

Chapter 24

Neural networks: from the perceptron to deep nets

Marylou Gabrié

**Department of Computing Sciences, Bocconi University, Milano, Italy*

†Bocconi Institute for Data Science and Analytics, Milano, Italy

Surya Ganguli

Department of Applied Physics, Stanford University, Stanford, USA

Carlo Lucibello and Riccardo Zecchina

Department of Computing Sciences, Bocconi University, Milano, Italy

Bocconi Institute for Data Science and Analytics, Milano, Italy

Artificial networks have been studied through the prism of statistical mechanics as disordered systems since the 80s, starting from the simple models of Hopfield's associative memory and the single-neuron perceptron classifier. Assuming data is generated by a teacher model, asymptotic generalisation predictions were originally derived using the replica method and the online learning dynamics has been described in the large system limit. In this chapter, we review the key original ideas of this literature along with their heritage in the ongoing quest to understand the efficiency of modern deep learning algorithms. One goal of current and future research is to characterize the bias of the learning algorithms toward well-generalising minima in a complex overparametrized loss landscapes with many solutions perfectly interpolating the training data. Works on perceptrons, two-layer committee machines and kernel-like learning machines shed light on these benefits of overparametrization. Another goal is to understand the

advantage of depth while models now commonly feature tens or hundreds of layers. If replica computations apparently fall short in describing general deep neural networks learning, studies of simplified linear or untrained models, as well as the derivation of scaling laws provide the first elements of answers.

24.1 Statistical physics approaches to learning problems

The replica method of spin glass theory has been successfully applied since the 80s to characterize the computational capacity of simple neural network models, starting with Hopfield's associative memory model and the simplest of neural classifiers, the perceptron. From the perspective of spin glass physics, the quenched disorder is given by the data, the patterns to be stored or classified. In the theoretical analysis of supervised learning scenarios, the data are often produced by a generative process, the teacher, of which the learning system, the student, may have prior statistical information. These studies concerned mainly the generalization error, that is, the expected error produced by a trained network when a previously unseen data item is presented. The results are derived in the asymptotic regime, in which the number of data points, the dimension of the input and the number of parameters are sent to infinity while maintaining a sensible scaling relationship between them. The statistical mechanics studies of the 80s and 90s have led to important advances, providing results on the storage and generalization capacity of nontrivial systems that were out of reach for mathematically rigorous techniques.

As far as the dynamics of learning processes is concerned, the contribution were mainly limited to the so called online setting, where the patterns are presented only once and the learning dynamics can be described in the continuous time limit by a set of differential equations. When multiple passes over the dataset are involved instead, gradient descent and stochastic gradient descent on the perceptron model have been analyzed using the much more complicated set of equations given by dynamical mean field theory.

In the conceptual framework of the 1980s-1990s, the role of replica symmetry breaking was key. It allowed to establish the limit of the learning capacity of non-convex classifiers and shed light on the geometrical structure of the loss landscape.

On the algorithmic side, the cavity method has led to the design of an efficient message-passing algorithm that can be used to obtain Bayesian or

maximum likelihood predictions in simple neural architectures for which the method is exact (perceptrons, tree-like networks). The fixed points of the algorithm are related to the stationary points of the replica free energy, and the dynamics can be tracked statistically with a simple set of scalar equations called state evolution. More complicated, and of much greater generality and relevance, is the analysis of algorithms such as gradient descent and stochastic gradient descent in non-convex landscapes and their link with good generalization. This is a complex challenge which requires understanding the interplay between complex algorithmic dynamics and the geometry of the learning landscape. The out-of-equilibrium dynamics of algorithms that implement learning processes represents an open conceptual challenge that is only minimally understood and that appears to be essential for a thorough understanding of contemporary neural systems. Effective algorithms do not uniformly sample the solution space but seem to be biased toward configurations that generalized well.

24.1.1 *The storage problem*

The modern application of statistical physics to artificial neural networks had its origins in 1982, with the seminal introduction of the Hopfield model (HM) [Hopfield (1982)]. The HM was created as a toy biological model of an associative memory, whose goal is to store P binary configurations, called *memories* or *patterns*, that represent the firing or non-firing state of N neurons. The prescribed memories $\mathbf{x}^\mu \in \{-1, +1\}^N$, for $\mu = 1, \dots, P$, are by definition successfully stored if the neural network dynamics has a fixed point very close to each pattern. In particular, Hopfield modeled the neural network dynamics as the zero temperature, greedy MCMC dynamics of an Ising spin system: $\sigma_i^{t+1} = \text{sign}(\sum_{j \neq i} J_{ij} \sigma_j^t)$. The synaptic connectivity matrix J is chosen to store the memories via the Hebb rule [Hebb (1949)]: $J_{ij} = \sum_{\mu=1}^P \xi_i^\mu \xi_j^\mu$. With a simple signal-to-noise argument, in Ref. [Hopfield (1982)] it is argued that network can store up to $P \approx 0.14N$ random patterns, yielding an extensive memory capacity proportional to the number of neurons.

Note that the Hopfield dynamics corresponds to minimizing the energy function $E(\sigma) = -\frac{1}{N} \sum_{i < j} \sum_{\mu=1}^P \xi_i^\mu \xi_j^\mu \sigma_i \sigma_j$. Amit, Gutfreund and Sompolinsky [Amit *et al.* (1987)] then precisely characterized the phase diagram of the statistical mechanics system associated with this energy function using replica theory, within a replica symmetric (RS) ansatz. In the temperature T vs memory load $\alpha = P/N$ plane, they found at low α

and low T a ferromagnetic region where the memories correspond to stable states that dominate the equilibrium Gibbs state. These states were called retrieval states. At higher α and low T , the retrieval states still exist as metastable states, but the equilibrium is dominated by mixtures of a finite number of patterns. Finally, for high α and low T there is a pure spin glass phase, while at high T there is a paramagnetic phase. In both cases the retrieval states are no longer present. Further analysis showing the existence of a replica symmetry breaking (RSB) instability and the application of the 1RSB formalism, only slightly refines the RS estimates [Crisanti *et al.* (1986)].

The Hebb rule, while neurobiologically motivated, is however only one of many possible ways to store the P memories. This raises the natural question of what might be the highest possible storage capacity over all possible choices of connectivity J . The answer to this question, under the statistical assumption of random and independently generated memory patterns, was given by the seminal calculations of Gardner and Derrida for continuous synaptic strengths [Gardner and Derrida (1988); Gardner (1988)] and Mézard and Krauth for discrete ones [Mezard (1989)]. The foundational idea was to consider the space of all possible connectivity matrices J that are consistent with the memory storage fixed point conditions $\xi_i^\mu = \text{sign}(\sum_{j \neq i} J_{ij} \xi_j^\mu)$ for every memory μ . Each memory therefore imposes a constraint on J . Moreover, these constraints decouple, or are independent, across every row of J . Now a single row of J corresponds to the incoming synaptic weights onto a single neuron. Thus the problem of calculating the storage capacity of an associative memory with pairwise interactions reduces to that of a single neuron (perceptron).

Let \mathbf{w} denote vector of synaptic weights onto any one neuron, corresponding to some row of J . The approach pioneered by Gardner was to compute the volume of allowed synaptic weight configurations consistent with the storage of a dataset of P examples, $\mathcal{D} = \{\mathbf{x}^\mu, y^\mu\}_{\mu=1}^P$. This volume can be computed via the partition function

$$Z_{\mathcal{D}} = \int dP(\mathbf{w}) \prod_{\mu=1}^P \Theta \left(\frac{1}{\sqrt{N}} y^\mu \sum_{i=1}^N w_i x_i^\mu \right). \quad (24.1)$$

Here $\Theta(x)$ is the Heaviside function, $\Theta(x) = 1$ if $x > 0$ and 0 otherwise and $P(\mathbf{w})$ is the uniform measure on the set of allowed perceptron weights (for continuous weights this is the hypersphere $\sum_i w_i^2 = N$, while for discrete weights this is the hypercube $\{-1, +1\}^N$). We are interested in the *typical* volume in the high dimensional-limit $P, N \rightarrow \infty$ with finite $\alpha = P/N$.

Therefore we consider the average *entropy*

$$S(\alpha) = \lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E}_{\mathcal{D}} \log Z_{\mathcal{D}}. \quad (24.2)$$

The expectation is over i.i.d. standard Gaussian inputs x_i^μ and over uniform i.i.d. $y^\mu \in \{-1, +1\}$ that are also independent of the inputs. As one requires storage of more patterns by increasing α , the entropy and typical volume both decrease. Importantly, at a critical capacity α_c , the volume of the solution space shrinks to zero, indicating more patterns cannot be typically stored for *any* choice of perceptron weights. For continuous spherical weights, a RS calculation gives $\alpha_c = 2$ and $\lim_{\alpha \rightarrow 2^-} S(\alpha) = -\infty$. For binary weights instead $\alpha_c \approx 0.83$, as can be obtained from the condition $S(\alpha_c) = 0$. See [Engel and Van den Broeck (2001)] for an extensive discussion of the storage problem.

Some generalizations of the Hopfield model achieve super-extensive capacity: they are able to store a number of patterns polynomial [Gardner (1987); Krotov and Hopfield (2016); Albanese *et al.* (2022)] or even exponential [Demircigil *et al.* (2017)] in the size N of the system. That is also true when continuous variables and memories are involved, as in the case of the modern Hopfield network of Ref. [Ramsauer *et al.* (2021)], a model linked to the wildly popular transformer architecture for deep learning [Vaswani *et al.* (2017)].

24.1.2 Teacher-student scenarios

From the theoretical analysis of the storage limits of neural networks, we move to learning problems, where the main interest is in characterizing the behavior of the generalization error. The theoretical analysis requires the statistical definition of: 1) a data generating process; 2) the model whose parameters have to be learned from the data. This framework is referred as the *teacher-student model* in the literature. In the simplest scenario, both the teacher and the student are perceptrons, characterized by weight vectors \mathbf{w}^* and \mathbf{w} respectively. In an ideal Bayesian framework, the learner has access to the probability distribution from which the teacher weights \mathbf{w}^* are generated and to the likelihood of producing a certain label y given \mathbf{w}^* and an input \mathbf{x} . The data generating process is the following:

$$\mathbf{w}^* \sim P_W, \quad (24.3)$$

$$\mathbf{x}^\mu \sim P_X, \quad (24.4)$$

$$y^\mu \sim P\left(y \mid \frac{1}{\sqrt{N}} \sum_i w_i^* x_i^\mu\right) \quad \mu = 1, \dots, P. \quad (24.5)$$

In this Bayesian setting, the associated free energy is the log-normalization factor of the posterior distribution of the weights:

$$\phi = \lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E}_{\mathcal{D}} \log \int d\mathbf{w} P(\mathbf{w}) P(\mathcal{D} | \mathbf{w}). \quad (24.6)$$

This can be computed using the replica method. The calculation involves the introduction of a $n \times n$ overlap matrix q_{ab} for the student and of a student-teacher overlap vector r_a :

$$q_{ab} = \frac{1}{N} \sum_i w_i^a w_i^b; \quad r_a = \frac{1}{N} \sum_i w_i^a w_i^*. \quad (24.7)$$

One then proceeds with a replica symmetry ansatz for these order parameters and sends the number of replicas n to 0 as usual. In this optimal Bayesian setting, where the student is statistically matched to the teacher, the Replica Symmetric ansatz is the correct one, thanks to the Nishimori condition [Nishimori (1980); Iba (1999); Zdeborová and Krzakala (2016)]. The free energy is simply expressed in terms of a few scalar integrals and obtained through saddle point evaluation of the order parameters.

The expected generalization error is defined as the expectation (over the realizations of the dataset, the test example (\mathbf{x}, y) , and the prediction $\hat{y}(\mathbf{x})$ from the model) of a cost function $c(\hat{y}, y)$ comparing the true target and the prediction:

$$\mathcal{E}_{gen} = \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}, y} \mathbb{E}_{\hat{y} | \mathbf{x}, \mathcal{D}} c(\hat{y}, y) \quad (24.8)$$

In classification tasks, the cost function is typically the 0-1 valued error-counting function, while in regression instead it is the mean square error. Crucially, the expected generalization error for the Bayesian prediction $P(\hat{y} | \mathbf{x}, \mathcal{D}) = \int d\mathbf{w} P(\hat{y} | \mathbf{w}) P(\mathbf{w} | \mathcal{D})$ can be simply expressed in terms of the saddle point order parameters. The whole RS replica picture for the Bayesian-optimal perceptron model has been rigorously established in Ref. [Barbier *et al.* (2018)].

In Section 24.2 we discuss the generalization of this approach to some simple models of two-layers multi-layer perceptrons [Monasson and Zecchina (1995); Schwarze and Hertz (1993)] and other simple models. Reaching out to more complex teacher and student architectures, e.g. deeper and with proportional widths for all layers, is a major challenge for the statistical physics analysis of neural networks.

24.1.3 Message passing algorithms

24.1.3.1 Belief Propagation equations

The Belief Propagation (BP) algorithm [Yedidia *et al.* (2002); Mézard and Montanari (2009)] is a message passing algorithm for computing marginal and free energies in sparse graphical models. Prominent applications are in coding [Kabashima and Saad (1998)] and combinatorial optimization [Krzakala *et al.* (2007)] among others. BP has been used to efficiently solve the problem of training a perceptron with binary weights [Braunstein and Zecchina (2006)] where the iterations of the node-to-factor and factor-to-node messages take the form:

$$m_{i \rightarrow \mu}^{(t+1)} = \tanh \left(\sum_{\nu \setminus \mu} \xi_i^\nu \hat{m}_{\nu \rightarrow i}^{(t)} \right) \quad (24.9)$$

$$\hat{m}_{\mu \rightarrow i}^{(t+1)} = f \left(\sum_{j \setminus i} \xi_j^\mu m_{j \rightarrow \mu}^{(t+1)}, \sum_j (m_{j \rightarrow \mu}^{(t+1)})^2 \right) \quad (24.10)$$

with f a simple function. With respect to the standard BP prescription, the equations here have been simplified exploiting the central limit theorem since the factor graph is dense. This approximation goes under the name of relaxed Belief Propagation [Mézard (2017)]. Fixed points of the algorithm give the estimated marginals and can be used to make a Bayesian prediction for a given input. In order to produce a single binary configuration instead, one has to apply a decimation or reinforcement heuristic [Braunstein and Zecchina (2006)] on top of BP.

24.1.3.2 Approximate Message Passing and State Evolution

It turns out that on dense graphical models such as the perceptron and under certain statistical assumptions on the disorder, the BP equations can be further simplified in what is known as the Approximate Message Passing algorithm (AMP). AMP was first proposed in a seminal paper by Donoho and Montanari [Donoho *et al.* (2009)] building on the rigorous analysis of Bolthausen of the TAP equations for the SK model [Bolthausen (2014)]. Compared to BP, AMP has lower memory complexity since it involves the computation of only node-related quantities instead of edge ones. Moreover, in the high-dimension limit, the statistics of the AMP messages can be rigorously tracked by a dynamical system involving only a few scalar quantities, known as State Evolution (SE). Remarkably, SE

involves the same quantities appearing in the RS replica calculation, and its fixed points correspond to the stationary points of the RS replica free energy [Barbier *et al.* (2018)]. Due to this connection, AMP has also been used as a proof technique for replica results in convex models [Loureiro *et al.* (2022a)]. See also [Advani and Ganguli (2016b,a)] for both replica and AMP perspectives on deriving optimal loss functions and regularizers for high dimensional regression.

While Ref. [Donoho (2006)] analyzed linear models, AMP was later extended to generalized linear models (where it is called GAMP) [Rangan (2011)], and committee machines with few hidden nodes (see Section 24.2).

In inference settings, AMP has been applied to multi-layer architectures with an extensive number of hidden nodes but fixed weights [Manoel *et al.* (2017); Fletcher *et al.* (2018)]. The deep learning setting is much more challenging though, and there have been only limited attempts so far [Lucibello *et al.* (2022)]. Being able to scale message passing to deep learning scenarios would allow to perform approximate Bayesian estimation, train discrete weights for computational efficiency and energy saving, have analytically trackable algorithms, and perform cheap hyperparameter selection.

24.1.4 *The geometry of the solution space*

Learning in Neural Networks (NNs) is in principle a difficult computational task: a non-convex optimization problem on a huge number of parameters. However, the problem seems to be relatively easy to solve, as even simple gradient-based algorithms convergence to solutions with good generalization capabilities. NNs models are evolving rapidly through a collective effort shared across many labs. It is thus difficult to define a unifying theoretical framework, and current NNs are in a sense similar to complex, highly evolved natural systems.

A major question in deep learning concerns understanding the non-convex geometry of the error landscape as a function of the parameters, and how this geometry might facilitate gradient based learning. Motivated by the geometry of random Gaussian landscapes [Bray and Dean (2007); Fyodorov and Williams (2007)] (derived via replica theoretic methods), Ref. [Dauphin *et al.* (2014)] numerically explored the statistics of extrema of the error landscape of deep and recurrent networks, finding that higher error extrema were typically higher index saddle points, not local minima. Thus high error local minima need not confound deep learning, though low index saddle points might, and [Dauphin *et al.* (2014)] developed an al-

gorithm to rapidly escape these saddle points. The authors of [Baity-Jesi *et al.* (2018)] undertook a careful comparison of the dynamics of stochastic gradient descent on neural networks versus the p-spin spherical spin glass energy function, finding interesting ageing phenomena indicative of the prevalence of more flat directions as once descends the training error. Ref. [Geiger *et al.* (2019)] found an interesting analogy between jamming and the error landscape of deep networks with a hinge loss, building on a prior analogy for the perceptron [Franz and Parisi (2016)]. As the network size transitions from overparameterized (with many weight configurations at zero error) to underparameterized (with many isolated minima), the error landscape undergoes a jamming transition. Finally Ref. [Maillard *et al.* (2020)] provided another interesting exploration by extending the Kac-Rice method to count critical points in generalized linear models.

A complementary view comes from some recent works that focus on the role played by rare attractive minima. At least for non-convex shallow neural networks classifying random patterns, it is possible to derive a theoretical description of the geometry of the zero training error configurations (so-called "solutions"). The most immediate result is that the solutions that dominate the zero-temperature Gibbs measure of the error loss (i.e., the most numerous) do not match those found by the efficient learning algorithms. Numerical evidence suggests that in fact the solutions found by the algorithms belong to particularly entropic regions, i.e., with a high density of other nearby solutions [Baldassi *et al.* (2015, 2016a)]. These types of solutions are often referred to as flat minima in the machine learning literature.

In deep learning settings, numerical results consistently show that flatness of a minimizer positively correlates with generalization ability (see e.g. [Jiang *et al.* (2020)]). Different algorithms explicitly targeting flat minima have been proposed [Chaudhari *et al.* (2016); Pittorino *et al.* (2021); Foret *et al.* (2021)].

The analytical study of flat minima and their generalization properties can be done by resorting to a large deviation technique introduced in Ref. [Baldassi *et al.* (2015)] and based on the Franz-Parisi potential [Franz and Parisi (1995)]. In order to introduce this framework, name Local Entropy (LE), it is useful to consider the simplest model displaying a rich geometry of the solution space, i.e. the binary perceptron. The LE framework can be applied to more complex architectures and continuous variables as well. The local entropy function is defined as the (normalized) logarithm of the number of solutions \mathbf{w}' at some intensive distance d from a reference

solution \mathbf{w} :

$$S_{LE}(d, \mathbf{w}) = \frac{1}{N} \ln \sum_{\mathbf{w}'} \mathbb{X}(\mathbf{w}') \delta(d_H(\mathbf{w}', \mathbf{w}) - dN) \quad (24.11)$$

where d_H is the Hamming distance and \mathbb{X} is the data-dependent indicator function for the solutions. The analysis of $S(d, \mathbf{w})$ for uniformly sampled solutions to the training set reveals that *typical solutions are isolated* [Huang *et al.* (2013)], meaning that no near solutions exist at small d . This is the analogous of sharp minima in continuous networks. It turns out that by sampling solutions according to the LE itself, that is from

$$P_{LE}(d, \mathbf{w}) \propto e^{yN S_{LE}(d, \mathbf{w})}, \quad (24.12)$$

it is possible to uncover the existence of *rare* (according to the flat measure) regions containing an exponential number of solutions (a large volume in the case of continuous variables). By analogy, these are the flat minima found in continuous and deep architectures. In the last equation y has the role of an inverse temperature conjugated to the local entropy. For large values of y the probability focuses on the \mathbf{w} which are surrounded by an exponential number of solutions at distance d , therefore suppressing isolated solutions when d is small enough. Fig. 24.1 displays the coexistence of isolated and dense solutions up to a certain value of α when the dense region (in blue) disappears. In Ref. [Baldassi *et al.* (2015)] it is analytically shown that dense solutions generalize better than isolated ones. It must be also noted that isolated solutions cannot be algorithmically accessed by polynomial algorithms. Rigorous results on the geometry of the symmetric variant of the binary perceptron have been obtained in Ref. [Abbe *et al.* (2022)].

Adapting the LE definition to generic architectures and continuous weights, one can try direct optimization of S_{LE} instead of the original loss. Estimation and optimization can be done within a double-loop algorithm known as Entropy-SGD [Chaudhari *et al.* (2016); Pittorino *et al.* (2021)]. Second order approximation of S_{LE} leads to the Sharpness Aware Minimization algorithm of Ref. [Foret *et al.* (2021)].

Other algorithmic approaches for the generic deep learning setting are obtained as follows [Baldassi *et al.* (2016a)]. Starting from Eq. (24.12), relax the constraint on the distance with a Lagrange multiplier γ , replace the indicator function in Eq. (24.11) with a Boltzmann weight involving a loss $\mathcal{L}(\mathbf{w}')$, take integer y , and unfold the exponential in Eq. (24.12) by

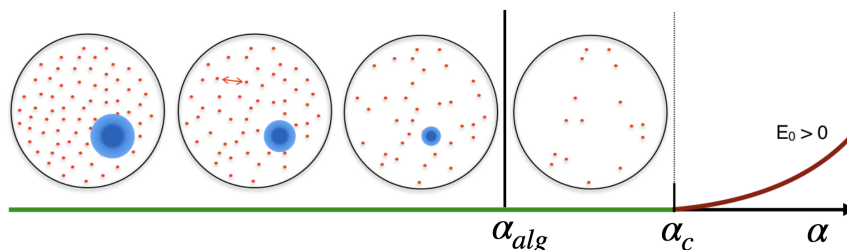


Fig. 24.1 Geometrical transitions in the solution space for the storage problem of the binary perceptron at increasing load $\alpha = P/N$. Finding a solution in the connected cluster is easy, while finding a solution after the cluster disappears at $\alpha = \alpha_{alg}$ becomes computationally hard since all of them are isolated. Above $\alpha = \alpha_c$ no zero energy configurations exist.

introducing y replicas of the original system. We end up with a replicated system with an effective energy

$$\mathcal{L}_R(\mathbf{w}, \{\mathbf{w}'^a\}_a) = \sum_{a=1}^y \mathcal{L}(\mathbf{w}'^a) + \gamma \sum_{a=1}^y d(\mathbf{w}, \mathbf{w}'^a). \quad (24.13)$$

We thus have a central replica interacting with y peripheral ones that also resent from the data-dependent loss function. As discussed in ref. [Baldassi *et al.* (2016a,b, 2020)], several algorithmic schemes can be derived in a straightforward way by optimizing the replicated loss function \mathcal{L}_R , for example with gradient-based algorithms, or by sampling the space of solutions with a Markovian process or with Belief Propagation equations. One of these, the flat-minima-seeking algorithm known as rSGD, has been applied to deep neural architectures in Ref. [Pittorino *et al.* (2021)].

24.1.5 Learning dynamics

With neural network models, the dynamics of gradient descent takes place in a high-dimensional non-convex landscape, therefore its theoretical description is highly non-trivial. For the perceptron and in the continuous time limit, dynamical mean-field theory has been used to provided such description in terms of low-dimensional integro-differential equations involving two-times correlations functions and responses [Agoritsas *et al.* (2018)]. Such description has been extended to stochastic gradient descent as well [Mignacco *et al.* (2020)]. Numerical solution of the equations is particularly challenging, and the framework has not been extended to deeper architectures so far.

In the simpler online setting, where only a single pass over the data points is allowed, a much simpler set of ODE can be derived. For the committee machine architecture detailed in Section 24.2 the quantities to be evolved are the student-student and teacher-student overlaps among the hidden perceptrons, $Q_{k,k'}$ and R_{k,k^*} respectively [Saad and Solla (1995a,b); Saad (1999); Goldt *et al.* (2019a)]. In terms of these overlaps, the evolution of the generalization error can be described at each time step. In the same setting, but assuming infinitely wide networks with finite input size, the authors of Ref. [Mei *et al.* (2018)] obtained a mean-field description in terms of a PDE characterizing a diffusion process for the perceptrons' weights.

24.2 Studying over-parametrized models

The repeated breakthroughs of deep learning starting from the 2010's propelled neural networks to the forefront of machine learning [LeCun *et al.* (2015)]. Moreover these neural networks increased rapidly in size; one of the first convolutional networks for document recognition, LeNet-5, had around 61,000 trainable parameters [LeCun *et al.* (1998)], while today's GPT-3 has around 175 billion [Brown *et al.* (2020)]. This proliferation of trainable parameters leads to highly flexible models, which begs the question of why they don't overfit to their training data and why they can still generalize well to new input examples. Explaining their success in this so-called *over-parametrized* regime constitutes a key theoretical question in machine learning, as it seemingly violates the classical picture of the bias-variance trade-off [Spigler *et al.* (2019); Belkin *et al.* (2019)]. In this section, we describe how over-parametrization has been tackled by the statistical mechanics approach.

24.2.1 *Committee machines*

Gardner's program for the perceptron, described above, was soon extended to multi-layer networks. The first analyzed model consisted of summing the outputs of multiple perceptrons with non-overlapping inputs [Mato and Parga (1992)]. This tree-like architecture was called *committee machine*. Soon after, this analysis was extended to full connectivity to the input [Schwarze and Hertz (1993); Schwarze (1993)]. These architectures are cases of a generic one-hidden layer neural-net, which passes a linear combination of the outputs of K perceptrons, each with non-linearity $\sigma(\cdot)$, through an

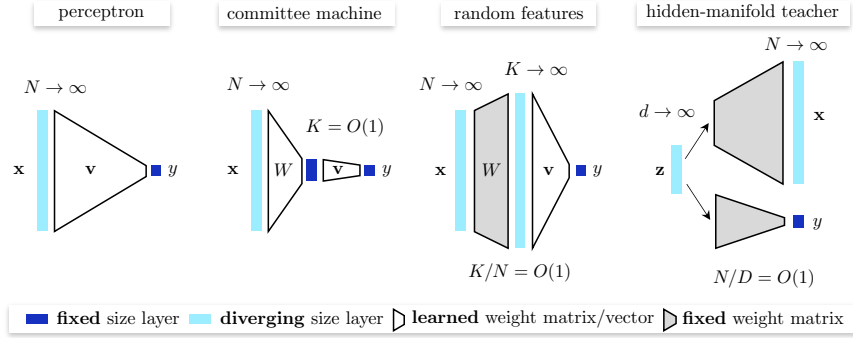


Fig. 24.2 Scaling comparison of models' dimensions in the thermodynamic limit.

output activation function $f(\cdot)$, yielding a final output y given by

$$y = f \left(\sum_{k=1}^K v_k \sigma \left(\sum_{i=1}^N W_{ki} x_i \right) \right). \quad (24.14)$$

The input to hidden layers weights are parameterized by a matrix $W \in \mathbb{R}^{K \times N}$ and the hidden to output weights are parameterized by a vector $\mathbf{v} \in \mathbb{R}^K$. K denotes the number of *hidden* neurons.

Remarkably, any smooth function on \mathbb{R}^N can be arbitrarily well approximated by (24.14) with a finite, but possibly large, K [Hornik (1991)]. This universal approximation theorem motivates the one-hidden layer neural network as a simple yet far-reaching learning model to be studied. Also, these one hidden layer networks allow a study of over-parametrization by considering large student networks with K hidden units learning from data generated by small teacher networks with $M < K$ hidden units.

Scaling and specialization transition – The statistical mechanics analysis of committee machines in the teacher-student scenario mirrors that of the perceptron. Given a dataset $\mathcal{D} = \{\mathbf{x}_\mu, y_\mu\}_{\mu=1}^P$ generated by a teacher committee machine with M hidden units of the form

$$y^\mu = f^* \left(\sum_{k=1}^M v_k^* \sigma \left(\sum_{i=1}^N W_{ki}^* x_i^\mu \right) \right), \quad (24.15)$$

one can derive learning curves in the usual high dimensional limit $P, N \rightarrow \infty$ with $\alpha = P/N$ held to be $O(1)$. Also the number of hidden units in the teacher (M) and student (K) are both kept $O(1)$ (see Figure 24.2). The

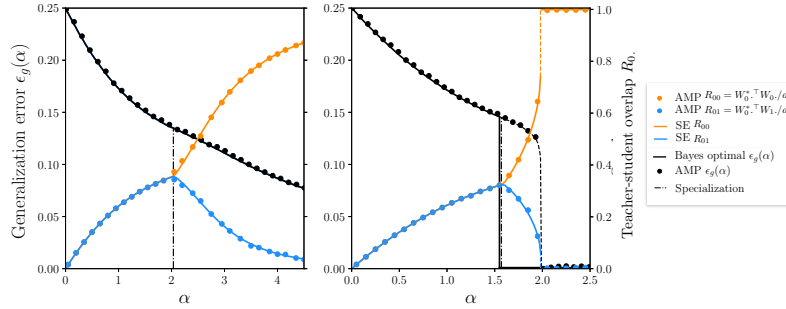


Fig. 24.3 Adapted from [Aubin *et al.* (2018)]. Overlap parameters and generalization for matched teacher-student committee machines with $M = K = 2$ hidden units as a function of the ratio $\alpha = P/N$ between sample size and input dimension, with first layer weights either Gaussian (left) or binary (right). As α grows beyond a critical value, the overlaps account for the specialization of the student hidden units to the teacher's.

order parameters are the teacher-student and student-student overlaps, as for the perceptron, except now they become matrices for the input layer:

$$R = \frac{WW^{*,\top}}{N} \in \mathbb{R}^{K \times M}, Q = \frac{WW^\top}{N} \in \mathbb{R}^{K \times K}. \quad (24.16)$$

For simplicity we assume the teacher output $\mathbf{v}^* \in \mathbb{R}^K$ to be the all 1 vector.

A remarkable phenomenology was identified using annealed and quenched replica computations [Schwarze and Hertz (1993); Schwarze (1993)] as well as an online learning analysis [Saad and Solla (1995a,b); Biehl and Schwarze (1995)]. Considering matched teacher and students ($M = K$) with fixed student output weights $\mathbf{v} = \mathbf{v}^*$, the Bayes optimal student weights were shown to specialize to the teacher weights only if a critical amount of data was available (i.e. $\alpha > \alpha_c$, see Figure 24.3). This specialization transition was derived rigorously recently in [Aubin *et al.* (2018)], through AMP and SE. Non-linear hidden units are necessary for specialization to set in, and in scarce data regimes ($\alpha < \alpha_c$) committee machines act no differently than linear models such as perceptrons.

Denosing solution – The online analysis of the specialization transition was also made rigorous recently by [Goldt *et al.* (2019a)]. In this work, Goldt and co-authors further extended the analysis to include learnable student output weights \mathbf{v} and uncovered one mechanism of improved generalization by over-parametrization. Focusing on committee machines with sigmoidal non-linearities, they showed that the generalization error

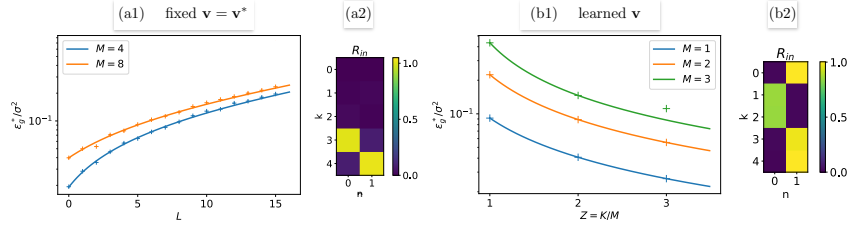


Fig. 24.4 Adapted from [Goldt *et al.* (2019a)]. (a1) and (b1) Generalization error as a function of the over-parametrization for sigmoidal committee machines with K hidden units learning from data generated with teachers with M hidden units ($L = K - M$). Generalization deteriorates with over-parametrization when the output weight vector \mathbf{v} is fixed (a1) and improves when \mathbf{v} is learned (b1). This behavior is traced back to the specialization of multiple hidden units for each of the teacher’s hidden unit in this second case only as testified by the reported teacher-student overlap matrices (a2) and (b2) for $M = 2$ and $K = 5$.

increases (decreases) with overparameterization $L = K - M$ if the output layer weight \mathbf{v} is fixed (learned) (see Figure 24.4). The natural order parameters, or overlap matrices R and Q , once again, elucidate the mechanism behind this phenomenology. With learnable output weights, several of the K hidden units of the student can (potentially weakly) specialize to one of the M teacher hidden units, and the student output weights can learn to combine, or denoise, these student hidden-unit outputs to mimic the overall teacher output. This *denoising* solution allows the student with more hidden units to build a larger ensemble of regressors of each of the teacher’s hidden units, thereby improving its accuracy with overparameterization. Conversely, when the student output weight \mathbf{v} is fixed, only M among the K student hidden units can effectively participate in the student committee, each specializing to one of the teacher’s units. The remaining $L = K - M$ units only contribute noise, so overparameterization hurts generalization.

24.2.2 Kernel-like learning

The analysis of committee machines revealed a fundamental mechanism of generalization through over-parametrization, yet, these models operate in an atypical regime where the input dimension diverges while the number of hidden units is fixed. Operating in a different scaling, kernel methods form another class of models that can be analyzed in depth (see Chap 16 [Shalev-Shwartz and Ben-David (2014)]). Kernel-based learners are non-

linear in the inputs but linear in the parameters. While deep neural networks are in essence more expressive than kernel methods¹, they sometimes behave not much differently from kernels, in particular in cases where they have infinitely wide hidden layers [Jacot *et al.* (2018)] (though see [Fort *et al.* (2020)] for empirical differences between kernel learners and deep networks at more natural widths). Moreover kernel models already exhibit key features to be understood in deep learning. Notably, the generalization performance of kernel learners undergoes a *double descent* [Advani *et al.* (2020); Spigler *et al.* (2019); Belkin *et al.* (2019)]: after an overfitting peak predicted by the classical bias-variance compromise, the generalization improves continuously as the number of parameters increases. Therefore, kernel methods also enjoy good generalization with over-parametrization.

Random features and Gaussian equivalence – The random feature model [Rahimi and Recht (2017)], is a simple kernel-like model that is closely related to the one-hidden layer model defined above in (24.14). It only differs by fixing first-layer weights W to random values from a given distribution, such as a Gaussian or select Fourier modes at random frequencies. The only learnable parameters are in the output weight $\mathbf{v} \in \mathbb{R}^K$, and the flexibility of the model can be controlled by the number of hidden units K (see Figure 24.2). In the scaling limit where K and the sample size P scale linearly with the dimension of input N , the generalization error can be characterized asymptotically through different methods. With the additional assumption that the input data is Gaussian distributed and the teacher input-output rule is a perceptron: $y^\mu = f^*(\mathbf{v}^{*\top} \mathbf{x}^\mu)$, learning curves were derived using random matrix theory [Mei and Montanari (2021)], replica computations [Gerace *et al.* (2020); D’Ascoli *et al.* (2021)] and Gaussian convex inequalities [Dhifallah and Lu (2020)]. An important insight common to these analyses is the equivalence of the non-linear features with a Gaussian model with matching moments. That is the generalization error averaged over the teacher data distribution

$$E_g = \mathbb{E}_{\mathbf{x}, y} \left[\left(f(\mathbf{v}^\top \sigma(W \mathbf{x})) - f^*(\mathbf{v}^{*\top} \mathbf{x}^\mu) \right)^2 \right], \quad (24.17)$$

concentrates in the thermodynamic limit around a function of the ratios $\alpha = P/N$ and $\gamma = K/N$ given by

$$\mathcal{E}_g(\alpha, \gamma) = \lim_{N \rightarrow \infty} E_g = \int_{\mathbb{R}^2} (f(\lambda) - f^*(\nu))^2 \mathcal{N}(\lambda, \gamma; 0, \Sigma_{\alpha, \gamma}) d\lambda d\nu, \quad (24.18)$$

¹Kernel models are not universal approximators.

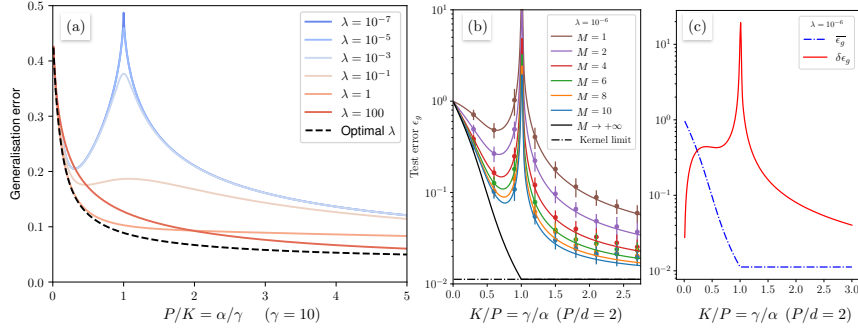


Fig. 24.5 The double descent analyzed in random-feature models through replica computation of the generalization error. (a) As an ℓ_2 regularization of strength λ is adjusted, the overfitting peak at $P/K = 1$ vanishes. (b) Similarly, averaging the prediction over an ensemble of M models mitigates the double descent pointing to the origin of the overfitting in exploding variance, as the computed decomposition shows in (c). Adapted from [Gerace *et al.* (2020)] (a) and [Loureiro *et al.* (2022b)] (b) and (c).

where the covariance of the jointly Gaussian scalars λ and ν are themselves function of the overlap order parameters of the problem and constants accounting for the choice of non-linearity σ and input distribution. This so-called *Gaussian equivalence principle* [Goldt *et al.* (2020)] was first described for square losses using rigorous random matrix theories [Pennington and Worah (2017); Mei and Montanari (2021)] and later extended to generic convex losses [Goldt *et al.* (2021); Hu and Lu (2020)]. This formulation finally allows to predict theoretically the origin of the overfitting peak in the double descent. At the interpolation threshold, where the learning model has just enough parameters to perfectly fit the training data, the variance of the learned predictor explodes [D’Ascoli *et al.* (2020); Mei and Montanari (2021); Loureiro *et al.* (2022b)]. Beyond the interpolation threshold, over-parametrized models feature many minimizers of the training loss, yet implicit or explicit regularization of the training typically selects a minimizer with good generalization (see Figure 24.5).

Universality – Beyond random features in a teacher-student scenario, kernel-learning has been analyzed in great generality using statistical mechanics tools. Loureiro and collaborators first replaced shallow random maps with learned deep feature maps [Loureiro *et al.* (2021a)], and then formulated a universal theory for any generalized linear model trained through the empirical risk minimization of a convex loss [Loureiro *et al.* (2022b)].

Interestingly, the latter result is mathematically rigorously proven through a novel proof technique relying on AMP [Gerbelot and Berthier (2021)].

Focusing instead on a generic formulation of kernel regression, Canatar and collaborators developed a widely applicable theory of learning encompassing arbitrary training data [Canatar *et al.* (2021)]. A replica computation yields the generalization error as a function of the kernel spectral decomposition, the training data distribution and the sample size. This analysis highlights the impact of the data distribution through the eigenmodes of the kernel and through the so-called *alignment* between the data and the task to be learned. Along these lines, one can show in simple linear settings that learning curves with respect to the amount of data can exhibit an arbitrary number of multiple descent peaks equal to the number of scales in the data [Mel and Ganguli (2021)]. This further highlights the importance of understanding how non-random structure in data impacts learning properties of neural networks.

24.2.3 *Studying the impact of structure in data*

The question of how and when structure in data drives generalization in neural networks is fundamental in learning theory. The sequence of above works can be viewed as the successive addition of structure in data, starting from random uncorrelated inputs and output labels in the storage problem, to random i.i.d inputs that are correlated to their corresponding outputs through the introduction of a teacher. More recent work has examined the important generalization of going beyond i.i.d. random inputs.

A first notable step is the introduction of the *hidden manifold* model assuming that the data and labels are generated from a low dimensional representation [Goldt *et al.* (2019b)] (see Figure 24.2). The analytical description of learning in committee machines, relying on the Gaussian equivalence principle, shows that generalization is tied to the dimension of the underlying manifold. These results are consistent with the analyses of random-feature models trained on anisotropic Gaussian data combining a high-variance *strong* subspace and a low-variance *weak* subspace [Ghorbani *et al.* (2020); D’Ascoli *et al.* (2021)]. The concept of data-task alignment is here again identified as a determinant of generalisation: it ensures that the effective dimension of the problem is small and allows for good generalisation at small sample size.

Another model of structure in data, especially relevant for classification problems, is the Gaussian mixture model of inputs. Limits of separability

and generalization behavior can be derived for classification of Gaussian clusters with generalized linear models [Loureiro *et al.* (2021b)] and committee machines [Refinetti *et al.* (2021)]. In particular, the latter work demonstrates that some configurations of clusters unlearnable with random features (an instance of a generalized linear model) can be separated by committee machines. The 2-layer neural network learns appropriate features to allow for a proper task-data representation alignment, whereas generalized linear models are limited by their fixed feature-maps, which may not be appropriate for the task.

24.2.4 *Non-convex overparametrized neural networks*

Finally, in ref. [Baldassi *et al.* (2022)] the computational consequences of overparameterization in non-convex neural network models are analyzed. In the simple case of discrete binary weights, a non-convex classifier is connected to a random features projection. The analysis shows that as the number of connection weights increases multiple phase transitions happen in the zero error landscape. A first transition happens at the so-called interpolation point, when solutions begin to exist (perfect fitting becomes possible). A second transition occurs with the discontinuous appearance of a different kind of “atypical” structures corresponding to high local entropy regions with good generalization properties.

24.3 Going deeper

Much of the remarkable progress in artificial intelligence over the last decade has been driven by our ability to train very deep networks with many successive nonlinear layers. Formally, the simplest version of a feedforward neural network with D layers is the multi-layer perceptron, defined through

$$\mathbf{x}^l = \phi(\mathbf{h}^l) \quad \mathbf{h}^l = W^l \mathbf{x}^{l-1} + \mathbf{b}^l \quad \text{for } l = 1, \dots, D. \quad (24.19)$$

Here, $\mathbf{x}^0 \in \mathbb{R}^{N_0}$ is the input, which propagates through the network to generate a sequence of activity vectors $\mathbf{x}^l \in \mathbb{R}^{N_l}$ in layer l with N_l neurons. W^l is an $N_l \times N_{l-1}$ weight matrix connecting neurons in layer $l-1$ to layer l , \mathbf{b}^l is a vector of biases to neurons in layer l , \mathbf{h}^l is the pattern of inputs to neurons at layer l , and ϕ is a single neuron scalar nonlinearity that acts component-wise to transform inputs \mathbf{h}^l to activities \mathbf{x}^l . The final output of the network is $\mathbf{y} = \mathbf{x}^D(\mathbf{x}^0, \mathbf{w})$ where \mathbf{w} collectively denotes all N neural network parameters $\{W^l, \mathbf{b}^l\}_{l=1}^D$.

The theoretical analysis of such deep networks in full generality raises several difficult challenges which remain open to this day [Bahri *et al.* (2020); Gabri e (2020); Maillard *et al.* (2022)]. For example, how can we describe their learning dynamics? What can deep networks express that their shallow counterparts cannot? How should we initialize the weights to optimize their learning dynamics? What does their training error landscape look like? How can we use them for generative modeling? How does their performance scale with network size? In the following we review some recent progress on these questions that involves ideas from both equilibrium and non-equilibrium statistical mechanics, nonlinear dynamical systems, random matrix theory, and free probability.

24.3.1 *Exact learning dynamics for deep linear networks*

A gold standard of understanding deep learning would be an exact solution to the learning dynamics of both training and generalization error as a function of training time, for arbitrarily deep networks and for arbitrarily structured data. While this is challenging for networks of the form in (24.19), remarkably it is possible in the case of linear networks with $\phi(x) = x$ and squared loss [Saxe *et al.* (2013); Lampinen and Ganguli (2018)]. Despite the linearity of the network, the learning dynamics is highly nonlinear. Exact solutions to these dynamics reveal that deep linear networks learn by successively approximating the singular value decomposition (SVD) of the input-output correlation matrix of the training data mode by mode. Each mode is learned on a time scale inversely related to its singular value.

Intriguingly, while deep nonlinear networks have long been used in psychology to model the developmental dynamics of semantic cognition infants [Rogers and McClelland (2004)], a recent analysis [Saxe *et al.* (2019)] showed that much of this developmental learning dynamics could be qualitatively captured by deep linear networks, thereby accounting for a diversity of phenomena, including progressive differentiation of semantics, semantic illusions, item typicality, category coherence, dynamic patterns of inductive projection, and the conservation of semantic similarity across species.

Additionally, recent advances in self-supervised learning [Chen *et al.* (2020); Grill *et al.* (2020)] yield new types of learning dynamics that can also be well modelled by deep linear networks [Tian *et al.* (2020, 2021)], even in settings where the learning dynamics do not correspond to gradient descent on any function [Grill *et al.* (2020)]. Moreover, these simpler learning models have prescriptive value: analytically derivable hyperparameter

choices that work well for training deep linear models, also work well for their highly nonlinear counterparts [Tian *et al.* (2021)], thereby opening the door to the use of mathematical analysis to drive practical design decisions.

24.3.2 Expressivity and signal propagation in random nets

Even before one trains a network, one has to choose the initial weights and biases, for example $\{W^l, \mathbf{b}^l\}_{l=1}^D$ in (24.19), and the choice of such an initialization can have a dramatic practical impact on subsequent learning dynamics. A common choice is a zero mean i.i.d. Gaussian initialization with variance σ_w^2/N_{l-1} for weights W_{ij}^l and variance σ_b^2 for biases b_i^l . This relative scaling ensures weights and biases exert similar control over a neuron in layer l as any previous layer width N_{l-1} becomes large.

In the limit of large N_l , [Poole *et al.* (2016)] analyzed the propagation of signals through (24.19) via dynamic mean field theory and found an order to chaos phase transition in the σ_w^2 by σ_b^2 plane for sigmoidal nonlinearities ϕ . In the ordered regime with small σ_w^2 relative to σ_b^2 , the network contracts nearby inputs as they propagate forward through the network and backpropagated error gradients vanish exponentially in depth. In the chaotic phase with σ_w^2 large relative to σ_b^2 , forward propagation chaotically amplifies and then folds small differences in inputs, leading to highly flexible and expressive input-output maps, while backpropagated error gradients explode exponentially in depth. Initializing networks near the edge of chaos in the σ_w^2 by σ_b^2 plane yields good, well-conditioned initializations. Indeed [Schoenholz *et al.* (2017)] found that the closer one initializes to the edge of chaos, the deeper a network one can train.

Subsequent work [Pennington *et al.* (2017, 2018)] showed that one can go beyond i.i.d. Gaussian initializations to speed up learning especially in very deep or recurrent networks. They considered the Jacobian

$$J = \frac{\partial \mathbf{x}^D}{\partial \mathbf{x}^0} = \prod_{l=1}^D D^l W^l. \quad (24.20)$$

Here D^l is a diagonal matrix with entries $D_{ij}^l = \phi'(h_i^l) \delta_{ij}$ and \mathbf{h}^l is defined in (24.19). J measures the susceptibility of the network's output to small changes in the input. Related susceptibilities play a fundamental role in the backpropagation of error that guides gradient based learning. Initializing at the edge of chaos controls the mean squared singular value of J to be 1. But for very deep networks, J could still become ill conditioned with the maximal singular value growing at a rate that is linear in the depth,

even at the edge of chaos. [Pennington *et al.* (2017, 2018)] controlled this growth by analytically computing the *entire* singular value spectrum of J using free probability, exploiting the fact that J is a product of random matrices, and then determined how and when we can ensure *dynamically isometric* initializations in which the *entire* singular value distribution can be tightly concentrated around 1. In particular, orthogonal weight matrices with sigmoidal nonlinearities achieve such dynamical isometry, which was shown to speed up training.

Building on these works [Xiao *et al.* (2018)] extended dynamically isometric initializations to convolutional neural networks, and showed how to train 10,000 layer models without any of the complex normalization tricks used in deep learning. Thus overall, studies of signal propagation and initialization in random deep networks provide an example of how statistical mechanics type analyses can guide engineering design decisions.

24.3.3 *Generative models via non-equilibrium dynamics*

A recent major advance in deep learning is the ability to generate remarkable, novel photorealistic images from text descriptions (see e.g. [Rombach *et al.* (2022); Saharia *et al.* (2022); Ramesh *et al.* (2022)]). Intriguingly one component of models that can do this well was inspired by ideas in non-equilibrium statistical mechanics [Sohl-Dickstein *et al.* (2015); Ho *et al.* (2020)]. In particular, being able to generate images requires modeling the probability distribution over natural images. Equilibrium methods for modeling complex distributions involve creating a Markov chain that obeys detailed balance with respect to the distribution. However, if the distribution has multiple modes, such Markov chains can take very long to mix.

[Sohl-Dickstein *et al.* (2015)] instead suggested training a finite time nonequilibrium stochastic process to generate images starting from noise. The method of training involved taking natural images and allowing them to diffuse in a high-dimensional pixel space, thereby destroying their structure and converting them to noise. Then a neural network was trained to reverse the flow of time in this otherwise irreversible structure-destroying diffusion process. The result is then a trained neural network that can sample images by converting any white noise image into a particular naturalistic image through an approximate reversal of the diffusion process. Thus this provides another example of how ideas in statistical mechanics can eventually lead to state of the art systems in artificial intelligence.

24.3.4 Neural scaling laws governing deep learning

Another major shift in deep learning over the last few years has been the immense scaling up of both dataset and model sizes. This has partially been motivated by empirically observed neural scaling laws [Hestness *et al.* (2017); Kaplan *et al.* (2020); Henighan *et al.* (2020); Gordon *et al.* (2021); Hernandez *et al.* (2021); Zhai *et al.* (2021); Hoffmann *et al.* (2022)] in many domains of machine learning, including vision, language, and speech, which demonstrate that test error often falls off as a power law with either the amount of training data, model size, or compute. Such power law scaling has motivated significant societal investments in data collection, compute, and associated energy consumption. However, obtaining a theoretical understanding of the origins of these scaling laws, and the dependence of their exponents on structure in data, model architecture, learning hyperparameters, etc... constitutes a major research question [Bahri *et al.* (2021)].

Also interesting is whether power law neural scaling can be beaten. Recent work [Sorscher *et al.* (2022)] explored scaling with respect to dataset size to analyze whether intelligent data pruning methods [Paul *et al.* (2021)] that involve careful subselection of training data, could lead to more efficient scaling of error w.r.t. pruned dataset sizes. Returning to the perceptron, [Sorscher *et al.* (2022)] developed a replica calculation to compute the test error for non-Gaussian pruned data, and showed it is possible to beat power law scaling, and even approach exponential scaling, provided one has access to a good data pruning metric. Then going beyond the perceptron [Sorscher *et al.* (2022)] showed how to beat power law scaling in modern architectures like ResNets trained on modern datasets, including ImageNet.

Overall, a deeper understanding of the origin of neural scaling laws, and the emergent capabilities that arise [Wei *et al.* (2022)], constitutes a major open question, which may benefit from the analysis of appropriate statistical mechanics models capturing the essence of these phenomena.

Bibliography

- Abbe, E., Li, S., and Sly, A. (2022). Binary perceptron: efficient algorithms can find solutions in a rare well-connected cluster, in *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, pp. 860–873.
- Advani, M. and Ganguli, S. (2016a). An equivalence between high dimensional bayes optimal inference and m-estimation, *Adv. Neural Inf. Process. Syst.* .
- Advani, M. and Ganguli, S. (2016b). Statistical mechanics of optimal convex inference in high dimensions, *Physical Review X* .
- Advani, M. S., Saxe, A. M., and Sompolinsky, H. (2020). High-dimensional dynamics of generalization error in neural networks, *Neural Networks* **132**, pp. 1–32, doi:10.1016/j.neunet.2020.08.022, arXiv:1710.03667, <https://doi.org/10.1016/j.neunet.2020.08.022>.
- Agoritsas, E., Biroli, G., Urbani, P., and Zamponi, F. (2018). Out-of-equilibrium dynamical mean-field equations for the perceptron model, *Journal of Physics A: Mathematical and Theoretical* **51**, 8, p. 085002, doi:10.1088/1751-8121/aaa68d, <https://doi.org/10.1088/1751-8121/aaa68d>, publisher: IOP Publishing.
- Albanese, L., Alemanno, F., Alessandrelli, A., and Barra, A. (2022). Replica Symmetry Breaking in Dense Hebbian Neural Networks, *Journal of Statistical Physics* **189**, 2, p. 24, doi:10.1007/s10955-022-02966-8, <https://doi.org/10.1007/s10955-022-02966-8>.
- Amit, D. J., Gutfreund, H., and Sompolinsky, H. (1987). Statistical mechanics of neural networks near saturation, *Ann. Phys.* **173**, 1, pp. 30–67.
- Aubin, B., Maillard, A., Barbier, J., Krzakala, F., Macris, N., and Zdeborová, L. (2018). The committee machine: Computational to statistical gaps in learning a two-layers neural network, in *Neural Information Processing Systems 2018*, NeurIPS, pp. 1–44, arXiv:1806.05451, <http://arxiv.org/abs/1806.05451>.
- Bahri, Y., Dyer, E., Kaplan, J., Lee, J., and Sharma, U. (2021). Explaining neural scaling laws, arXiv:2102.06701 [cs.LG].
- Bahri, Y., Kadmon, J., Pennington, J., Schoenholz, S. S., Sohl-Dickstein, J., and Ganguli, S. (2020). Statistical mechanics of deep learning, *Annual Review of Condensed Matter Physics* .

- Baity-Jesi, M., Sagun, L., Geiger, M., Spigler, S., Arous, G. B., Cammarota, C., LeCun, Y., Wyart, M., and Biroli, G. (2018). Comparing dynamics: Deep neural networks versus glassy systems, *International Conference on Machine Learning (ICML)*.
- Baldassi, C., Borgs, C., Chayes, J. T., Ingrosso, A., Lucibello, C., Saglietti, L., and Zecchina, R. (2016a). Unreasonable effectiveness of learning neural networks: From accessible states and robust ensembles to basic algorithmic schemes, *Proceedings of the National Academy of Sciences* **113**, 48, pp. E7655–E7662, doi:10.1073/pnas.1608103113, <http://arxiv.org/abs/1605.06444>, arXiv: 1605.06444.
- Baldassi, C., Ingrosso, A., Lucibello, C., Saglietti, L., and Zecchina, R. (2015). Subdominant Dense Clusters Allow for Simple Learning and High Computational Performance in Neural Networks with Discrete Synapses, *Physical Review Letters* **115**, 12, p. 128101, doi:10.1103/PhysRevLett.115.128101, <http://link.aps.org/doi/10.1103/PhysRevLett.115.128101>.
- Baldassi, C., Ingrosso, A., Lucibello, C., Saglietti, L., and Zecchina, R. (2016b). Local entropy as a measure for sampling solutions in constraint satisfaction problems, *Journal of Statistical Mechanics: Theory and Experiment* **2016**, 2, p. 023301.
- Baldassi, C., Lauditi, C., Malatesta, E. M., Pacelli, R., Perugini, G., and Zecchina, R. (2022). Learning through atypical phase transitions in over-parameterized neural networks, *Physical Review E* **106**, 1, p. 014116.
- Baldassi, C., Pittorino, F., and Zecchina, R. (2020). Shaping the learning landscape in neural networks around wide flat minima, *Proceedings of the National Academy of Sciences* **117**, 1, pp. 161–170.
- Barbier, J., Krzakala, F., Macris, N., Miolane, L., and Zdeborová, L. (2018). Phase Transitions, Optimal Errors and Optimality of Message-Passing in Generalized Linear Models, *Proceedings of the 31st Conference On Learning Theory PMLR* **75**, pp. 728–731, arXiv:1708.03395, <http://arxiv.org/abs/1708.03395>.
- Belkin, M., Hsu, D., Ma, S., and Mandal, S. (2019). Reconciling modern machine-learning practice and the classical bias–variance trade-off, *Proceedings of the National Academy of Sciences* **116**, 32, pp. 15849–15854, doi:10.1073/pnas.1903070116, arXiv:1812.11118, <http://www.pnas.org/lookup/doi/10.1073/pnas.1903070116><http://arxiv.org/abs/1812.11118>.
- Biehl, M. and Schwarze, H. (1995). Learning by on-line gradient descent, *J. Phys. A. Math. Gen.* **28**, 3, pp. 643–656, doi:10.1088/0305-4470/28/3/018.
- Bolthausen, E. (2014). An iterative construction of solutions of the tap equations for the Sherrington-Kirkpatrick model, *Communications in Mathematical Physics* **325**, 1, pp. 333–366.
- Braunstein, A. and Zecchina, R. (2006). Learning by Message Passing in Networks of Discrete Synapses, *Physical Review Letters* **96**, 3, p. 030201, doi:10.1103/PhysRevLett.96.030201, <http://link.aps.org/doi/10.1103/PhysRevLett.96.030201>, arXiv: cond-mat/0511159v2 Genre: Disordered Systems and Neural Networks; Learning; Neurons and Cognition.
- Bray, A. J. and Dean, D. S. (2007). Statistics of critical points of gaussian fields

- on large-dimensional spaces, *Physical review letters* **98**, 15, p. 150201.
- Brown, T. B., Krueger, G., Mann, B., Askell, A., Herbert-voss, A., Winter, C., Ziegler, D. M., Radford, A., and Mccandlish, S. (2020). Language Models are Few-Shot Learners, in *Advances in Neural Information Processing Systems*.
- Canatar, A., Bordelon, B., and Pehlevan, C. (2021). Spectral bias and task-model alignment explain generalization in kernel regression and infinitely wide neural networks, *Nature Communications* **12**, 1, p. 2914, doi: 10.1038/s41467-021-23103-1, arXiv:2006.13198, <http://arxiv.org/abs/2006.13198><http://www.nature.com/articles/s41467-021-23103-1><http://dx.doi.org/10.1038/s41467-021-23103-1>.
- Chaudhari, P., Choromanska, A., Soatto, S., LeCun, Y., Baldassi, C., Borgs, C., Chayes, J., Sagun, L., and Zecchina, R. (2016). Entropy-SGD: Biasing Gradient Descent Into Wide Valleys, *ArXiv e-prints*, pp. 1–14<http://arxiv.org/abs/1611.01838>, arXiv: 1611.01838.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. (2020). A simple framework for contrastive learning of visual representations, in H. D. Iii and A. Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning, Proceedings of Machine Learning Research*, Vol. 119 (PMLR), pp. 1597–1607.
- Crisanti, A., Amit, D. J., and Gutfreund, H. (1986). Saturation Level of the Hopfield Model for Neural Network, *Europhysics Letters (EPL)* **2**, 4, pp. 337–341, doi:10.1209/0295-5075/2/4/012, <https://doi.org/10.1209/0295-5075/2/4/012>, publisher: IOP Publishing.
- D’Ascoli, S., Gabri e, M., Sagun, L., and Biroli, G. (2021). On the interplay between data structure and loss function in classification problems, in *Neural Information Processing Systems 2021*, arXiv:2103.05524, <http://arxiv.org/abs/2103.05524>.
- D’Ascoli, S., Refinetti, M., Biroli, G., and Krzakala, F. (2020). Double trouble in double descent: Bias and variance(s) in the lazy regime, in *International Conference on Machine Learning (ICML)*, arXiv:2003.01054.
- Dauphin, Y. N., Pascanu, R., Gulcehre, C., Cho, K., Ganguli, S., and Bengio, Y. (2014). Identifying and attacking the saddle point problem in high-dimensional non-convex optimization, in *Advances in Neural Information Processing Systems*, pp. 2933–2941.
- Demircigil, M., Heusel, J., L owe, M., Uppang, S., and Vermet, F. (2017). On a Model of Associative Memory with Huge Storage Capacity, *Journal of Statistical Physics* **168**, 2, pp. 288–299, doi:10.1007/s10955-017-1806-y, <http://arxiv.org/abs/1702.01929>, arXiv: 1702.01929.
- Dhifallah, O. and Lu, Y. M. (2020). A Precise Performance Analysis of Learning with Random Features, *arXiv preprint* **2008.11904**, 3, arXiv:2008.11904, <http://arxiv.org/abs/2008.11904>.
- Donoho, D. L. (2006). Compressed sensing, *IEEE Transactions on Information Theory* **52**, 4, pp. 1289–1306, doi:10.1109/TIT.2006.871582, <http://ieeexplore.ieee.org/document/1614066/>.
- Donoho, D. L., Maleki, A., and Montanari, A. (2009). Message-passing algo-

- rithms for compressed sensing, *Proceedings of the National Academy of Sciences* **106**, 45, pp. 18914–18919, doi:10.1073/pnas.0909892106, <http://www.pnas.org/lookup/doi/10.1073/pnas.0909892106>.
- Engel, A. and Van den Broeck, C. (2001). *Statistical Mechanics of Learning* (Cambridge University Press, Cambridge), ISBN 9781139164542, doi:10.1017/CBO9781139164542, <http://ebooks.cambridge.org/ref/id/CBO9781139164542>.
- Fletcher, A. K., Rangan, S., and Schniter, P. (2018). Inference in Deep Networks in High Dimensions, *2018 IEEE International Symposium on Information Theory (ISIT)* **1**, 8, pp. 1884–1888, doi:10.1109/ISIT.2018.8437792, [arXiv:1706.06549](http://arxiv.org/abs/1610.03082), <http://arxiv.org/abs/1610.03082><https://ieeexplore.ieee.org/document/8437792><https://arxiv.org/abs/1706.06549>.
- Foret, P., Kleiner, A., Mobahi, H., and Neyshabur, B. (2021). Sharpness-aware minimization for efficiently improving generalization, in *International Conference on Learning Representations*, <https://openreview.net/forum?id=6Tm1mposlrM>.
- Fort, S., Dziugaite, G. K., Paul, M., Kharaghani, S., Roy, D. M., and Ganguli, S. (2020). Deep learning versus kernel learning: an empirical study of loss landscape geometry and the time evolution of the neural tangent kernel, in *NeurIPS*.
- Franz, S. and Parisi, G. (1995). Recipes for Metastable States in Spin Glasses, *Journal de Physique I* **5**, 11, pp. 1401–1415, doi:10.1051/jp1:1995201, <http://dx.doi.org/10.1051/jp1:1995201>, publisher: EDP Sciences.
- Franz, S. and Parisi, G. (2016). The simplest model of jamming, *J. Phys. A: Math. Theor.* **49**, 14, p. 145001.
- Fyodorov, Y. V. and Williams, I. (2007). Replica symmetry breaking condition exposed by random matrix calculation of landscape complexity, *Journal of Statistical Physics* **129**, 5-6, pp. 1081–1116.
- Gabrielé, M. (2020). Mean-field inference methods for neural networks, *Journal of Mathematical and Theoretical Physics A: Mathematical and Theoretical Physics* **53**, 22, p. 223002, doi:10.1088/1751-8121/ab7f65, [arXiv:1911.00890](http://arxiv.org/abs/1911.00890), <http://arxiv.org/abs/1911.00890><https://iopscience.iop.org/article/10.1088/1751-8121/ab7f65>.
- Gardner, E. (1987). Multiconnected neural network models, *Journal of Physics A: Mathematical and General* **20**, 11, pp. 3453–3464, doi:10.1088/0305-4470/20/11/046, <https://iopscience.iop.org/article/10.1088/0305-4470/20/11/046>.
- Gardner, E. (1988). The space of interactions in neural network models, *Journal of Physics A: Mathematical and General* **21**, 1, pp. 257–270, doi:10.1088/0305-4470/21/1/030, <http://stacks.iop.org/0305-4470/21/i=1/a=030?key=crossref.d0cc5de21f82288ec6af227f5aeac887>, ISBN: 0305-4470.
- Gardner, E. and Derrida, B. (1988). Optimal storage properties of neural network models, *Journal of Physics A: Mathematical and General* **21**, 1, pp. 271–284, doi:10.1088/0305-4470/21/1/031, <http://iopscience.iop.org/0305-4470/21/1/031><http://stacks.iop.org/0305-4470/21/i=1/a=031>

- 031?key=crossref.3a7fdc7f43597f0a05464dd78e14c207.
- Geiger, M., Spigler, S., d'Ascoli, S., Sagun, L., Baity-Jesi, M., Biroli, G., and Wyart, M. (2019). Jamming transition as a paradigm to understand the loss landscape of deep neural networks, *Phys. Rev. E* **100**, p. 012115, doi: 10.1103/PhysRevE.100.012115.
- Gerace, F., Loureiro, B., Krzakala, F., Mézard, M., and Zdeborová, L. (2020). Generalisation error in learning with random features and the hidden manifold model, in *Proceedings of the 37th International Conference on Machine Learning*, PMLR 119:3452-3462, arXiv:2002.09339, <http://arxiv.org/abs/2002.09339>.
- Gerbelot, C. and Berthier, R. (2021). Graph-based Approximate Message Passing Iterations, *arXiv preprint* **2109.11905**, arXiv:2109.11905, <http://arxiv.org/abs/2109.11905>.
- Ghorbani, B., Mei, S., Misiakiewicz, T., and Montanari, A. (2020). When Do Neural Networks Outperform Kernel Methods? in *Neural Information Processing Systems*, arXiv:2006.13409, <http://arxiv.org/abs/2006.13409>.
- Goldt, S., Advani, M. S., Saxe, A. M., Krzakala, F., and Zdeborová, L. (2019a). Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student setup, in *Neural Information Processing Systems*, Vol. 1906.08632, arXiv:1906.08632, <http://arxiv.org/abs/1906.08632>.
- Goldt, S., Loureiro, B., Reeves, G., Mézard, M., Krzakala, F., Zdeborová, L., Loureiro, B., Reeves, G., Krzakala, F., Mézard, M., and Zdeborová, L. (2021). The Gaussian equivalence of generative models for learning with two-layer neural networks, *Mathematical and Scientific Machine Learning* **145**, pp. 1–37, arXiv:2006.14709, <http://arxiv.org/abs/2006.14709>.
- Goldt, S., Mézard, M., Krzakala, F., and Zdeborová, L. (2019b). Modelling the influence of data structure on learning in neural networks, *arXiv preprint* **1909.11500**, 4, pp. 1–32, doi:10.1103/PhysRevX.10.041044, arXiv:1909.11500, <https://link.aps.org/doi/10.1103/PhysRevX.10.041044><http://arxiv.org/abs/1909.11500>.
- Goldt, S., Mézard, M., Krzakala, F., and Zdeborová, L. (2020). Modeling the Influence of Data Structure on Learning in Neural Networks: The Hidden Manifold Model, *Physical Review X* **10**, 4, p. 041044, doi:10.1103/PhysRevX.10.041044, arXiv:1909.11500, <https://link.aps.org/doi/10.1103/PhysRevX.10.041044><http://arxiv.org/abs/1909.11500>.
- Gordon, M. A., Duh, K., and Kaplan, J. (2021). Data and parameter scaling laws for neural machine translation, in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing* (Association for Computational Linguistics, Online and Punta Cana, Dominican Republic), pp. 5915–5922.
- Grill, Strub, Altché, Tallec, and others (2020). Bootstrap your own latent—a new approach to self-supervised learning, *Adv. Neural Inf. Process. Syst.* .
- Hebb, D. O. (1949). *The organization of behavior* (Wiley, New York).
- Henighan, T., Kaplan, J., Katz, M., Chen, M., Hesse, C., Jackson, J., Jun, H., Brown, T. B., Dhariwal, P., Gray, S., et al. (2020). Scaling laws for autoregressive generative modeling, *arXiv preprint arXiv:2010.14701* .

- Hernandez, D., Kaplan, J., Henighan, T., and McCandlish, S. (2021). Scaling laws for transfer, *arXiv preprint arXiv:2102.01293* .
- Hestness, J., Narang, S., Ardalani, N., Diamos, G., Jun, H., Kianinejad, H., Patwary, M., Ali, M., Yang, Y., and Zhou, Y. (2017). Deep learning scaling is predictable, empirically, *arXiv preprint arXiv:1712.00409* .
- Ho, J., Jain, A., and Abbeel, P. (2020). Denoising diffusion probabilistic models, *Advances in Neural Information Processing Systems* **33**, pp. 6840–6851.
- Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., Casas, D. d. L., Hendricks, L. A., Welbl, J., Clark, A., *et al.* (2022). Training compute-optimal large language models, *arXiv preprint arXiv:2203.15556* .
- Hopfield, J. J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences of the United States of America* **79**, 8, pp. 2554–2558, doi:10.1073/pnas.79.8.2554.
- Hornik, K. (1991). Approximation Capabilities of Multilayer Neural Network, *Neural Networks* **4**, 1989, pp. 251–257, doi:10.1016/0893-6080(91)90009-T.
- Hu, H. and Lu, Y. M. (2020). Universality Laws for High-Dimensional Learning with Random Features, *arXiv preprint* **2009.07669**, 1, [arXiv:arXiv:2009.07669v1](https://arxiv.org/abs/2009.07669).
- Huang, H., Wong, K. Y. M., and Kabashima, Y. (2013). Entropy landscape of solutions in the binary perceptron problem, *Journal of Physics A: Mathematical and Theoretical* **46**, 37, p. 375002, doi:10.1088/1751-8113/46/37/375002, [arXiv:arXiv:1304.2850v2](https://arxiv.org/abs/1304.2850v2), <https://iopscience.iop.org/article/10.1088/1751-8113/46/37/375002>.
- Iba, Y. (1999). The Nishimori line and Bayesian statistics, *Journal of Physics A: Mathematical and General* **32**, 21, pp. 3875–3888, doi:10.1088/0305-4470/32/21/302, <https://doi.org/10.1088/0305-4470/32/21/302>, publisher: IOP Publishing.
- Jacot, A., Gabriel, F., and Hongler, C. (2018). Neural Tangent Kernel: Convergence and Generalization in Neural Networks, in *Advances in neural information processing systems*, 5, [arXiv:1806.07572](https://arxiv.org/abs/1806.07572), <http://arxiv.org/abs/1806.07572>.
- Jiang, Y., Neyshabur, B., Mobahi, H., Krishnan, D., and Bengio, S. (2020). Fantastic generalization measures and where to find them, in *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020* (OpenReview.net), <https://openreview.net/forum?id=SJgIPJBFvH>.
- Kabashima, Y. and Saad, D. (1998). Belief propagation vs. TAP for decoding corrupted messages, *Europhysics Letters* **44**, 5, p. 668, doi:10.1209/epl/i1998-00524-7, <https://iopscience.iop.org/article/10.1209/epl/i1998-00524-7/meta>, publisher: IOP Publishing.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. (2020). Scaling laws for neural language models, *arXiv preprint arXiv:2001.08361* .
- Krotov, D. and Hopfield, J. J. (2016). Dense Associative Memory for Pattern

- Recognition, *Advances in Neural Information Processing Systems*, pp. 1180–1188 <http://arxiv.org/abs/1606.01164>.
- Krzakala, F., Montanari, A., Ricci-Tersenghi, F., Semerjian, G., and Zdeborová, L. (2007). Gibbs states and the set of solutions of random constraint satisfaction problems. *Proceedings of the National Academy of Sciences of the United States of America* **104**, 25, pp. 10318–23, doi:10.1073/pnas.0703685104, <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1965511&tool=pmcentrez&rendertype=abstract>, iSBN: 0703685104.
- Lampinen, A. K. and Ganguli, S. (2018). An analytic theory of generalization dynamics and transfer learning in deep linear networks, in *International Conference on Learning Representations (ICLR)*.
- LeCun, Y., Bengio, Y., and Hinton, G. E. (2015). Deep learning, *Nature* **521**, 7553, pp. 436–444, doi:10.1038/nature14539, <http://www.nature.com/doifinder/10.1038/nature14539>.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition, *Proceedings of the IEEE* **86**, 11, pp. 2278–2323, doi:10.1109/5.726791.
- Loureiro, B., Gerbelot, C., Cui, H., Goldt, S., Krzakala, F., Mézard, M., and Zdeborová, L. (2021a). Learning curves of generic features maps for realistic datasets with a teacher-student model, *Advances in Neural Information Processing Systems* **22**, NeurIPS, pp. 18137–18151, [arXiv:2102.08127](https://arxiv.org/abs/2102.08127).
- Loureiro, B., Gerbelot, C., Refinetti, M., Sicuro, G., and Krzakala, F. (2022a). Fluctuations, Bias, Variance & Ensemble of Learners: Exact Asymptotics for Convex Losses in High-Dimension, in *Proceedings of the 39th International Conference on Machine Learning (PMLR)*, pp. 14283–14314, <https://proceedings.mlr.press/v162/loureiro22a.html>, iSSN: 2640-3498.
- Loureiro, B., Gerbelot, C., Refinetti, M., Sicuro, G., and Krzakala, F. (2022b). Fluctuations, bias, variance & ensemble of learners: Exact asymptotics for convex losses in high-dimension, in K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning, Proceedings of Machine Learning Research*, Vol. 162 (PMLR), pp. 14283–14314, <https://proceedings.mlr.press/v162/loureiro22a.html>.
- Loureiro, B., Sicuro, G., Gerbelot, C., Pacco, A., Krzakala, F., and Zdeborová, L. (2021b). Learning Gaussian Mixtures with Generalised Linear Models: Precise Asymptotics in High-dimensions, in *Advances in Neural Information Processing Systems* **35**, [arXiv:2106.03791](https://arxiv.org/abs/2106.03791), <http://arxiv.org/abs/2106.03791>.
- Lucibello, C., Pittorino, F., Perugini, G., and Zecchina, R. (2022). Deep learning via message passing algorithms based on belief propagation, *Machine Learning: Science and Technology* **3**, 3, p. 035005, doi:10.1088/2632-2153/ac7d3b, <https://doi.org/10.1088/2632-2153/ac7d3b>, publisher: IOP Publishing.
- Maillard, A., Ben Arous, G., and Biroli, G. (2020). Landscape complexity for the empirical risk of generalized linear models, in J. Lu and R. Ward (eds.),

- Proceedings of The First Mathematical and Scientific Machine Learning Conference, Proceedings of Machine Learning Research*, Vol. 107 (PMLR), pp. 287–327.
- Maillard, A., Krzakala, F., Mézard, M., and Zdeborová, L. (2022). Perturbative construction of mean-field equations in extensive-rank matrix factorization and denoising, *Journal of Statistical Mechanics: Theory and Experiment* **2022**, 8, p. 083301, doi:10.1088/1742-5468/ac7e4c, <https://dx.doi.org/10.1088/1742-5468/ac7e4c>.
- Manoel, A., Krzakala, F., Mézard, M., and Zdeborová, L. (2017). Multi-layer generalized linear estimation, in *2017 IEEE International Symposium on Information Theory (ISIT)* (IEEE), ISBN 978-1-5090-4096-4, pp. 2098–2102, doi:10.1109/ISIT.2017.8006899, arXiv:1701.06981, <http://arxiv.org/abs/1701.06981><http://dx.doi.org/10.1109/ISIT.2017.8006899><http://ieeexplore.ieee.org/document/8006899/>.
- Mato, G. and Parga, N. (1992). Generalization properties of multilayered neural networks, *Journal of Physics A: Mathematical and General* **25**, 19, pp. 5047–5054, doi:10.1088/0305-4470/25/19/017.
- Mei, S. and Montanari, A. (2021). The Generalization Error of Random Features Regression: Precise Asymptotics and the Double Descent Curve, *Communications on Pure and Applied Mathematics* **75**, 4, pp. 667–766, doi:10.1002/cpa.22008, arXiv:arXiv:1908.05355v5, <https://onlinelibrary.wiley.com/doi/10.1002/cpa.22008>.
- Mei, S., Montanari, A., and Nguyen, P.-M. (2018). A mean field view of the landscape of two-layer neural networks, *Proceedings of the National Academy of Sciences* **115**, 33, pp. E7665–E7671, doi:10.1073/pnas.1806579115, <https://www.pnas.org/doi/abs/10.1073/pnas.1806579115>, publisher: Proceedings of the National Academy of Sciences.
- Mel, G. and Ganguli, S. (2021). A theory of high dimensional regression with arbitrary correlations between input features and target functions: sample complexity, multiple descent curves and a hierarchy of phase transitions, in M. Meila and T. Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, Proceedings of Machine Learning Research*, Vol. 139 (PMLR), pp. 7578–7587.
- Mezard, M. (1989). The space of interactions in neural networks: Gardner’s computation with the cavity method, *Journal of Physics A: Mathematical and General* **22**, 12, pp. 2181–2190, doi:10.1088/0305-4470/22/12/018, <http://stacks.iop.org/0305-4470/22/i=12/a=018?key=crossref.f0820bd8e5653d6fdacb0f2beb1f5def>.
- Mézard, M. and Montanari, A. (2009). *Information, Physics, and Computation* (Oxford University Press).
- Mignacco, F., Krzakala, F., Urbani, P., and Zdeborová, L. (2020). Dynamical mean-field theory for stochastic gradient descent in Gaussian mixture classification, in *Advances in Neural Information Processing Systems*, Vol. 33 (Curran Associates, Inc.), pp. 9540–9550, <https://proceedings.neurips.cc/paper/2020/hash/6c81c83c4bd0b58850495f603ab45a93-Abstract.html>.

- Monasson, R. and Zecchina, R. (1995). Weight Space Structure and Internal Representations: a Direct Approach to Learning and Generalization in Multilayer Neural Networks, *Physical Review Letters* **75**, 12, pp. 2432–2435, doi:10.1103/PhysRevLett.75.2432, arXiv:9501082v1 [arXiv:cond-mat].
- Mézard, M. (2017). Mean-field message-passing equations in the Hopfield model and its generalizations, *Physical Review E* **95**, 2, p. 022117, doi:10.1103/PhysRevE.95.022117, <https://link.aps.org/doi/10.1103/PhysRevE.95.022117>, publisher: American Physical Society.
- Nishimori, H. (1980). Exact results and critical properties of the Ising model with competing interactions, *Journal of Physics C: Solid State Physics* **13**, 21, pp. 4071–4076, doi:10.1088/0022-3719/13/21/012, <https://doi.org/10.1088/0022-3719/13/21/012>, publisher: IOP Publishing.
- Paul, M., Ganguli, S., and Dziugaite, G. K. (2021). Deep learning on a data diet: Finding important examples early in training, *Adv. Neural Inf. Process. Syst.* **34**.
- Pennington, J., Schoenholz, S. S., and Ganguli, S. (2017). Resurrecting the sigmoid in deep learning through dynamical isometry: Theory and practice, in *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS'17* (Curran Associates Inc., Red Hook, NY, USA), ISBN 9781510860964, p. 4788–4798.
- Pennington, J., Schoenholz, S. S., and Ganguli, S. (2018). The emergence of spectral universality in deep networks, in *Artificial Intelligence and Statistics (AISTATS)*.
- Pennington, J. and Worah, P. (2017). Nonlinear random matrix theory for deep learning, in *Neural Information Processing Systems*, 31, pp. 2634–2643.
- Pittorino, F., Lucibello, C., Feinauer, C., Perugini, G., Baldassi, C., Demyanenko, E., and Zecchina, R. (2021). Entropic gradient descent algorithms and wide flat minima, *Journal of Statistical Mechanics: Theory and Experiment* **2021**, 12, p. 124015, doi:10.1088/1742-5468/ac3ae8, <https://doi.org/10.1088/1742-5468/ac3ae8>, publisher: IOP Publishing.
- Poole, B., Lahiri, S., Raghu, M., Sohl-Dickstein, J., and Ganguli, S. (2016). Exponential expressivity in deep neural networks through transient chaos, *Adv. Neural Inf. Process. Syst.* .
- Rahimi, A. and Recht, B. (2017). Random Features for Large-Scale Kernel Machines, in *Neural Information Processing Systems*.
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen, M. (2022). Hierarchical text-conditional image generation with CLIP latents, *CoRR abs/2204.06125*, doi:10.48550/arXiv.2204.06125, 2204.06125, <https://doi.org/10.48550/arXiv.2204.06125>.
- Ramsauer, H., Schäfl, B., Lehner, J., Seidl, P., Widrich, M., Gruber, L., Holzleitner, M., Adler, T., Kreil, D., Kopp, M. K., Klambauer, G., Brandstetter, J., and Hochreiter, S. (2021). Hopfield networks is all you need, in *International Conference on Learning Representations*, <https://openreview.net/forum?id=tL89RnzIiCd>.
- Rangan, S. (2011). Generalized approximate message passing for estimation with random linear mixing, in *2011 IEEE International Symposium on Infor-*

- mation Theory Proceedings* (IEEE), ISBN 978-1-4577-0596-0, pp. 2168–2172, doi:10.1109/ISIT.2011.6033942, arXiv:1010.5141, <http://arxiv.org/abs/1010.5141><http://ieeexplore.ieee.org/document/6033942/>.
- Refinetti, M., Goldt, S., Krzakala, F., and Zdeborova, L. (2021). Classifying high-dimensional gaussian mixtures: Where kernel methods fail and neural networks succeed, in M. Meila and T. Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, Proceedings of Machine Learning Research*, Vol. 139 (PMLR), pp. 8936–8947, <https://proceedings.mlr.press/v139/refinetti21b.html>.
- Rogers, T. T. and McClelland, J. L. (2004). *Semantic cognition: A parallel distributed processing approach* (The MIT Press).
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. (2022). High-resolution image synthesis with latent diffusion models, in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10684–10695.
- Saad, D. (1999). *On-line learning in neural networks* (Cambridge University Press).
- Saad, D. and Solla, S. A. (1995a). Exact solution for on-line learning in multilayer neural networks, *Physical Review Letters* **74**, 21, pp. 4337–4340, doi:10.1103/PhysRevLett.74.4337.
- Saad, D. and Solla, S. A. (1995b). On-line learning in soft committee machines, *Physical Review E* **52**, 4, pp. 4225–4243.
- Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E. L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., *et al.* (2022). Photorealistic text-to-image diffusion models with deep language understanding, *Advances in Neural Information Processing Systems* **35**, pp. 36479–36494.
- Saxe, A. M., McClelland, J. L., and Ganguli, S. (2013). Exact solutions to the nonlinear dynamics of learning in deep linear neural networks, *Advances in Neural Information Processing Systems*, pp. 1–9 arXiv:1312.6120, <http://arxiv.org/abs/1312.6120>.
- Saxe, A. M., McClelland, J. L., and Ganguli, S. (2019). A mathematical theory of semantic development in deep neural networks, *Proc. Natl. Acad. Sci. U. S. A.*
- Schoenholz, S. S., Brain, G., Gilmer, J., Brain, G., Ganguli, S., Sohl-Dickstein, J., and Brain, G. (2017). Deep Information Propagation, in *International Conference on Learning Representations 2017*, pp. 1–18, arXiv:1611.01232, <http://arxiv.org/abs/1611.01232>.
- Schwarze, H. (1993). Learning a Rule in a Multilayer Neural-Network, *Journal of Physics a-Mathematical and General* **26**, 21, pp. 5781–5794, doi:10.1088/0305-4470/26/21/017.
- Schwarze, H. and Hertz, J. (1993). Generalization in Fully Connected Committee Machines, *Europhysics Letters (EPL)* **21**, 7, pp. 785–790, doi:10.1209/0295-5075/21/7/012, <http://stacks.iop.org/0295-5075/21/i=7/a=012?key=crossref.6ea69b4612ff4e290227b05c508f1b8b>.
- Shalev-Shwartz, S. and Ben-David, S. (2014). *Understanding Machine Learn-*

- ing: From Theory to Algorithms* (Cambridge University Press, Cambridge), ISBN 9781107298019, doi:10.1017/CBO9781107298019, <http://ebooks.cambridge.org/ref/id/CBO9781107298019>.
- Sohl-Dickstein, J., Weiss, E. A., Maheswaranathan, N., and Ganguli, S. (2015). Deep unsupervised learning using nonequilibrium thermodynamics, *International Conference on Machine Learning (ICML)* .
- Sorscher, B., Geirhos, R., Shekhar, S., Ganguli, S., and Morcos, A. (2022). Beyond neural scaling laws: beating power law scaling via data pruning, *Advances in Neural Information Processing Systems* **35**, pp. 19523–19536.
- Spigler, S., Geiger, M., D’Ascoli, S., Sagun, L., Biroli, G., and Wyart, M. (2019). A jamming transition from under- to over-parametrization affects generalization in deep learning, *Journal of Physics A: Mathematical and Theoretical* **52**, 47, p. 474001, doi:10.1088/1751-8121/ab4c8b, arXiv:1810.09665, <http://arxiv.org/abs/1810.09665><https://iopscience.iop.org/article/10.1088/1751-8121/ab4c8b>.
- Tian, Y., Chen, X., and Ganguli, S. (2021). Understanding self-supervised learning dynamics without contrastive pairs, in M. Meila and T. Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, Proceedings of Machine Learning Research*, Vol. 139 (PMLR), pp. 10268–10278.
- Tian, Y., Yu, L., Chen, X., and Ganguli, S. (2020). Understanding self-supervised learning with dual deep networks, arXiv:2010.00578 [cs.LG].
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention is all you need, in *Advances in neural information processing systems*, pp. 5998–6008.
- Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., Chi, E. H., Hashimoto, T., Vinyals, O., Liang, P., Dean, J., and Fedus, W. (2022). Emergent abilities of large language models, *Transactions on Machine Learning Research* <https://openreview.net/forum?id=yzkSU5zdWd>, survey Certification.
- Xiao, L., Bahri, Y., Sohl-Dickstein, J., Schoenholz, S., and Pennington, J. (2018). Dynamical isometry and a mean field theory of cnns: How to train 10,000-layer vanilla convolutional neural networks, in *International Conference on Machine Learning*, pp. 5393–5402.
- Yedidia, J. S., Freeman, W. T., and Weiss, Y. (2002). Understanding Belief Propagation and its Generalizations, *Intelligence* **8**, pp. 236–239, doi:10.1348/026151002166325, <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.164.7114&rep=rep1&type=pdf>.
- Zdeborová, L. and Krzakala, F. (2016). Statistical physics of inference: Thresholds and algorithms, *Advances in Physics* **65**, 5, pp. 453–552, doi:10.1080/00018732.2016.1211393, arXiv:1511.02476, <http://arxiv.org/abs/1511.02476><https://www.tandfonline.com/doi/full/10.1080/00018732.2016.1211393>.
- Zhai, X., Kolesnikov, A., Houlsby, N., and Beyer, L. (2021). Scaling vision transformers, *arXiv preprint arXiv:2106.04560* .