Vulnerability of Deep Learning

Richard Kenway*

School of Physics and Astronomy, University of Edinburgh, James Clerk Maxwell Building, Peter Guthrie Tait Road, Edinburgh EH9 3FD, UK (Dated: March 19, 2018)

The Renormalisation Group (RG) provides a framework in which it is possible to assess whether a deep-learning network is sensitive to small changes in the input data and hence prone to error, or susceptible to adversarial attack. Distinct classification outputs are associated with different RG fixed points and sensitivity to small changes in the input data is due to the presence of relevant operators at a fixed point. A numerical scheme, based on Monte Carlo RG ideas, is proposed for identifying the existence of relevant operators and the corresponding directions of greatest sensitivity in the input data. Thus, a trained deep-learning network may be tested for its robustness and, if it is vulnerable to attack, dangerous perturbations of the input data identified.

I. INTRODUCTION

A. The Threat to Deep Learning

Despite many successful applications, deep learning [1–3] is poorly understood. The absence of an underlying theory means that we cannot be certain under what circumstances a given trained network will operate correctly and this uncertainty calls into question the use of such networks in safety-critical applications. There is also the possibility that slightly perturbed input data can be constructed that would fool the network, but not a human, making the network susceptible to adversarial attack, or simply cause it to misclassify and thereby conceal information in the data [4].

It is important to distinguish between ambiguous data, which any system, machine or human, would find difficult to classify and data that has been imperceptibly altered to cause the deep network to misclassify it. This paper introduces a theoretical framework, based on the Renormalisation Group (RG), which specifies the conditions under which the latter may occur and proposes a calculational method for testing a trained network for this vulnerability. Further, it determines the directions in the space of input data along which small perturbations are amplified by such a vulnerable deep network, causing it to generate a qualitatively incorrect output.

The theoretical framework captures the salient features of deep learning: the multi-layer structure, trained one layer at a time. It is demonstrated for the specific case in which each layer is a Restricted Boltzmann Machine (RBM) [3] and the network is trained to classify a data set $\{\mathbf{x}\}$ in terms of outputs $\{\mathbf{y}\}$, whose exact conditional probability distribution is $p(\mathbf{y}|\mathbf{x})$. However, the approach may be extended straightforwardly to other deep networks and to data generation as well as classification.

B. Analogy with the Renormalisation Group

The formal analogy between deep learning and the RG has been noted before, although with some controversy about whether the context should be unsupervised, or supervised learning [5–7]. There are more fundamental issues.

The first is that the RG applied to critical phenomena describes how features of the microscopic degrees of freedom of a system (analogous to the input data $\{x\}$) depend on length scale, successively eliminating short-distance fluctuations and generating a sequence of effective Hamiltonians, which describe the surviving features on larger and larger length scales. A critical point is associated with scale-invariant fluctuations and convergence of the sequence of effective Hamiltonians to a fixed-point Hamiltonian which describes them. Although many applications of deep learning do involve data which exhibit structure on different length scales, e.g., data drawn from the physical world, and this appears to be captured by deep learning so that it is "cheap" [6], this need not be

The second is that the RG transformation effected by each layer of the network may be different, reflecting the emergence of different types of features in the data, whereas in critical phenomena each application of the RG transformation changes the length scale by a fixed amount.

We will assume that training of the deep network converges to a conditional probability distribution that is a good approximation to $p(\mathbf{y}|\mathbf{x})$, provided there is sufficient number of layers. Thus, the sequence of RG transformations corresponding to applying the RBM in each layer to the output from the hidden nodes in the previous layer, starting with the input data, converges.

Presumably, some more generalised form of scaling takes place in deep learning and we will assume that this has the effect of successively scaling to smaller values all but a finite number of eigenvalues of the Fisher Information Matrix (FIM) of the trained deep network, so that its spectrum is sloppy [8]. We will say more about this in Section IV.

In Section II we introduce the RG formalism in the

^{*} Email address: r.d.kenway@ed.ac.uk

context of a deep RBM network used for classification and, in particular, define the sequence of effective Hamiltonians whose fixed points and universality classes define distinct outputs. In Section III we develop the computational scheme, based on Monte Carlo RG [9, 10], for estimating the stability of these fixed points with respect to the depth of the network. Section IV relates the instability of a particular fixed point to direction(s) in the input data space $\{x\}$, along which small changes in the data cause major changes in the output probability distribution and, hence, are sources of vulnerability.

II. RENORMALISATION GROUP FOR DEEP LEARNING

Consider a layered network of N RBMs with input nodes \mathbf{h}_0 , N-1 layers of hidden nodes \mathbf{h}_k , $k=1,\ldots,N-1$, and output nodes \mathbf{h}_N . For example, the hidden nodes may be binary vectors of the same dimension as the input data $\{\mathbf{x}\}$ and the outputs $\{\mathbf{y}\}$, and the joint probability distribution for the k^{th} layer is of the form [3]

$$t_k(\mathbf{h}_k, \mathbf{h}_{k-1}) = \frac{e^{\mathbf{h}_k^{\mathrm{T}} \mathbf{W}_k \mathbf{h}_{k-1} + \mathbf{a}_k^{\mathrm{T}} \mathbf{h}_k + \mathbf{b}_k^{\mathrm{T}} \mathbf{h}_{k-1}}}{z_k},$$

$$z_k = \mathrm{Tr}_{\mathbf{h}_k, \mathbf{h}_{k-1}} e^{\mathbf{h}_k^{\mathrm{T}} \mathbf{W}_k \mathbf{h}_{k-1} + \mathbf{a}_k^{\mathrm{T}} \mathbf{h}_k + \mathbf{b}_k^{\mathrm{T}} \mathbf{h}_{k-1}}, (1)$$

where \mathbf{W}_k , \mathbf{a}_k and \mathbf{b}_k are weights determined by the layerwise training. We will not require the precise form of t_k .

The trained deep RBM network generates the probability distribution for the output \mathbf{y} given input data \mathbf{x} , $q_N(\mathbf{y}|\mathbf{x})$, iteratively for k = 1, ..., N through

$$q_k(\mathbf{h}_k|\mathbf{x}) = \text{Tr}_{\mathbf{h}_{k-1}} t_k(\mathbf{h}_k|\mathbf{h}_{k-1}) q_{k-1}(\mathbf{h}_{k-1}|\mathbf{x}),$$

$$q_0(\mathbf{h}_0|\mathbf{x}) = \delta_{\mathbf{h}_0,\mathbf{x}},$$

$$q_N(\mathbf{y}|\mathbf{x}) \approx p(\mathbf{y}|\mathbf{x}),$$
(2)

where, according to Bayes Theorem,

$$t_k(\mathbf{h}_k|\mathbf{h}_{k-1}) = \frac{t_k(\mathbf{h}_k, \mathbf{h}_{k-1})}{\operatorname{Tr}_{\mathbf{h}_k} t_k(\mathbf{h}_k, \mathbf{h}_{k-1})}.$$
 (3)

We define the effective Hamiltonian for the k^{th} layer by

$$q_k(\mathbf{h}|\mathbf{x}) = \frac{e^{-H_{\mathbf{x}}^{(k)}(\mathbf{h})}}{Z_{\mathbf{x}}^{(k)}},$$
$$Z_{\mathbf{x}}^{(k)} = \text{Tr}_{\mathbf{h}}e^{-H_{\mathbf{x}}^{(k)}(\mathbf{h})}.$$
 (4)

Then, the sequence of Hamiltonians is generated by the RBMs in each layer via

$$\frac{e^{-H_{\mathbf{x}}^{(k+1)}(\mathbf{h})}}{Z_{\mathbf{x}}^{(k+1)}} = \operatorname{Tr}_{\mathbf{h}'} t_{k+1}(\mathbf{h}|\mathbf{h}') \frac{e^{-H_{\mathbf{x}}^{(k)}(\mathbf{h}')}}{Z_{\mathbf{x}}^{(k)}}, \tag{5}$$

such that each Hamiltonian depends only on the previous one in the sequence. This has the same form as an RG transformation with kernel t_{k+1} [5].

At this point, it is worth noting that the input data \mathbf{x} determines the couplings in the space of effective Hamiltonians via the first layer,

$$\frac{e^{-H_{\mathbf{x}}^{(1)}(\mathbf{h})}}{Z_{\mathbf{x}}^{(1)}} = t_1(\mathbf{h}|\mathbf{x}),\tag{6}$$

and the couplings in each subsequent layer are functions solely of the couplings in the previous layer.

We parametrise the space of effective Hamiltonians in terms of a complete set of operators $O_{\alpha}(\mathbf{h})$, which, in the case of binary variables, consist of all possible products of the components of \mathbf{h} , with couplings g_{α} . Thus, the effective Hamiltonian for layer k is

$$H_{\mathbf{x}}^{(k)}(\mathbf{h}) = \sum_{\alpha} g_{\alpha}^{(k)} O_{\alpha}(\mathbf{h}). \tag{7}$$

Then the effect of the RG transformation in eq. (5) is to define $\{g^{(k+1)}\}$ solely in terms of $\{g^{(k)}\}$ and, thereby, to generate a flow in the coupling-constant space of the effective Hamiltonians.

The key assumption we make is that the training of the deep RBM network converges, so that, for N large enough, the sequence of effective Hamiltonians converges to a fixed point:

$$H_{\mathbf{x}}^*(\mathbf{h}) = \sum_{\alpha} g_{\alpha}^* O_{\alpha}(\mathbf{h}), \tag{8}$$

such that

$$\frac{e^{-H_{\mathbf{x}}^*(\mathbf{y})}}{Z_{\mathbf{x}}^*} = q_N(\mathbf{y}|\mathbf{x}) \approx p(\mathbf{y}|\mathbf{x}). \tag{9}$$

The analogue of the RG universality class of the fixed point $H_{\mathbf{x}}^*$ is the subset of the input data, $C_{\mathbf{x}}$, which are associated with the same output distribution, $p(\mathbf{y}|\mathbf{x})$. Thus, assuming the training correctly classifies the data,

$$\mathbf{x}, \mathbf{x}' \in C_{\mathbf{x}} \Rightarrow H_{\mathbf{x}}^* = H_{\mathbf{x}'}^* \tag{10}$$

and

$$p(\mathbf{y}|\mathbf{x}) \neq p(\mathbf{y}|\mathbf{x}') \Rightarrow H_{\mathbf{y}}^* \neq H_{\mathbf{y}'}^*,$$
 (11)

so that there must be more than one fixed point. In any non-trivial classification problem, there is a distinct fixed point of the correctly trained deep network for every distinct output.

Next, we consider the flow in coupling-constant space close to a fixed point and, in particular, the stability properties of the fixed point as the number of layers is increased. We assume that the deeper the network the more accurately it is able to learn the classification problem. Close to the fixed point $\{g^*\}$, the relationship between the couplings in adjacent layers is, to leading order,

$$g_{\alpha}^{(k+1)}(\{g^{(k)}\}) = g_{\alpha}^{(k+1)}(\{g^*\}) + \sum_{\beta} (g_{\beta}^{(k)} - g_{\beta}^*) \left. \frac{\partial g_{\alpha}^{(k+1)}}{\partial g_{\beta}^{(k)}} \right|_{\{g^*\}}. (12)$$

Since $g_{\alpha}^{(k+1)}(\{g^*\}) = g_{\alpha}^*$, this becomes

$$g_{\alpha}^{(k+1)} - g_{\alpha}^* = \sum_{\beta} T_{\alpha\beta}^* (g_{\beta}^{(k)} - g_{\beta}^*),$$
 (13)

where

$$T_{\alpha\beta}^* = \left. \frac{\partial g_{\alpha}^{(k+1)}}{\partial g_{\beta}^{(k)}} \right|_{\{q^*\}} \tag{14}$$

is the stability matrix of the fixed point. The stability properties of a fixed point determine whether it becomes sensitive to small changes in the input data as the number of layers is increased.

Define the eigenvalues and eigenvectors of the stability matrix by

$$\sum_{\beta} T_{\alpha\beta}^* \phi_{\beta}^{(\mu)} = \Lambda_{\mu} \phi_{\alpha}^{(\mu)} \tag{15}$$

and express the couplings, g_{α} , and the corresponding operators, $O_{\alpha}(\mathbf{h})$, in the basis of eigenvectors of the stability matrix as

$$g_{\alpha} = \sum_{\mu} \phi_{\alpha}^{(\mu)} \tilde{g}_{\mu}, \tag{16}$$

$$\tilde{O}_{\mu}(\mathbf{h}) = \sum_{\alpha} \phi_{\alpha}^{(\mu)} O_{\alpha}(\mathbf{h}). \tag{17}$$

Then the effective Hamiltonian in this basis is

$$H_{\mathbf{x}}^{(k)} = \sum_{\mu} \tilde{g}_{\mu}^{(k)} \tilde{O}_{\mu}(\mathbf{h}) \tag{18}$$

and, close to the fixed point, from eqs (13) and (15),

$$H_{\mathbf{x}}^{(k+1)} - H_{\mathbf{x}}^* = \sum_{\mu} (\tilde{g}_{\mu}^{(k+1)} - \tilde{g}_{\mu}^*) \tilde{O}_{\mu}(\mathbf{h})$$
$$= \sum_{\mu} \Lambda_{\mu} (\tilde{g}_{\mu}^{(k)} - \tilde{g}_{\mu}^*) \tilde{O}_{\mu}(\mathbf{h}). \tag{19}$$

In the language of the RG, if $\Lambda_{\mu} > 1$, $\tilde{O}_{\mu}(\mathbf{h})$ is relevant and, if $\Lambda_{\mu} < 1$, it is irrelevant. In the context of deep learning, the existence of a relevant operator at a fixed point creates a sensitivity to input data which generate a non-zero coupling to this operator. Such data will lead to a probability distribution for the output which increasingly deviates from the fixed-point distribution as the network is made deeper (more layers are added). This is the source of vulnerability of deep networks to tiny changes in the input data which excite a relevant operator, leading to the sequence of effective Hamiltonians converging to the wrong fixed point, i.e., misclassifying the data.

As it stands, this provides an explanation of the potential vulnerability of deep networks, but doesn't help to diagnose whether a particular network is susceptible. This is because, in practice, we only have access to the RBM weights, e.g., eq. (1), which enable us to sample

the hidden nodes and output nodes for a given input, but don't enable us to construct the sequence of effective Hamiltonians, or the fixed-point Hamiltonian. In the next section, we will explain how Monte Carlo RG methods may be used to estimate the stability matrix $T_{\alpha\beta}^*$, defined in eq. (14), and how its eigenvalues behave in deep networks close to the fixed point. This will be sufficient to determine whether a particular trained deep network has any fixed points with relevant operators.

It is worth noting that the existence of multiple fixed points (associated with a non-trivial classification problem) does not imply the existence of relevant operators at any of the fixed points. In the robust situation, where no fixed point has a relevant direction, data close to the boundary between two universality classes are simply ambiguous and do not represent a source of adversarial attack (because a human would be equally likely to misclassify them).

III. STABILITY OF DEEP LEARNING

Having trained a deep network of RBMs on a classification problem whose exact solution is $p(\mathbf{y}|\mathbf{x})$, we have available a sequence of RBMs given, for example, by eq. (1), which enables us to use the conditional probability distributions in eq. (2) to sample the hidden nodes in each layer and the output layer. Using this capability, we apply methods developed for Monte Carlo RG [9, 10] to estimate the stability matrix for each fixed point in the space of effective Hamiltonians.

Denote by $\langle O_{\gamma} \rangle^{(k)}$ the expectation value of the operator O_{γ} with respect to the effective Hamiltonian in layer k, eq. (4), i.e.,

$$\langle O_{\gamma} \rangle^{(k)} = \operatorname{Tr}_{\mathbf{h}} O_{\gamma}(\mathbf{h}) q_{k}(\mathbf{h} | \mathbf{x})$$

$$= \frac{\operatorname{Tr}_{\mathbf{h}} O_{\gamma}(\mathbf{h}) e^{-H_{\mathbf{x}}^{(k)}(\mathbf{h})}}{Z_{\mathbf{x}}^{(k)}}.$$
(20)

We start from the chain rule applied to $\langle O_{\gamma} \rangle^{(k+1)}$ close to a fixed point (i.e., $k \approx N$ for N large),

$$\frac{\partial \langle O_{\gamma} \rangle^{(k+1)}}{\partial g_{\beta}^{(k)}} = \sum_{\alpha} \frac{\partial g_{\alpha}^{(k+1)}}{\partial g_{\beta}^{(k)}} \frac{\partial \langle O_{\gamma} \rangle^{(k+1)}}{\partial g_{\alpha}^{(k+1)}}
= \sum_{\alpha} T_{\alpha\beta}^{(k+1)} \frac{\partial \langle O_{\gamma} \rangle^{(k+1)}}{\partial g_{\alpha}^{(k+1)}}.$$
(21)

Here $T_{\alpha\beta}^{(k+1)}$ is the approximation to the stability matrix $T_{\alpha\beta}^*$ in eq. (14) obtained from the $k^{\rm th}$ layer, which should converge to $T_{\alpha\beta}^*$ for large k. Differentiating eq. (20) and using eq. (7), we obtain for the term on the RHS of eq. (21)

$$\frac{\partial \langle O_{\gamma} \rangle^{(k+1)}}{\partial g_{\alpha}^{(k+1)}} = \langle O_{\gamma} \rangle^{(k+1)} \langle O_{\alpha} \rangle^{(k+1)} - \langle O_{\gamma} O_{\alpha} \rangle^{(k+1)}. \tag{22}$$

To obtain a similar expression for the term on the LHS of eq. (21), we need to use t_{k+1} , the RBM for layer k+1, to relate $H_{\mathbf{x}}^{(k+1)}$ to $H_{\mathbf{x}}^{(k)}$, as in eq. (5). In Monte Carlo RG, this is where the RG blocking step enters. We obtain

$$\frac{\partial \langle O_{\gamma} \rangle^{(k+1)}}{\partial g_{\beta}^{(k)}} = \langle O_{\gamma} \rangle^{(k+1)} \langle O_{\beta} \rangle^{(k)}
- \operatorname{Tr}_{\mathbf{h}\mathbf{h}'} O_{\gamma}(\mathbf{h}) t_{k+1}(\mathbf{h}|\mathbf{h}') O_{\beta}(\mathbf{h}') q_{k}(\mathbf{h}'|\mathbf{x})
= \langle O_{\gamma} \rangle^{(k+1)} \langle O_{\beta} \rangle^{(k)} - \langle O_{\gamma} t_{k+1} O_{\beta} \rangle^{(k)}.$$
(23)

Thus, with a suitable choice of operators, given a specific input \mathbf{x} , by sampling the hidden nodes in the upper layers of the network (large k) and computing a set of expectation values of these operators, we can construct a set of linear equations for the elements of the stability matrix, $T_{\alpha\beta}^{(k+1)} \approx T_{\alpha\beta}^*$,

$$\sum_{\alpha} T_{\alpha\beta}^{(k+1)} \left[\langle O_{\gamma} \rangle^{(k+1)} \langle O_{\alpha} \rangle^{(k+1)} - \langle O_{\gamma} O_{\alpha} \rangle^{(k+1)} \right]$$
$$= \langle O_{\gamma} \rangle^{(k+1)} \langle O_{\beta} \rangle^{(k)} - \langle O_{\gamma} t_{k+1} O_{\beta} \rangle^{(k)}. \tag{24}$$

Hence, we may obtain a set of successive approximations for the eigenvalues Λ_{μ} and eigenvectors $\phi_{\alpha}^{(\mu)}$ in eq. (15).

By varying the subset of operators used in eq. (24), we can test the convergence of the estimates of the largest eigenvalues of the stability matrix associated with a given layer of the network and, by computing the eigenvalues from successive layers, we can investigate whether increasing depth of the network is associated with scaling behaviour of the corresponding operators, as in eq. (19).

If this scaling behaviour is due to an eigenvalue of magnitude greater than one, then even a tiny component of the initial data, which creates a non-zero coupling to the corresponding eigenvector, will cause instability of the fixed point if the network is deep. In the next section, we will explain how to identify the specific perturbation of the data which does this.

IV. GENERATION OF ADVERSARIAL DATA

We consider the exact conditional probability distribution, $p(\mathbf{y}|\mathbf{x})$, and hence the trained network, $q_N(\mathbf{y}|\mathbf{x})$, to be parametrised by the input data, \mathbf{x} , via the fixed-point Hamiltonian in eqs (8) and (9). The Fisher Information Matrix (FIM) [8, 11] is a metric which measures how a probability distribution changes along different directions in parameter space. If the fixed-point Hamiltonian has an unstable direction due to a relevant operator, then this will correspond to the direction in parameter space along which the probability distribution changes fastest, i.e., the stiffest direction. Hence, to generate adversarial data we need to compute the eigenvector associated with the largest eigenvalue of the Fisher information matrix for the fixed-point probability distribution. A tiny

admixture of this component in the data will lead to erroneous classification if the network is too deep.

The FIM is defined in terms of the Kullback-Liebler divergence, D_{KL} , which is itself a measure of how distinguishable two probability distributions p_1 and p_2 are, using data sampled from p_1 :

$$D_{\mathrm{KL}}(p_1||p_2) = \mathrm{Tr}_{\mathbf{x}} p_1(\mathbf{x}) \log \left[\frac{p_1(\mathbf{x})}{p_2(\mathbf{x})} \right].$$
 (25)

This enables us to measure the dependence of $q_N(\mathbf{y}|\mathbf{x})$ on \mathbf{x} by considering two nearby data points, \mathbf{x} and \mathbf{x}' , and expanding the Kullback-Liebler divergence as a Taylor series in their difference [8],

$$D_{\mathrm{KL}}(q_{N}(\mathbf{y}|\mathbf{x})||q_{N}(\mathbf{y}|\mathbf{x}')) = \mathrm{Tr}_{\mathbf{y}}q_{N}(\mathbf{y}|\mathbf{x})\log\left[\frac{q_{N}(\mathbf{y}|\mathbf{x})}{q_{N}(\mathbf{y}|\mathbf{x}')}\right]$$
$$= \frac{1}{2}\sum_{ij}(x'_{i}-x_{i})F_{ij}^{(N)}(x'_{j}-x_{j})$$
$$+ \dots, \tag{26}$$

where, using eqs (4), (7) and (20),

$$F_{ij}^{(N)} = \frac{\partial^{2}}{\partial x_{i}' \partial x_{j}'} D_{KL} \left(q_{N}(\mathbf{y}|\mathbf{x}) || q_{N}(\mathbf{y}|\mathbf{x}') \right) \Big|_{\mathbf{x}' = \mathbf{x}}$$

$$= \left\langle \frac{\partial H_{\mathbf{x}}^{(N)}}{\partial x_{i}} \frac{\partial H_{\mathbf{x}}^{(N)}}{\partial x_{j}} \right\rangle^{(N)} - \left\langle \frac{\partial H_{\mathbf{x}}^{(N)}}{\partial x_{i}} \right\rangle^{(N)} \left\langle \frac{\partial H_{\mathbf{x}}^{(N)}}{\partial x_{j}} \right\rangle^{(N)}$$

$$= \sum_{\alpha \alpha'} \frac{\partial g_{\alpha}^{(N)}}{\partial x_{i}} \frac{\partial g_{\alpha'}^{(N)}}{\partial x_{j}} \left[\left\langle O_{\alpha} O_{\alpha'} \right\rangle^{(N)} - \left\langle O_{\alpha} \right\rangle^{(N)} \left\langle O_{\alpha'} \right\rangle^{(N)} \right]$$
(27)

is the FIM.

Having chosen a subspace of operators, $\{O_{\alpha}\}$, in which to express the effective Hamiltonians in each layer of the network (via eq. (7)), such that sufficiently precise estimates of the largest eigenvalues of the stability matrix associated with each layer, $T_{\alpha\beta}^{(k+1)}$, are obtained by the method in Section III, we can compute the FIM using the chain rule:

$$\frac{\partial g_{\alpha}^{(N)}}{\partial x_i} = \sum_{\beta} \left[T^{(N)} \dots T^{(2)} \right]_{\alpha\beta} \frac{\partial g_{\beta}^{(1)}}{\partial x_i}.$$
 (28)

The last derivative may be computed from eq. (6) and the explicit form of the RBM for the first layer, e.g., eq. (1), by expressing $\mathbf{h}_{1}^{\mathrm{T}}\mathbf{W}_{1}\mathbf{x} + \mathbf{a}_{1}^{\mathrm{T}}\mathbf{h}$ in terms of $\{O_{\alpha}\}$ to determine $g_{\beta}^{(1)}(\mathbf{x})$. Alternatively, we may solve the system of linear equations in eq. (24) for the first layer:

$$\sum_{\alpha} \frac{\partial g_{\alpha}^{(1)}}{\partial x_{i}} \left[\langle O_{\gamma} \rangle^{(1)} \langle O_{\alpha} \rangle^{(1)} - \langle O_{\gamma} O_{\alpha} \rangle^{(1)} \right]$$

$$= \operatorname{Tr}_{\mathbf{h}} O_{\gamma}(\mathbf{h}) \frac{\partial}{\partial x_{i}} t_{1}(\mathbf{h}|\mathbf{x}). \tag{29}$$

Thus, computing the FIM involves translating the input data into the chosen coupling-constant subspace via the

RBM in the first layer, taking the product of the stability matrices associated with each layer in eq. (28), and computing a set of expectation values of operators in this subspace using the full network. The fact that the FIM is built from the product of stability matrices connects the scaling behaviour in eqs (13) and (19) with the sloppiness of the FIM spectrum observed in [8], i.e., all but a fixed number of eigenvalues of the FIM are scaled to very small values, resulting in a fixed number of stiff (relevant) directions and many sloppy (irrelevant) directions.

We may then obtain the eigenvector corresponding to the largest eigenvalue of the FIM. For deep networks which have an unstable fixed point, small perturbations of the corresponding input data in the direction of this eigenvector of the FIM are likely to result in the data being misclassified.

V. CONCLUSIONS

In this paper we have extended the formal analogy between deep learning and the RG. This enabled us to interpret the classification problem in terms of a sequence of effective Hamiltonians associated with the layers of the network, whose couplings are determined by the input data, and the convergence of this sequence to a distinct fixed point for each distinct output of the classification problem. Input data with the same classification should all result in a flow in the coupling-constant space of these effective Hamiltonians which converges to the same fixed

point.

This exposed the possibility that a fixed point might have an unstable direction, associated with a relevant operator in the language of the RG. Small perturbations of the input data that create a non-zero overlap with such a relevant operator tend to be amplified exponentially by successive layers of the network, so that the flow diverges from the correct fixed point and the data may be misclassified. Fixed points with no relevant operators are not prone to this sort of error and will only misclassify data that is ambiguous even to a human.

Convergence of the sequence of effective Hamiltonians to a fixed point for a given subset of the input data is guaranteed by the training (assuming this converges to a good approximation to the exact conditional probability distribution). The existence of a relevant direction at the fixed point is associated with sensitivity to perturbations which take the data outside of this subset. The size of the amplification factor for the perturbation is greater the greater the depth of the network. Hence, the vulnerability to misclassification is directly due to the depth of a network.

The key result of this paper is a method for diagnosing whether a given trained network is vulnerable to adversarial attack, or simply misclassification due to noise in the data, based only on computing expectation values of operators in the hidden and output layers. The computations are fairly demanding, but may be justified for a network that is intended for use in safety-critical applications.

G.E. Hinton, S. Osindero and Y-W. Teh, Neural Comput., 18, 1527 (2006).

^[2] Y. LeCun, Y. Bengio and G. Hinton, Nature, 521, 436 (2015).

^[3] R. Salakhutdinov, Annu. Rev. Stat. Appl., 2, 361 (2015).

^[4] X. Yuan, P. He, Q. Zhu, R.R. Bhat and X. Li, Adversarial Examples: Attacks and Defenses for Deep Learning, arXiv:1712.07107.

^[5] P. Mehta and D.J. Schwab, An exact mapping between the Variational Renormalization Group and Deep Learning, arXiv:1410.3831.

^[6] H.W. Lin, M. Tegmark and D. Rolnick, J. Stat. Phys., 168, 1223 (2017).

^[7] D.J. Schwab and P. Mehta, Comment on "Why does deep and cheap learning work so well?", arXiv:1609.03541.

^{8]} B.B. Machta, R. Chachra, M. Transtrum and J.P. Sethna, Science **342**, 604 (2013).

^[9] S-K. Ma, Phys. Rev. Lett., **37**, 461 (1976).

^[10] R.H. Swendsen, Phys. Rev. Lett., 42, 859 (1979).

^[11] A. Raju, B.B. Machta and J.P. Sethna, Information geometry and the renormalization group, arXiv:1710.05787.