p, e.g., the "curse of dimensionality"
- complete change of understanding of many ML algorithms
- <u>RMT</u> provides the tools!

## From low to high dimensional machine learning



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

• **low dimension**: data vectors  $\mathbf{x}_i \in \mathbb{R}^p$ , p = 2, 3, gathered in different "groups" can be classified using distance-based approach

### high dimension:

(i) **easy** or **trivial** scenario where low dimensional intuition holds and a pairwise distance-based classification approach via, e.g.,

Johnson–Lindenstrauss lemma, is efficient;

 (ii) hard or non-trivial scenario where such intuition collapses: data vectors at approximately the same Euclidean distance, regardless their arising from same or different classes.

## Non-trivial high dimensional classification beyond the JL regime

In the high dimensional regime where data dimension p and sample size n both large, a **dual** phenomenon:

- (i) data points not pairwise classifiable: Euclidean distance between any two data points  $\mathbf{x}_i \in C_a$  and  $\mathbf{x}_j \in C_b$ approximately constant  $\approx \tau = O(1)$  independent of their classes  $C_a, C_b$ :  $\|\mathbf{x}_i - \mathbf{x}_j\|^2 / p = \tau + o(1)$  as  $n, p \to \infty$  and data pairs *neither close nor far* from each other;
- (ii) classification remains possible by exploiting the spectral information of large Euclidean distance matrix  $\mathbf{E} = \{ \|\mathbf{x}_i \mathbf{x}_j\|^2 / p\}_{i,j=1}^n$ , thanks to a collective behavior of all data belonging to same (and large) classes.





Figure: Euclidean distance matrices **E**, the histogram of the entries of **E**, and the second top eigenvectors  $\mathbf{v}_2$ , for small (left, p = 5) and large (right, p = 250) dimensional data  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$  with  $\mathbf{x}_1, \dots, \mathbf{x}_{n/2} \in C_1$  and  $\mathbf{x}_{n/2+1}, \dots, \mathbf{x}_n \in C_2$  for  $n = 5\,000$  and different values of p.

## Sample covariance matrix in the large *n*, *p* regime

For  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$ , estimate population covariance  $\mathbf{C} \in \mathbb{R}^{p \times p}$  from *n* data samples  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$ .

Maximum likelihood sample covariance matrix with entry-wise convergence

$$\hat{\mathbf{C}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}} = \frac{1}{n} \mathbf{X} \mathbf{X}^{\mathsf{T}} \in \mathbb{R}^{p \times p}, \quad [\hat{\mathbf{C}}]_{ij} \to [\mathbf{C}]_{ij}$$

almost surely as  $n \to \infty$ : 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.

$$\|\hat{\mathbf{C}} - \mathbf{C}\| \not\to 0, \quad n, p \to \infty$$

 $\Rightarrow$  eigenvalue mismatch and not consistent!  $\Rightarrow$  matrix norms not equivalent in large dimensions!

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

## Quantitative spectral characterization of sample covariance

### Theorem (Concentration of sample covariance, [Ver18, Theorem 4.6.1])

Let  $\mathbf{X} \in \mathbb{R}^{p \times n}$  be a random matrix with i.i.d. sub-gaussian columns  $\mathbf{x}_i \in \mathbb{R}^p$  such that  $\mathbb{E}[\mathbf{x}_i] = \mathbf{0}$  and  $\mathbb{E}[\mathbf{x}_i \mathbf{x}_i^{\mathsf{T}}] = \mathbf{I}_p$ , one has, with probability at least  $1 - 2\exp(-t^2)$  for any  $t \ge 0$  that

$$\|\hat{\mathbf{C}} - \mathbf{I}_p\| \le C_1 \max(\delta, \delta^2), \quad \delta = C_2(\sqrt{p/n} + t/\sqrt{n})$$
(1)

for some constants  $C_1, C_2 > 0$  independent of n, p.

non-asymptotic and high probability characterization

• however, not precise in the  $p \sim n$  regime, since  $\delta = O(\sqrt{p/n}) = O(1)$ 

#### Theorem (Marčenko-Pastur law, [MP67])

Under the same setting of Theorem 1, as  $n, p \to \infty$  with  $p/n \to c \in (0, \infty)$ , with probability one, the empirical spectral measure  $\mu_{\hat{\mathbf{C}}} \equiv \frac{1}{p} \sum_{i=1}^{p} \delta_{\lambda_i(\hat{\mathbf{C}})}$  of  $\hat{\mathbf{C}} \equiv \frac{1}{n} \mathbf{X} \mathbf{X}^{\mathsf{T}}$  converges weakly to a probability measure  $\mu$  given explicitly by

$$\mu(dx) = (1 - c^{-1})^+ \delta_0(x) + \frac{1}{2\pi cx} \sqrt{(x - E_-)^+ (E_+ - x)^+} \, dx \tag{2}$$

where  $E_{\pm} = (1 \pm \sqrt{c})^2$  and  $(x)^+ = \max(0, x)$ , and is known as the Marčenko-Pastur law. M.W. Mahoney (UC Berkeley)

June 24, 2022 14 / 33

## Two ways of spectral characterization of sample covariance

matrix concentration-type characterization

$$\|\hat{\mathbf{C}} - \mathbf{I}_p\| \le C_1 \max(\delta, \delta^2), \quad \delta = C_2(\sqrt{p/n} + t/\sqrt{n})$$

⇒ non-asymptotic characterization of small dimensional intuition: how Ĉ concentrates around I<sub>p</sub>;
 random matrix-type characterization of precise eigenvalue distribution

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

 $\Rightarrow$  asymptotic characterization (as  $n, p \rightarrow \infty$ ) of large dimensional intuition: how  $\hat{\mathbf{C}}$  differs from  $\mathbf{I}_p$ !



Figure: Histogram of the eigenvalues of  $\hat{C}$  (blue) versus the Marčenko-Pastur law (red), for X having standard Gaussian entries in different settings: (left: small versus large dimensional intuition) p = 20, n = 100p versus p = 20, n = 100p; and (right: non-asymptotic versus asymptotic MP law) p = 20, n = 100p versus p = 500, n = 100p.

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

What about n = 100p? For  $\mathbf{C} = \mathbf{I}_p$ , as  $n, p \to \infty$  with  $p/n \to c \in (0, \infty)$ : the Marčenko–Pastur 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 = 50000.

• 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.

## Beyond eigenvalue distribution: a modern RMT approach via the resolvent

This **change-of-intuition** leads to very **different** behavior for small- versus large-dimensional ML:

- linear models: low-rank approximation, spectral classification/clustering, and linear least squares regression in high dimensions different from their small dimensional counterparts
- as well as more involved **nonlinear** models: kernel spectral clustering, nonlinear neural nets, etc. Technical challenges:
  - classical RMT focuses on eigenvalue distribution
  - ML applications need eigenvectors and more complex matrix functionals!



Figure: Different objects of interest and their corresponding technical tools for "old" and "new school" RMT.

## "Curse of dimensionality": loss of relevance of Euclidean distance

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

$$C_1$$
:  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_1, \mathbf{C}_1)$ , versus  $C_2$ :  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_2, \mathbf{C}_2)$ ;

Neyman-Pearson test: classification is possible only when [CLM18]

$$\|\mu_1 - \mu_2\| \ge C_{\mu}$$
, or  $\|\mathbf{C}_1 - \mathbf{C}_2\| \ge C_{\mathbf{C}} \cdot p^{-1/2}$ 

for some constants  $C_{\mu}$ ,  $C_{\mathbf{C}} > 0$ .

▶ In this non-trivial setting, for  $\mathbf{x}_i \in C_a, \mathbf{x}_j \in C_b$ :

$$\max_{1 \le i \ne j \le n} \left\{ \frac{1}{p} \| \mathbf{x}_i - \mathbf{x}_j \|^2 - \frac{2}{p} \operatorname{tr} \mathbf{C}^{\circ} \right\} \xrightarrow{a.s.} 0$$

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

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

<sup>&</sup>lt;sup>1</sup>Romain 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

## Reminder on kernel spectral clustering

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



## Reminder on kernel spectral clustering



 $\Downarrow$  *K*-dimensional representation  $\Downarrow$ 



Eigenvector 1

↓ EM or k-means clustering. (Three classes/clusters in this example.)

## Visualization of kernel matrices for large dimensional Gaussian data

**Objective**: "cluster" Gaussian data  $\mathbf{x}_1, \ldots, \mathbf{x}_n \in \mathbf{R}^p$  into  $C_1$  or  $C_2$ . Consider Gaussian kernel matrix  $\mathbf{K}_{ij} = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2p)$  and the second top eigenvectors  $\mathbf{v}_2$  for small (**left**) and large (**right**) dimensional data.



M. W. Mahoney (UC Berkeley)

## Kernel matrices for large dimensional real-world data



## A spectral viewpoint of large kernel matrices in large dimensions

► "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.,  $C_1 : \mathbf{x}_i = \boldsymbol{\mu}_1 + \mathbf{z}_i$  versus  $C_2 : \mathbf{x}_j = \boldsymbol{\mu}_2 + \mathbf{z}_j$ ) so that  $\|\mathbf{x}_i - \mathbf{x}_j\|^2/p \xrightarrow{a.s.} 2$ , and  $\mathbf{K} = \exp\left(-\frac{2}{2}\right) \left(\mathbf{1}_n \mathbf{1}_n^\mathsf{T} + \frac{1}{p} \mathbf{Z}^\mathsf{T} \mathbf{Z}\right) + g(\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|) \frac{1}{p} \mathbf{j} \mathbf{j}^\mathsf{T} + * + o_{\|\cdot\|}(1)$ 

with Gaussian matrix  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_n] \in \mathbb{R}^{p \times n}$  and  $\mathbf{j} = [\mathbf{1}_{n/2}; -\mathbf{1}_{n/2}]$ , the class-information vector **accumulated effect** of small "hidden" statistical information ( $\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|$  in this case)

Therefore

entry-wise:

$$\mathbf{K}_{ij} = \exp(-1)\left(1 + \underbrace{\frac{1}{p}\mathbf{z}_{i}^{\mathsf{T}}\mathbf{z}_{j}}_{O(p^{-1/2})}\right) \pm \underbrace{\frac{1}{p}g(\|\mu_{1} - \mu_{2}\|)}_{O(p^{-1})} + *, \text{ so that } \frac{1}{p}g(\|\mu_{1} - \mu_{2}\|) \ll \frac{1}{p}\mathbf{z}_{i}^{\mathsf{T}}\mathbf{z}_{j},$$

**spectrum-wise:** (i)  $\|\mathbf{K} - \exp(-1)\mathbf{1}_n\mathbf{1}_n^{\mathsf{T}}\| \neq 0$ ; (ii)  $\|\frac{1}{p}\mathbf{Z}^{\mathsf{T}}\mathbf{Z}\| = O(1)$  and  $\|g(\|\mu_1 - \mu_2\|)\frac{1}{p}\mathbf{j}\mathbf{j}^{\mathsf{T}}\| = O(1)!$ 

Same phenomenon as the sample covariance example:  $[\hat{\mathbf{C}} - \mathbf{C}]_{ij} \to 0 \Rightarrow ||\hat{\mathbf{C}} - \mathbf{C}|| \to 0!$ 

 $\Rightarrow$  With **modern RMT**, we **understand** kernel spectral clustering (eigenvectors!) for large dimensional data!

## Reminder on random features and neural networks

- kernel matrices  $\mathbf{K} \in \mathbb{R}^{n \times n}$  from pairwise comparison of *n* data points: expansive for *n* large
- <u>idea</u>: find easy-to-compute  $\hat{\mathbf{K}}$  to approximate  $\mathbf{K}$ , e.g.,  $\|\hat{\mathbf{K}} \mathbf{K}\|$  is small
- **example**: random Fourier feature [RR08]  $\Sigma^{\mathsf{T}} = [\cos(\mathsf{WX})^{\mathsf{T}}, \sin(\mathsf{WX})^{\mathsf{T}}] \in \mathbb{R}^{2N \times n}$  of data  $\overline{\mathsf{X}} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$  with standard Gaussian  $\mathbf{W} \in \mathbb{R}^{N \times p}$ , i.e.,  $\mathbf{W}_{ij} \sim \mathcal{N}(0, 1)$
- ▶ approximates Gaussian kernel  $\exp(||\mathbf{x}_i \mathbf{x}_j||^2/2)$ : entry-wise convergence of RFF Gram  $\frac{1}{N} [\boldsymbol{\Sigma}^{\mathsf{T}} \boldsymbol{\Sigma}]_{ij} \rightarrow [\mathbf{K}_{\text{Gauss}}]_{ij}$  Gaussian kernel matrix as number of features  $N \rightarrow \infty$
- **proof**: (strong) law of large numbers:

$$\frac{1}{N} [\boldsymbol{\Sigma}^{\mathsf{T}} \boldsymbol{\Sigma}]_{ij} = \frac{1}{N} \sum_{k=1}^{N} \cos(\mathbf{x}_{i}^{\mathsf{T}} \mathbf{w}_{k}) \cos(\mathbf{w}_{k}^{\mathsf{T}} \mathbf{x}_{j}) + \sin(\mathbf{x}_{i}^{\mathsf{T}} \mathbf{w}_{k}) \sin(\mathbf{w}_{k}^{\mathsf{T}} \mathbf{x}_{j})$$

$$\rightarrow \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{p})} [\cos(\mathbf{x}_{i}^{\mathsf{T}} \mathbf{w}) \cos(\mathbf{w}^{\mathsf{T}} \mathbf{x}_{j}) + \sin(\mathbf{x}_{i}^{\mathsf{T}} \mathbf{w}) \sin(\mathbf{w}^{\mathsf{T}} \mathbf{x}_{j})] = [\mathbf{K}_{\cos} + \mathbf{K}_{\sin}]_{ij} = [\mathbf{K}_{Gauss}]_{ij}$$
for  $\mathbf{K}_{\cos} = e^{-\frac{1}{2}(||\mathbf{x}_{i}||^{2} + ||\mathbf{x}_{j}||^{2})} \cosh(\mathbf{x}_{i}^{\mathsf{T}} \mathbf{x}_{j})$  and  $\mathbf{K}_{\sin} = e^{-\frac{1}{2}(||\mathbf{x}_{i}||^{2} + ||\mathbf{x}_{j}||^{2})} \sinh(\mathbf{x}_{i}^{\mathsf{T}} \mathbf{x}_{j}).$ 

<sup>&</sup>lt;sup>3</sup>Ali Rahimi and Benjamin Recht. "Random Features for Large-Scale Kernel Machines". In: Advances in Neural Information Processing Systems. Vol. 20. NIPS'08. Curran Associates, Inc., 2008, pp. 1177–1184

## Random features-based ridge regression and neural networks

$$\begin{array}{c} & & \mathbf{W} \in \mathbb{R}^{N \times p} & \underbrace{\sin}_{\mathbf{Cos}} & \boldsymbol{\beta} \in \mathbb{R}^{2N} \text{ in } (3) \\ \hline & & \mathbf{X} \in \mathbb{R}^{p \times n} \\ & & \mathbf{X} \in \mathbb{R}^{p \times \hat{n}} \end{array} \xrightarrow{ \boldsymbol{\Sigma}_{\mathbf{X}}^{\mathsf{T}} = \boldsymbol{\Sigma}^{\mathsf{T}} = [\cos(\mathbf{W}\mathbf{X})^{\mathsf{T}}, \sin(\mathbf{W}\mathbf{X})^{\mathsf{T}}] \\ & & \boldsymbol{\Sigma}_{\hat{\mathbf{X}}}^{\mathsf{T}} = [\cos(\mathbf{W}\hat{\mathbf{X}})^{\mathsf{T}}, \sin(\mathbf{W}\hat{\mathbf{X}})^{\mathsf{T}}] \end{array}$$

Figure: Illustration of random Fourier features regression model.

▶ RFF ridge regressor  $\beta \in \mathbb{R}^{2N}$  given by, for regularization penalty  $\gamma \ge 0$ ,

$$\boldsymbol{\beta} \equiv \frac{1}{n} \boldsymbol{\Sigma} (\frac{1}{n} \boldsymbol{\Sigma}^{\mathsf{T}} \boldsymbol{\Sigma} + \gamma \mathbf{I}_n)^{-1} \mathbf{y} \cdot \mathbf{1}_{2N>n} + (\frac{1}{n} \boldsymbol{\Sigma} \boldsymbol{\Sigma}^{\mathsf{T}} + \gamma \mathbf{I}_{2N})^{-1} \frac{1}{n} \boldsymbol{\Sigma} \, \mathbf{y} \cdot \mathbf{1}_{2N(3)$$

▶ **Performance**: training and test Mean Squared Error (MSE):  $E_{\text{train}} = \frac{1}{n} ||\mathbf{y} - \mathbf{\Sigma}_{\mathbf{X}}^{\mathsf{T}} \boldsymbol{\beta}||^2$  and  $E_{\text{test}} = \frac{1}{\hat{n}} ||\hat{\mathbf{y}} - \mathbf{\Sigma}_{\mathbf{X}}^{\mathsf{T}} \boldsymbol{\beta}||^2$ , with  $\mathbf{\Sigma}_{\mathbf{X}}^{\mathsf{T}} \in \mathbb{R}^{\hat{n} \times 2N}$  RFFs of a test set  $(\hat{\mathbf{X}}, \hat{\mathbf{y}})$  of size  $\hat{n}$ .

single-hidden-layer neural network with cos + sin activations, connected to neural tangent kernel (NTK)

<sup>&</sup>lt;sup>3</sup>Arthur 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

## Random Fourier features approximate Gaussian kernel, but in which sense?

- ► [RR08]: entry-wise convergence of RFF Gram  $\frac{1}{N} [\Sigma^{\mathsf{T}} \Sigma]_{ij} \rightarrow [\mathbf{K}_{\text{Gauss}}]_{ij}$  Gaussian kernel matrix as  $N \rightarrow \infty$
- ► again, **not true** in spectral norm sense, i.e.,  $\|\mathbf{\Sigma}^{\mathsf{T}}\mathbf{\Sigma}/N \mathbf{K}_{\text{Gauss}}\| \neq 0$  unless  $N \gg n$ 
  - − e.g.,  $\Sigma^{\mathsf{T}}\Sigma \in \mathbb{R}^{n \times n}$  of rank at most *N* if *N* ≤ *n*, while **K**<sub>Gauss</sub> of rank *n* (for distinct **x**<sub>*i*</sub>)

- significant impact on various RFF-based algorithms



Figure: Training MSEs of RFF ridge regression on MNIST data (class 3 versus 7) as a function of regression penalty  $\lambda$ .

- effective kernel can be derived with RMT in the large *n*, *p*, *N* regime
- ▶ provides precise training and test performances of RFF for any ratio N/n, more practical and more flexible, recover Gaussian kernel result with  $N/n \rightarrow \infty$
- data-dependent theory with no strong assumption on data

## Sharp analysis of RFF ridge regression performance via RMT



Figure: MSEs of RFF ridge regression on Fashion- (left two) and Kannada-MNIST (right two).



Figure: Test MSEs of RFF regression as a function of the ratio N/n, on MNIST data set.

## "Recap" for double descent phenomenon for over-parameterized models



Figure: Comparison between training risk (blue) and true/test risk (red).

- empirically observed for various large-scale machine learning models, e.g., RF-based methods, decision trees, ensemble methods, and deep NNs
- **proved** here for RFF on real-world data!

**b** phase transition from under- to over-param of resolvent  $(\Sigma^{\mathsf{T}}\Sigma + \lambda \mathbf{I}_n)^{-1}$  in the ridgeless  $\lambda \to 0$  limit

<sup>4</sup>Mikhail Belkin et al. "Reconciling modern machine-learning practice and the classical bias–variance trade-off". In: *Proceedings of the National Academy of Sciences* 116.32 (2019), pp. 15849–15854

<sup>5</sup>Trevor Hastie et al. "Surprises in High-Dimensional Ridgeless Least Squares Interpolation". In: arXiv (2019). eprint: 1903.08560

## Take-away messages and references

#### Take-away messages:

- ▶ RF methods: classical statistical learning theory provides performance guarantee for  $N \gg n, p$
- ▶ here we derive (limiting) kernel in the more practical large *n*, *p*, *N* regime
- fast tuning of regularization parameter  $\lambda$
- double descent theory for novel understanding of over-parameterized neural networks

#### **References**:

- Zhenyu Liao, Romain Couillet, and Michael W Mahoney. "A random matrix analysis of random Fourier features: beyond the Gaussian kernel, a precise phase transition, and the corresponding double descent". In: Advances in Neural Information Processing Systems (NeurIPS). vol. 33. Curran Associates, Inc., 2020, pp. 13939–13950
- Cosme Louart, Zhenyu Liao, and Romain Couillet. "A random matrix approach to neural networks". In: Annals of Applied Probability 28.2 (2018), pp. 1190–1248
- Song Mei and Andrea Montanari. "The Generalization Error of Random Features Regression: Precise Asymptotics and the Double Descent Curve". In: Communications on Pure and Applied Mathematics (2021)
- Zhenyu Liao and Michael W. Mahoney. "Hessian Eigenspectra of More Realistic Nonlinear Models". In: Advances in Neural Information Processing Systems (NeurIPS). 2021

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!

# Thank you! Q & A?