Why Deep Learning Works: Traditional and Heavy-Tailed Implicit Self-Regularization in Deep Neural Networks

Michael W. Mahoney

ICSI and Dept of Statistics, UC Berkeley

http://www.stat.berkeley.edu/~mmahoney/

June 2019

(Joint work with Charles H. Martin, Calculation Consulting, charles@calculationconsulting.com)

(4) (3) (4) (4) (4)

Outline



Background

- 3 Preliminary Empirical Results
- 5 More detailed empirical results
- 6 An RMT-based Theory for Deep Learning
- Tikhonov regularization versus Heavy-tailed regularization
- 8 Varying the Batch Size: Explaining the Generalization Gap
- Output Description (1998) Using the Theory

Motivations: Theoretical AND Practical

Theoretical: deeper insight into Why Deep Learning Works?

- convex versus non-convex optimization?
- explicit/implicit regularization?
- is / why is / when is deep better?
- VC theory versus Statistical Mechanics theory?

• . . .

Practical: use insights to improve engineering of DNNs?

- when is a network fully optimized?
- can we use labels and/or domain knowledge more efficiently?
- large batch versus small batch in optimization?
- designing better ensembles?

```
• . . .
```

通 ト イ ヨ ト イ ヨ ト

Motivations: towards a Theory of Deep Learning



FIG. 1. Plot of the foldability distribution $p(\mathcal{F})$ for different numbers of compact states (n = 10, 100, 1081, 103 346), calculated using the random energy model.

Raises broad questions about Why Deep Learning Works

ation in DNNs

A D N A B N A B N A B N

Energy Landscape Theory

June 2019

Implicit Self-regularization in DNNs

Motivations: regularization in DNNs?

ICLR 2017 Best paper

- Large neural network models can easily overtrain/overfit on randomly labeled data
- Popular ways to regularize (basically $\min_x f(x) + \lambda g(x)$, with "control parameter" λ) may or may not help.

Understanding deep learning requires rethinking generalization?? https://arxiv.org/abs/1611.03530

Rethinking generalization requires revisiting old ideas: statistical mechanics approaches and complex learning behavior!! https://arxiv.org/abs/1710.09553 (Martin & Mahoney)

く 目 ト く ヨ ト く ヨ ト

Set up: the Energy Landscape

Energy/Optimization function:

 $E_{DNN} = h_L(\mathbf{W}_L \times h_{L-1}(\mathbf{W}_{L-1} \times h_{L-2}(\cdots) + \mathbf{b}_{L-1}) + \mathbf{b}_L)$

Train this on labeled data $\{d_i, y_i\} \in \mathcal{D}$, using Backprop, by minimizing loss \mathcal{L} :

$$\min_{W_i, b_i} \mathcal{L}\left(\sum_i E_{DNN}(d_i) - y_i\right)$$

*E*_{DNN} is "the" *Energy Landscape*:

- The part of the optimization problem parameterized by the heretofore unknown elements of the weight matrices and bias vectors, and as defined by the data {d_i, y_i} ∈ D
- Pass the data through the Energy function E_{DNN} multiple times, as we run Backprop training
- The Energy Landscape* is *changing* at each epoch

*i.e., the optimization function that is nominally being optimized $\langle a \rangle \langle a \rangle \langle a \rangle$

Mahoney (UC Berkeley)

Implicit Self-regularization in DNNs

Problem: How can this possibly work?



It has been known for a long time that local minima are not the issue.

Mahoney (UC Berkeley)

June 2019 7 / 62

∃ >

Problem: Local Minima?



Duda, Hart and Stork, 2000

Whereas in low-dimensional spaces, local minima can be plentiful, in high dimension, the problem of local minima is different: The high-dimensional space may afford more ways (dimensions) for the system to "get around" a barrier or local maximum during learning. The more superfluous the weights, the less likely it is a network will get trapped in local minima. However, networks with an unnecessarily large number of weights are undesirable because of the dangers of overfitting, as we shall see in Section 6.11.

Solution: add more capacity and regularize, i.e., over-parameterization

Motivations: what is regularization?



Every adjustable *knob* and *switch*—and there are *many*[†]—is regularization.

[†]https://arxiv.org/pdf/1710.10686.pdf

Outline

- Background
- 2 Regularization and the Energy Landscape
- 3 Preliminary Empirical Results
- 4 Gaussian and Heavy-tailed Random Matrix Theory
- 5 More detailed empirical results
- 6 An RMT-based Theory for Deep Learning
- 🕜 Tikhonov regularization versus Heavy-tailed regularization
- 8 Varying the Batch Size: Explaining the Generalization Gap
- Output Description Using the Theory
- More General Implications
- Conclusions

Basics of Regularization

Ridge Regression / Tikhonov-Phillips Regularization

$$\begin{split} \hat{\mathbf{W}} \mathbf{x} &= \mathbf{y} \\ \mathbf{x} &= \left(\hat{\mathbf{W}}^T \hat{\mathbf{W}} + \alpha I \right)^{-1} \hat{\mathbf{W}}^T \mathbf{y} \\ \min_{\mathbf{x}} \| \hat{\mathbf{W}} \mathbf{x} - \mathbf{y} \|_2^2 + \alpha \| \hat{\mathbf{x}} \|_2^2 \end{split}$$

Moore-Penrose pseudoinverse (1955) Ridge regularization (Phillips, 1962)

familiar optimization problem

Softens the rank of $\hat{\mathbf{W}}$ to focus on large eigenvalues.

Related to Truncated SVD, which does hard truncation on rank of \hat{W}

Early stopping, truncated random walks, etc. often implicitly solve regularized optimiation problems.

How we will study regularization

The Energy Landscape is *determined* by layer weight matrices W_L :

$$E_{DNN} = h_L(\mathbf{W}_L \times h_{L-1}(\mathbf{W}_{L-1} \times h_{L-2}(\cdots) + \mathbf{b}_{L-1}) + \mathbf{b}_L)$$

Traditional regularization is applied to \mathbf{W}_L :

$$\min_{W_l, b_l} \mathcal{L}\left(\sum_i E_{DNN}(d_i) - y_i\right) + \alpha \sum_l \|\mathbf{W}_l\|$$

Different types of regularization, e.g., different norms $\|\cdot\|$, leave different empirical signatures on \mathbf{W}_L .

What we do:

- Turn off "all" regularization.
- Systematically turn it back on, explicitly with α or implicitly with knobs/switches.
- Study empirical properties of \mathbf{W}_L .

< □ > < □ > < □ > < □ > < □ > < □ >

Lots of DNNs Analyzed

Question: What happens to the layer weight matrices W_L ?

(Don't evaluate your method on one/two/three NN, evaluate it on a dozen/hundred.)

Retrained LeNet5 on MINST using Keras.

Two other small models:

- 3-Layer MLP
- Mini AlexNet



Wide range of state-of-the-art pre-trained models:

• AlexNet, Inception, etc.

Mahoney (UC Berkeley)

Outline

- Background
- 2 Regularization and the Energy Landscape
- O Preliminary Empirical Results
 - 4 Gaussian and Heavy-tailed Random Matrix Theory
- 5 More detailed empirical results
- 6 An RMT-based Theory for Deep Learning
- 🕜 Tikhonov regularization versus Heavy-tailed regularization
- 8 Varying the Batch Size: Explaining the Generalization Gap

(日) (문) (문) (문) (문)

- Using the Theory
- More General Implications
- Conclusions

Matrix complexity: Matrix Entropy and Stable Rank

$$\mathbf{W} = \mathbf{U} \Sigma \mathbf{V}^T \qquad \nu_i = \Sigma_{ii} \qquad p_i = \nu_i^2 / \sum_i \nu_i^2$$
$$\mathcal{S}(\mathbf{W}) = \frac{-1}{\log(R(\mathbf{W}))} \sum_i p_i \log p_i \qquad \mathcal{R}_s(\mathbf{W}) = \frac{\|\mathbf{W}\|_F^2}{\|\mathbf{W}\|_2^2} = \frac{\sum_i \nu_i^2}{\nu_{max}^2}$$

A warm-up: train a 3-Layer MLP:



Figure: Matrix Entropy & Stable Rank show transition during Backprop training.

Mahoney (UC Berkeley)

Matrix complexity: Scree Plots

$$\mathbf{W} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{T} \qquad \nu_{i} = \Sigma_{ii} \qquad p_{i} = \nu_{i}^{2} / \sum_{i} \nu_{i}^{2}$$
$$\mathcal{S}(\mathbf{W}) = \frac{-1}{\log(R(\mathbf{W}))} \sum_{i} p_{i} \log p_{i} \qquad \mathcal{R}_{s}(\mathbf{W}) = \frac{\|\mathbf{W}\|_{F}^{2}}{\|\mathbf{W}\|_{2}^{2}} = \frac{\sum_{i} \nu_{i}^{2}}{\nu_{max}^{2}}$$

A warm-up: train a 3-Layer MLP:



Figure: Scree plots for initial and final configurations for MLP3.

Mahoney (UC Berkeley)

June 2019 16 / 62

Matrix complexity: Singular/Eigen Value Densities

$$\mathbf{W} = \mathbf{U} \Sigma \mathbf{V}^{T} \qquad \nu_{i} = \Sigma_{ii} \qquad p_{i} = \nu_{i}^{2} / \sum_{i} \nu_{i}^{2}$$
$$\mathcal{S}(\mathbf{W}) = \frac{-1}{\log(R(\mathbf{W}))} \sum_{i} p_{i} \log p_{i} \qquad \mathcal{R}_{s}(\mathbf{W}) = \frac{\|\mathbf{W}\|_{F}^{2}}{\|\mathbf{W}\|_{2}^{2}} = \frac{\sum_{i} \nu_{i}^{2}}{\nu_{max}^{2}}$$

A warm-up: train a 3-Layer MLP:



Figure: Histograms of the Singular Values ν_i and associated Eigenvalues $\lambda_i = \nu_i^2$.

Mahoney (UC Berkeley)

ESD: detailed insight into W_L

Empirical Spectral Density (ESD: eigenvalues of $X = \mathbf{W}_{L}^{T}\mathbf{W}_{L}$)

import keras

```
import numpy as np
import matplotlib.pyplot as plt
...
W = model.layers[i].get_weights()[0]
...
X = np.dot(W, W.T)
evals, evecs = np.linalg.eig(W, W.T)
plt.hist(X, bin=100, density=True)
```

ESD: detailed insight into W_L

Empirical Spectral Density (ESD: eigenvalues of $X = \mathbf{W}_L^T \mathbf{W}_L$)



Entropy decrease corresponds to:

- modification (later, breakdown) of random structure and
- onset of a new kind of self-regularization.

・ 何 ト ・ ヨ ト ・ ヨ ト

Outline

- Background
- 2 Regularization and the Energy Landscape
- In the second second
- Gaussian and Heavy-tailed Random Matrix Theory
- 5 More detailed empirical results
- 6 An RMT-based Theory for Deep Learning
- Tikhonov regularization versus Heavy-tailed regularization
- 8 Varying the Batch Size: Explaining the Generalization Gap
- Output Description Using the Theory
- More General Implications
- 11 Conclusions

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

Let **W** be an $N \times M$ random matrix, with elements $W_{ij} \sim N(0, \sigma_{mp}^2)$. Then, the ESD of **X** = **W**^T**W**, converges to a deterministic function:

$$\rho_{N}(\lambda) := \frac{1}{N} \sum_{i=1}^{M} \delta(\lambda - \lambda_{i})$$

$$\xrightarrow{N \to \infty}_{Q \text{ fixed}} \begin{cases} \frac{Q}{2\pi\sigma_{mp}^{2}} \frac{\sqrt{(\lambda^{+} - \lambda)(\lambda - \lambda^{-})}}{\lambda} & \text{if } \lambda \in [\lambda^{-}, \lambda^{+}] \\ 0 & \text{otherwise.} \end{cases}$$

with well-defined edges (which depend on Q, the aspect ratio):

$$\lambda^{\pm} = \sigma_{mp}^2 \left(1 \pm rac{1}{\sqrt{Q}}
ight)^2 \qquad Q = N/M \ge 1.$$

Random Matrix Theory 102': Marchenko-Pastur



(a) Vary aspect ratios



Figure: Marchenko-Pastur (MP) distributions.

Important points:

- Global bulk stats: The overall shape is deterministic, fixed by Q and σ .
- Local edge stats: The edge λ^+ 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.

Fitting Heavy-tailed Distributions



Figure: The log-log histogram plots of the ESD for three Heavy-Tailed random matrices **M** with same aspect ratio Q = 3, with $\mu = 1.0, 3.0, 5.0$, corresponding to the three Heavy-Tailed Universality classes ($0 < \mu < 2$ vs $2 < \mu < 4$ and $4 < \mu$).

Non-negligibe finite size effects



Figure: Dependence of α (the fitted PL parameter) on μ (the hypothesized limiting PL parameter).

June 2019 26 / 62

< □ > < □ > < □ > < □ > < □ > < □ >

Outline

- Background
- 2 Regularization and the Energy Landscape
- In the second second
- Gaussian and Heavy-tailed Random Matrix Theory
- 5 More detailed empirical results
- 6 An RMT-based Theory for Deep Learning
- Tikhonov regularization versus Heavy-tailed regularization
- 8 Varying the Batch Size: Explaining the Generalization Gap
- Output Description Using the Theory
- More General Implications
- Conclusions

Experiments: just apply this to pre-trained models

https://medium.com/@siddharthdas_32104/cnns-architectures-lenet-alexnet-vgg-googlenet-resnet-and-more-...

Year	CNN	Developed by	Place	Top-5 error rate	No. of parameters
1998	LeNet(8)	Yann LeCun et al			60 thousand
2012	AlexNet(7)	Alex Krizhevsky, Geoffrey Hinton, Ilya Sutskever	1st	15.3%	60 million
2013	ZFNet()	Matthew Zeiler and Rob Fergus	1st	14.8%	
2014	GoogLeNet(1 9)	Google	1st	6.67%	4 million
2014	VGG Net(16)	Simonyan, Zisserman	2nd	7.3%	138 million
2015	ResNet(152)	Kaiming He	1st	3.6%	



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Mahoney (UC Berkeley)

Implicit Self-regularization in DNNs

Experiments: just apply this to pre-trained models

Model	Lovor	0	$(M \times N)$	0	n	Best
Woder	Layer	^v	$(M \times N)$	a		Fit
alexnet	17/FC1	2.25	(4096×9216)	2.29	0.0527	PL
	20/FC2	1	(4096×4096)	2.25	0.0372	PL
	22/FC3	4.1	(1000×4096)	3.02	0.0186	PL
densenet121	432	1.02	(1000×1024)	3.32	0.0383	PL
densenet121	432	1.02	(1000×1024)	3.32	0.0383	PL
densenet161	572	2.21	(1000×2208)	3.45	0.0322	PL
densenet169	600	1.66	(1000×1664)	3.38	0.0396	PL
densenet201	712	1.92	(1000×1920)	3.41	0.0332	PL
inception v3	L226	1.3	(768×1000)	5.26	0.0421	PL
	L302	2.05	(1000×2048)	4.48	0.0275	PL
resnet101	286	2.05	(1000×2048)	3.57	0.0278	PL
resnet152	422	2.05	(1000×2048)	3.52	0.0298	PL
resnet18	67	1.95	(512×1000)	3.34	0.0342	PL
resnet34	115	1.95	(512×1000)	3.39	0.0257	PL
resnet50	150	2.05	(1000×2048)	3.54	0.027	PL
vgg11	24	6.12	(4096×25088)	2.32	0.0327	PL
	27	1	(4096×4096)	2.17	0.0309	TPL
	30	4.1	(1000×4096)	2.83	0.0398	PL
vgg11 bn	32	6.12	(4096×25088)	2.07	0.0311	TPL
	35	1	(4096×4096)	1.95	0.0336	TPL
	38	4.1	(1000×4096)	2.99	0.0339	\mathbf{PL}
vgg16	34	6.12	(4096×25088)	2.3	0.0277	PL
	37	1	(4096×4096)	2.18	0.0321	TPL
	40	4.1	(1000×4096)	2.09	0.0403	TPL
vgg16 bn	47	6.12	(4096×25088)	2.05	0.0285	TPL
	50	1	(4096×4096)	1.97	0.0363	TPL
	53	4.1	(1000×4096)	3.03	0.0358	\mathbf{PL}
vgg19	40	6.12	(4096×25088)	2.27	0.0247	PL
	43	1	(4096×4096)	2.19	0.0313	PL
	46	4.1	(1000×4096)	2.07	0.0368	TPL
vgg19 bn	56	6.12	(4096×25088)	2.04	0.0295	TPL
	59	1	(4096×4096)	1.98	0.0373	TPL
	62	4.1	(1000×4096)	3.03	0.035	PL

Mahoney (UC Berkeley)

June 2019 29 / 62

RMT: LeNet5 (an old/small NN example)



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

Marchenko-Pastur Bulk + Spikes

Mahoney (UC Berkeley)

Implicit Self-regularization in DNNs

RMT: AlexNet (a typical modern/large DNN example)



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

Marchenko-Pastur Bulk-decay + Heavy-tailed

Mahoney (UC Berkeley)

Implicit Self-regularization in DNNs



Figure: ESD for Layers L226 and L302 in InceptionV3, as distributed w/ pyTorch.

Marchenko-Pastur bulk decay, onset of Heavy Tails

Outline

- Background
- 2 Regularization and the Energy Landscape
- In the second second
- 4 Gaussian and Heavy-tailed Random Matrix Theory
- 5 More detailed empirical results
- 6 An RMT-based Theory for Deep Learning
- 7 Tikhonov regularization versus Heavy-tailed regularization
- 8 Varying the Batch Size: Explaining the Generalization Gap
- Output Description Using the Theory
- More General Implications
- Conclusions

RMT-based 5+1 Phases of Training



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

Mahoney (UC Berkeley)

June 2019 34 / 62

RMT-based 5+1 Phases of Training

We model "noise" and also "signal" with random matrices:

$$\mathbf{W} \simeq \mathbf{W}^{rand} + \Delta^{sig}. \tag{1}$$

	Operational Definition	Informal Description via Eqn. (1)	Edge/tail Fluctuation Comments	Illustration and Description
Random-like	ESD well-fit by MP with appropriate λ^+	\mathbf{W}^{rand} random; $\ \Delta^{sig}\ $ zero or small	$\lambda_{max} pprox \lambda^+$ is sharp, with TW statistics	Fig. 10(a)
Bleeding-out	ESD RANDOM-LIKE, excluding eigenmass just above λ^+	W has eigenmass at bulk edge as spikes "pull out"; Δ ^{sig} medium	$\begin{array}{c} BPP \text{ transition,} \\ \lambda_{max} \text{ and} \\ \lambda^+ \text{ separate} \end{array}$	Fig. 10(b)
Bulk+Spikes	$\begin{array}{l} ESD \mathrm{Random-Like} \\ plus \geq 1 spikes \\ well above \lambda^+ \end{array}$	\mathbf{W}^{rand} well-separated from low-rank Δ^{sig} ; $\ \Delta^{sig}\ $ larger	λ^+ is TW, λ_{max} is Gaussian	Fig. 10(c)
Bulk-decay	ESD less RANDOM-LIKE; Heavy-Tailed eigenmass above λ^+ ; some spikes	Complex ∆ ^{sig} with correlations that don't fully enter spike	Edge above λ^+ is not concave	Fig. 10(d)
Heavy-Tailed	ESD better-described by Heavy-Tailed RMT than Gaussian RMT	\mathbf{W}^{rand} is small; Δ^{sig} is large and strongly-correlated	No good λ^+ ; $\lambda_{max} \gg \lambda^+$	Fig. 10(e)
RANK-COLLAPSE	ESD has large-mass spike at $\lambda = 0$	W very rank-deficient; over-regularization	_	Fig. 10(f)

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

RMT-based 5+1 Phases of Training

Lots of technical issues ...

< ∃ ►

Outline

- Background
- 2 Regularization and the Energy Landscape
- 3 Preliminary Empirical Results
- 4 Gaussian and Heavy-tailed Random Matrix Theory
- 5 More detailed empirical results
- 6 An RMT-based Theory for Deep Learning
- Tikhonov regularization versus Heavy-tailed regularization
- 8 Varying the Batch Size: Explaining the Generalization Gap
- Output Description Using the Theory
- More General Implications
- Conclusions

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

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

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

Implicit Self-regularization in DNNs

June 2019 38 / 62

Heavy-tailed Self-regularization

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

- We model strongly-correlated systems by heavy-tailed random matrices
- I.e., we model signal (not noise) by heavy-tailed random matrices

Then RMT/MP ESD will also have heavy tails

Known results from RMT / polymer theory (Bouchaud, Potters, etc)





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

Heavy-tailed Self-regularization

Summary of what we "suspect" today

- No single scale threshold.
- No simple low rank approximation for \mathbf{W}_L .
- Contributions from correlations at all scales.
- Can not be treated perturbatively.

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

Spikes: carry more "information" than the Bulk

Spikes have less entropy, are more localized than bulk.







(a) Vector Entropies.

(b) Localization Ratios.

(c) Participation Ratios.

() < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < ()

Figure: Eigenvector localization metrics for the FC1 layer of MiniAlexNet.

Information begins to concentrate in the spikes.

Power Law Universality: ImageNet and AllenNLP



Figure 12: Distribution of power law exponents α for linear layers in pre-trained models trained on ImageNet, available in pyTorch, and for those NLP models, available in AllenNLP.

All these models display remarkable Heavy Tailed Universality

< ⊒ >

Power Law Universality: ImageNet



All these models display remarkable Heavy Tailed Universality

Power Law Universality: GPT versus GPT2

GPT and GPT2 Layer Weight Matrix Power Law Exponents α , $\rho(\lambda) \sim \lambda^{-\alpha}$



GPT versus GPT2: example of a class of models that "improves" over time.

< 3 >

Outline

- Background
- 2 Regularization and the Energy Landscape
- 3 Preliminary Empirical Results
- 4 Gaussian and Heavy-tailed Random Matrix Theory
- 5 More detailed empirical results
- 6 An RMT-based Theory for Deep Learning
- Tikhonov regularization versus Heavy-tailed regularization
- 8 Varying the Batch Size: Explaining the Generalization Gap

< □ > < (四 > < (回 >) < (u >

- Using the Theory
- More General Implications
- Conclusions

Self-regularization: Batch size experiments

A theory should make predictions:

- We predict the existence of 5+1 phases of increasing implicit self-regularization
- We characterize their properties in terms of HT-RMT

Do these phases exist? Can we find them?

There are *many* knobs. Let's vary one—batch size.

- Tune the batch size from very large to very small
- A small (i.e., retrainable) model exhibits all 5+1 phases
- Large batch sizes => decrease generalization accuracy
- Large batch sizes => decrease implicit self-regularization

Generalization Gap Phenomena: all else being equal, small batch sizes lead to more implicitly self-regularized models.

・ 何 ト ・ ヨ ト ・ ヨ ト

Batch Size Tuning: Generalization Gap



Figure: Varying Batch Size: Stable Rank and MP Softrank for FC1 and FC2 Training and Test Accuracies versus Batch Size for MiniAlexNet.

- Decreasing batch size leads to better results—it induces strong correlations in **W**.
- Increasing batch size leads to worse results—it washes out strong correlations in **W**.

∃ >

Batch Size Tuning: Generalization Gap



Figure: Varying Batch Size. ESD for Layer FC1 of MiniAlexNet. We exhibit all 5 of the main phases of training by varying only the batch size.

- Decreasing batch size induces strong correlations in **W**, leading to a more implicitly-regularized model.
- Increasing batch size washes out strong correlations in **W**, leading to a less implicitly-regularized model.

Mahoney (UC Berkeley)

Implicit Self-regularization in DNNs

Summary so far

applied Random Matrix Theory (RMT)

self-regularization \sim entropy / information decrease

5+1 phases of learning

small models \sim Tinkhonov-like regularization

modern DNNs \sim heavy-tailed self-regularization

Remarkably ubiquitous

How can this be used?

Why does deep learning work?

Mahoney (UC Berkeley)

Implicit Self-regularization in DNNs

June 2019 49 / 62

Outline

- Background
- 2 Regularization and the Energy Landscape
- In the second second
- 4 Gaussian and Heavy-tailed Random Matrix Theory
- 5 More detailed empirical results
- 6 An RMT-based Theory for Deep Learning
- Tikhonov regularization versus Heavy-tailed regularization
- 8 Varying the Batch Size: Explaining the Generalization Gap

(日) (문) (문) (문) (문)

- O Using the Theory
 - More General Implications
- Conclusions

Predicting test accuracies: Product norms

M&M: "Heavy-Tailed Universality Predicts Trends in Test Accuracies ... Pre-Trained ..." https://arxiv.org/abs/1901.08278

Test Accuracy vs Average Log Norm (log W/II)

The product norm is a VC-like data-dependent capacity metric for DNNs. We can predict trends in the test accuracy—*without peeking at the test data!*

"pip install weightwatcher"

Mahoney (UC Berkelev	
manoney (OC Derkeley	

Implicit Self-regularization in DNNs

June 2019 51 / 62

Universality, capacity control, and norm-powerlaw relations

M&M: "Heavy-Tailed Universality Predicts Trends in Test Accuracies ... Pre-Trained ..." https://arxiv.org/abs/1901.08278

- "Universality" suggests the power law exponent α would make a good, Universal, DNN capacity control metric
- Consider a weighted average

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

- To get weights b_{1,i}, relate Frobenius norm and Power Law exponent.
- Create a random Heavy Tailed (Pareto) matrix:

$$\mathsf{Pr}\left(W^{\mathit{rand}}_{i,j}
ight)\sim rac{1}{x^{1+\mu}}$$

• Examine the norm-powerlaw relations:

$$\frac{\log \|\mathbf{W}\|_F^2}{\log \lambda_{\max}} \quad \text{versus} \quad \alpha$$

• Argue that:

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

Predicting test accuracies better: Weighted Power Laws

M&M: "Heavy-Tailed Universality Predicts Trends in Test Accuracies ... Pre-Trained ..." https://arxiv.org/abs/1901.08278





We can predict trends in the test accuracy-without peeking at the test data!

"pip install weightwatcher"

Mahoney (UC Berkeley)

Outline

- Background
- 2 Regularization and the Energy Landscape
- 3 Preliminary Empirical Results
- 4 Gaussian and Heavy-tailed Random Matrix Theory
- 5 More detailed empirical results
- 6 An RMT-based Theory for Deep Learning
- Tikhonov regularization versus Heavy-tailed regularization
- 8 Varying the Batch Size: Explaining the Generalization Gap
- Using the Theory
- 10 More General Implications
 - Conclusions

Rethinking generalization requires revisiting old ideas

Martin and Mahoney https://arxiv.org/abs/1710.09553

Very Simple Deep Learning (VSDL) model:

- DNN is a black box, load-like parameters lpha, & temperature-like parameters au
- $\bullet\,$ Adding noise to training data decreases α
- Early stopping increases τ

Nearly any non-trivial model^{\ddagger} exhibits "phase diagrams," with *qualitatively* different generalization properties, for different parameter values.



(a) Training/generalization (b) Learning phases in τ - α (c) Noisifying data and adjusterror in the VSDL model. In knobs.

[‡]when analyzed via the Statistical Mechanics Theory of Generalization (SMToG) 🛌 💿 🔍

Mahoney (UC Berkeley)	Implicit Self-regularization in DNNs	June 2019 55 / 62
-----------------------	--------------------------------------	-------------------

Remembering Regularization

Martin and Mahoney https://arxiv.org/abs/1710.09553

Statistical Mechanics (1990s): (this) Overtraining \rightarrow Spin Glass Phase



Binary Classifier with N Random Labelings:

 $2^{\it N}$ over-trained solutions: locally (ruggedly) convex, very high barriers, all unable to generalize

Implications: Minimizing Frustration and Energy Funnels

As simple as can be?, Wolynes, 1997



Energy Landscape Theory: "random heteropolymer" versus "natural protein" folding

Mahoney (UC Berkeley

Implications: The Spin Glass of Minimal Frustration

https://calculatedcontent.com/2015/03/25/why-does-deep-learning-work/



< ⊒ >

Implications: Rugged Energy Landscapes of Heavy-tailed Models

Martin and Mahoney https://arxiv.org/abs/1710.09553



Spin Glasses with Heavy Tails?

- Local minima do not concentrate near the ground state (Cizeau and Bouchaud 1993)
- Configuration space with a "rugged convexity"

Contrast with (Gaussian) Spin Glass model of Choromanska et al. 2015

If Energy Landscape is ruggedly funneled, then no "problems" with local minima!

Outline

- Background
- 2 Regularization and the Energy Landscape
- In the second second
- 4 Gaussian and Heavy-tailed Random Matrix Theory
- 5 More detailed empirical results
- 6 An RMT-based Theory for Deep Learning
- 🕜 Tikhonov regularization versus Heavy-tailed regularization
- 8 Varying the Batch Size: Explaining the Generalization Gap
- Output Description Using the Theory
 - More General Implications



Conclusions: "pip install weightwatcher"

Main Empirical Results:

- Small/old NNs: Tikhonov-like self-regularization
- Modern DNNs: Heavy-tailed self-regularization

Main Modeling Results: $\mathbf{W} \simeq \mathbf{W}^{rand} + \Delta^{sig}$:

- Small/old NNs: model "noise" W^{rand} with Gaussian random matrices
- Modern DNNs: model strongly-correlated "signal" Δ^{sig} with Heavy-tailed random matrices

Main Theoretical Results: Use Heavy-tailed RMT to:

- Using global bulk stats and local edge stats, construct a operational/phenomenological theory of DNN learning
- Hypothesize 5+1 phases of learning

Evaluating the Theory:

- Effect of implicit versus explicit regularization
- Exhibit all 5+1 phases by adjusting batch size: explain the generalization gap

Main Methodological Contribution:

• Observations \rightarrow Hypotheses \rightarrow Build a Theory \rightarrow Test the Theory.

Many Implications:

• E.g., justify claims about rugged convexity of Energy Landscape

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
 B
 A
 B
 A
 B
 A
 B
 A

If you want more ... "pip install weightwatcher" ...

Background paper:

 Rethinking generalization requires revisiting old ideas: statistical mechanics approaches and complex learning behavior (https://arxiv.org/abs/1710.09553)

Main paper (full):

- Implicit Self-Regularization in Deep Neural Networks: Evidence from Random Matrix Theory and Implications for Learning (https://arxiv.org/abs/1810.01075)
- Code: https://github.com/CalculatedContent/ImplicitSelfRegularization

Main paper (abridged):

- Traditional and Heavy-Tailed Self Regularization in Neural Network Models (https://arxiv.org/abs/1901.08276)
- Code: https://github.com/CalculatedContent/ImplicitSelfRegularization

Applying the theory paper:

- Heavy-Tailed Universality Predicts Trends in Test Accuracies for Very Large Pre-Trained Deep Neural Networks (https://arxiv.org/abs/1901.08278)
- Code: https://github.com/CalculatedContent/PredictingTestAccuracies
- https://github.com/CalculatedContent/WeightWatcher
- "pip install weightwatcher"

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの