

UNDERSTANDING WHY NEURAL NETWORKS GENERALIZE WELL THROUGH GSNR OF PARAMETERS

Anonymous authors

Paper under double-blind review

ABSTRACT

1 As deep neural networks (DNNs) achieve tremendous success across many appli-
 2 cation domains, researchers tried to explore in many aspects on why they gener-
 3 alize well. In this paper, we provide a novel perspective on these issues using
 4 the gradient signal to noise ratio (GSNR) of parameters during training process
 5 of DNNs. The GSNR of a parameter is simply defined as the ratio between its
 6 gradient’s squared mean and variance, over the data distribution. Based on several
 7 approximations, we establish a quantitative relationship between model param-
 8 eters’ GSNR and the generalization gap. This relationship indicates that larger
 9 GSNR during training process leads to better generalization performance. Futher,
 10 we show that, different from that of shallow models (e.g. logistic regression, sup-
 11 port vector machines), the gradient descent optimization dynamics of DNNs nat-
 12 urally produces large GSNR during training, which is probably the key to DNNs’
 13 remarkable generalization ability.

14 1 INTRODUCTION

15 Deep neural networks typically contain far more trainable parameters than training samples, which
 16 seems to easily cause a poor generalization performance. However, in fact they usually exhibit re-
 17 markably small generalization gaps. Traditional generalization theory such as VC dimension (Vap-
 18 nik & Chervonenkis, 1991) or Rademacher complexity (Bartlett P L, 2002) cannot explain its mech-
 19 anism. Extensive research focuses on the generalization ability of DNNs (Neyshabur et al., 2017; S.
 20 et al., 2018; Keskar et al., 2016; Dinh et al., 2017; Hoffer et al., 2017; R et al., 2018; Dziugaite &
 21 Roy, 2017; etc, 2018; Kawaguchi et al., 2017; Advani & Saxe, 2017).

22 Unlike that of shallow models such as logistic regression or support vector machines, the global
 23 minimum of high-dimensional and non-convex DNNs cannot be found analytically, and can only
 24 be approximated by gradient descent and its variants (Zeiler, 2012; Kingma & Ba, 2014; Graves,
 25 2013). Previous work (Zhang et al., 2016; Hardt et al., 2015; Dziugaite & Roy, 2017) suggests that
 26 the generalization ability of DNNs is closely related to gradient descent optimization. For example,
 27 Hardt et al. (2015) claims that any model trained with stochastic gradient descent (SGD) for reason-
 28 able epochs would exhibit small generalization error. Their analysis is based on the smoothness of
 29 loss function. In this work, we attempt to understand the generalization behavior of DNNs through
 30 GSNR and reveal how GSNR affects the training dynamics of gradient descent.

31 The GSNR of a parameter is defined as the ratio between its gradient’s squared mean and variance
 32 over the data distribution. Previous work tried to use GSNR to conduct theoretical analysis on
 33 deep learning. For example, Rainforth et al. (2018) used GSNR to analyze variational bounds in
 34 unsupervised DNNs such as variational auto-encoder (VAE). Here we focus on analyzing the relation
 35 between GSNR and the generalization gap.

36 Intuitively, GSNR measures the similarity of a parameter’s gradients among different training sam-
 37 ples. Large GSNR implies that most training samples agree on the optimization direction of this
 38 parameter, thus the parameter is more likely to be associated with a meaningful “pattern” and we
 39 assume its update could lead to a better generalization. In this work, we prove that the GSNR is
 40 strongly related to the generalization performance, and larger GSNR means a better generalization.

41 To reveal the mechanism of DNNs’ good generalization ability, we show that the gradient descent
 42 optimization dynamics of DNN naturally leads to large GSNR of model parameters, thus a good

43 generalization. Further, we give a complete analysis and a detailed interpretation to this phenomenon.
44 We believe this is probably the key to DNNs remarkable generalization ability.

45 In the remainder of this paper we first analyze the relation between GSNR and generalization (Sec-
46 tion 2). We then show how the training dynamics lead to large GSNR of model parameters experi-
47 mentally and analytically in Section 3.

48 2 LARGER GSNR LEADS TO BETTER GENERALIZATION

49 In this section, we establish a quantitative relation between the GSNR of model parameters and the
50 generalization gap based on several proper approximations. This relationship indicates that larger
51 GSNR during training process leads to better generalization performance.

52 2.1 GRADIENTS SIGNAL TO NOISE RATIO

53 Consider a data distribution $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ from which each sample $s = (x, y)$ is drawn, and a model
54 $\hat{y} = f(x, \theta)$ trying to predict y from x , parameterized by θ . We use $\mathbf{g}_s(\theta)$ to denote parameters'
55 gradient w.r.t sample s . Given the data distribution \mathcal{Z} , we can evaluate the (sample-wise) mean
56 and variance of $\mathbf{g}_s(\theta)$. We denote them as $\mu(\theta) = \mathbb{E}_{s \sim \mathcal{Z}}(\mathbf{g}_s(\theta))$ and $\rho^2(\theta) = \text{Var}_{s \sim \mathcal{Z}}(\mathbf{g}_s(\theta))$,
57 respectively.

58 The gradient signal to noise ratio (GSNR) of q -th model parameter is defined as:

$$r_q(\theta) \equiv \frac{\mu_q^2(\theta)}{\rho_q^2(\theta)} \quad (1)$$

59 We can see, at a particular point of the parameter space, GSNR measures the similarity of the param-
60 eter gradients among different data samples.

61 2.2 ONE-STEP GENERALIZATION RATIO

62 In this section we introduce a new concept to help us measure the generalization performance during
63 gradient descent optimization, which we call one-step generalization ratio (OSGR). We denote the
64 training set as $D = \{(x_0, y_0), \dots, (x_n, y_n)\} \sim \mathcal{Z}^n$ consists of n samples drawn from \mathcal{Z} , and the
65 test set as $D' = \{(x'_0, y'_0), \dots, (x'_{n'}, y'_{n'})\} \sim \mathcal{Z}^{n'}$. In practice we use the loss on D' to estimate the
66 generalization loss. We denote the empirical training loss and test loss respectively as

$$L[D] = \frac{1}{n} \sum_{i=0}^n L(y_i, f(x_i, \theta)), \quad L[D'] = \frac{1}{n'} \sum_{i=0}^{n'} L(y'_i, f(x'_i)) \quad (2)$$

67 Then the empirical generalization gap is defined by $e = L[D'] - L[D]$.

68 In gradient descent optimization, both the training and test loss would decline step by step. We
69 use $\Delta L[D]$ and $\Delta L[D']$ to denote the one-step training and test loss decrease during training, re-
70 spectively. Let's consider the ratio between the expectations of $\Delta L[D']$ and $\Delta L[D]$ of one single
71 training step, which we denote as $\mathcal{R}(\mathcal{Z})$.

$$\mathcal{R}(\mathcal{Z}) \equiv \frac{E_{D, D' \sim \mathcal{Z}^n}(\Delta L[D'])}{E_{D \sim \mathcal{Z}^n}(\Delta L[D])} \quad (3)$$

72 Note that the expectation of $\Delta L[D']$ is over D and D' , and it's because the optimization step, aka
73 training process, is on D . We refer $\mathcal{R}(\mathcal{Z})$ as OSGR of gradient descent optimization. Statistically
74 the training loss decreases faster than the test loss and OSGR would be less than 1, which usually
75 results in a non-zero generalization gap at the end of training. If OSGR is large in the whole training
76 process, generalization gap will be small when training completes, implying good generalization
77 ability of the model.

78 2.3 RELATION BETWEEN GSNR AND OSGR

79 In this section, we focus on deriving the relation between GSNR and OSGR, further to prove the
80 GSNR is strongly related to model generalization performance.

81 For simplicity, we assume the sizes of training and test datasets are equal, $n = n'$. In gradient
82 descent optimization, we compute the mean of gradient over the training dataset, which we denote
83 as $\mathbf{k}(\theta, D) = \frac{1}{n} \sum_{i=0}^n \mathbf{g}_i(\theta)$, and use it as the opposite descent direction. Note the difference
84 between $\mu(\theta)$ and $\mathbf{k}(\theta, D)$, where $\mu(\theta)$ is the gradient mean over a data distribution and $\mathbf{k}(\theta, D)$ is
85 the empirical gradient mean of a training dataset.

86 The gradient mean of the training and test datasets are, respectively

$$\mathbf{k}(\theta, D) = \frac{1}{n} \sum_{i=0}^n \mathbf{g}_i(\theta) = \frac{\partial L[D]}{\partial \theta} \quad , \quad \mathbf{k}'(\theta, D') = \frac{1}{n} \sum_{i=0}^n \mathbf{g}'_i(\theta) = \frac{\partial L[D']}{\partial \theta} \quad (4)$$

87 where $\mathbf{g}'_i(\theta)$ is the gradient w.r.t the i -th sample of the test dataset.

88 Both the training and test dataset are randomly generated from the same distribution \mathcal{Z}^n , so we
89 can treat $\mathbf{k}(\theta, D)$ and $\mathbf{k}'(\theta, D')$ as random variables. At the beginning of the optimization process,
90 θ is randomly initialized thus independent of D , so $\mathbf{k}(\theta, D)$ and $\mathbf{k}'(\theta, D')$ would obey the same
91 distribution. After a period of training, the model parameters begin to fit the training dataset and
92 become a function of D , $\theta = \theta(D)$, then the distribution of $\mathbf{k}(\theta(D), D)$ and $\mathbf{k}'(\theta(D), D')$ becomes
93 different. However we will make our derivation under the non-overfitting limit approximation 2.3.1
94 stated as below.

95 **Assumption 2.3.1 (Non-overfitting limit approximation)** *In the early training stage, the mean of*
96 *gradients w.r.t the training dataset and test dataset $\mathbf{k}(\theta(D), D)$ and $\mathbf{k}'(\theta(D), D')$ obey the same*
97 *distribution, and we denote their mean and variance as $\mu(\theta)$ and $\sigma^2(\theta)$, i.e.*

$$\mathbb{E}_{D \sim \mathcal{Z}^n} [\mathbf{k}(\theta(D), D)] = \mathbb{E}_{D, D' \sim \mathcal{Z}^n} [\mathbf{k}'(\theta(D), D')] = \mu(\theta) \quad (5)$$

98

$$\text{Var}_{D \sim \mathcal{Z}^n} [\mathbf{k}(\theta(D), D)] = \text{Var}_{D, D' \sim \mathcal{Z}^n} [\mathbf{k}'(\theta(D), D')] = \sigma^2(\theta) \quad (6)$$

99 For simplicity, we denote $\mathbf{k}(\theta(D), D)$ and $\mathbf{k}(\theta(D), D')$ as $\mathbf{k}(\theta)$ and $\mathbf{k}'(\theta)$ respectively. We can get
100 the following relation:

$$\mu(\theta) = \mathbb{E}_{D \sim \mathcal{Z}^n} [\mathbf{k}(\theta)] = \mathbb{E}_{D \sim \mathcal{Z}^n} \left[\frac{1}{n} \sum_{i=0}^n \mathbf{g}_i(\theta) \right] = \mathbb{E}_{s \sim \mathcal{Z}} [\mathbf{g}_s(\theta)] \quad (7)$$

$$\sigma^2(\theta) = \text{Var}_{D \sim \mathcal{Z}^n} [\mathbf{k}(\theta)] = \text{Var}_{D \sim \mathcal{Z}^n} \left[\frac{1}{n} \sum_{i=0}^n \mathbf{g}_i(\theta) \right] = \frac{1}{n} \text{Var}_{s \sim \mathcal{Z}} [\mathbf{g}_s(\theta)] \equiv \frac{1}{n} \rho^2(\theta) \quad (8)$$

101 where $\sigma^2(\theta)$ is the variance of the gradient mean of a training dataset while $\rho^2(\theta)$ is the variance of
102 the gradient of a single sample.

103 In one gradient descent step, the model parameter is updated by $\Delta\theta = \theta_{t+1} - \theta_t = -\lambda \mathbf{k}(\theta)$ where
104 λ is the learning rate. If λ is small enough, the one-step training and test loss decrease can be
105 approximated as

$$\Delta L[D] \approx -\Delta\theta \cdot \frac{\partial L[D]}{\partial \theta} + O(\lambda^2) = \lambda \mathbf{k}(\theta) \cdot \mathbf{k}(\theta) + O(\lambda^2) \quad (9)$$

$$\Delta L[D'] \approx -\Delta\theta \cdot \frac{\partial L[D']}{\partial \theta} + O(\lambda^2) = \lambda \mathbf{k}(\theta) \cdot \mathbf{k}'(\theta) + O(\lambda^2) \quad (10)$$

106 Usually there are some differences between the directions of $\mathbf{k}(\theta)$ and $\mathbf{k}'(\theta)$, so statistically $\Delta L[D]$
107 tends to be larger than $\Delta L[D']$ and the generalization gap would increase during training. When
108 $\lambda \rightarrow 0$, in one single training step the empirical generalization gap increases by:

$$\Delta e \approx \lambda \mathbf{k}(\theta) \cdot \mathbf{k}(\theta) - \lambda \mathbf{k}(\theta) \cdot \mathbf{k}'(\theta) = \lambda (\mu(\theta) + \epsilon) (\mu(\theta) + \epsilon - \mu(\theta) - \epsilon') \quad (11)$$

$$= \lambda (\mu(\theta) + \epsilon) (\epsilon - \epsilon') \quad (12)$$

109 We replace the random variables, $\mathbf{k}(\theta) = \mu(\theta) + \epsilon$, $\mathbf{k}'(\theta) = \mu(\theta) + \epsilon'$. ϵ and ϵ' are random
 110 variables with zero mean and variance $\sigma^2(\theta)$. Since $E(\epsilon') = E(\epsilon) = 0$, ϵ and ϵ' are independent,
 111 the expectation of Δe is

$$E_{D, D' \sim \mathcal{Z}^n}(\Delta e) = E(\lambda \epsilon \cdot \epsilon) + O(\lambda^2) = \lambda \sum_{q=1}^Q \sigma_q^2 + O(\lambda^2) \quad (13)$$

112 where Q is the number of parameters and σ_q^2 is the variance of mean gradient of the q -th parameter.
 113 Consider the expectation of $\Delta L[D]$ and $\Delta L[D']$ when $\lambda \rightarrow 0$

$$E_{D \sim \mathcal{Z}^n}(\Delta L[D]) \approx \lambda E_{D \sim \mathcal{Z}^n}(\mathbf{k}(D) \cdot \mathbf{k}(D)) = \lambda \sum_{q=1}^Q E_{D \sim \mathcal{Z}^n}(k_q^2(D)) \quad (14)$$

$$E_{D, D' \sim \mathcal{Z}^n}(\Delta L[D']) = E_{D, D' \sim \mathcal{Z}^n}(\Delta L[D] - \Delta e) \approx \lambda \sum_{q=1}^Q (E_{D \sim \mathcal{Z}^n}(k_q^2(D)) - \sigma_q^2) \quad (15)$$

$$= \lambda \sum_{q=1}^Q (E_{D \sim \mathcal{Z}^n}(k_q^2(D)) - \rho_q^2/n) + O(\lambda^2) \quad (16)$$

114 Now we are ready to derive the relation between GSNR and OSGR, by substituting 16 and 14 into
 115 3:

$$\mathcal{R}(\mathcal{Z}) = 1 - \frac{\sum_{q=1}^Q \rho_q^2}{n \sum_{q=1}^Q E_{D \sim \mathcal{Z}^n}(k_q^2(D))} \quad (17)$$

116 Now the right hand side of equation (17) can be estimated using only the samples in the training
 117 datasets. Because the gradient average $k_q(D)$ and sample-wise gradient variance ρ_q^2 can both be
 118 computed within one training dataset. We will elaborate on this estimation method in section 2.4.

119 Reformulate equation (17) as:

$$\mathcal{R}(\mathcal{Z}) = 1 - \frac{1}{n} \sum_{q=1}^Q \frac{E_{D \sim \mathcal{Z}^n}(k_q^2(D))}{\sum_{q'=1}^Q E_{D \sim \mathcal{Z}^n}(k_{q'}^2(D))} \frac{\rho_q^2}{E_{D \sim \mathcal{Z}^n}(k_q^2(D))} \quad (18)$$

$$= 1 - \frac{1}{n} \sum_{q=1}^Q \frac{E_{D \sim \mathcal{Z}^n}(k_q^2(D))}{\sum_{q'=1}^Q E_{D \sim \mathcal{Z}^n}(k_{q'}^2(D))} \frac{1}{r_q + \frac{1}{n}} \quad (19)$$

120 where $E_{D \sim \mathcal{Z}^n}(k_q^2(D)) = \text{Var}_{D \sim \mathcal{Z}^n}(k_q(D)) + E_{D \sim \mathcal{Z}^n}^2(k_q(D)) = \frac{1}{n} \rho_q^2 + \mu_q^2$.

121 We define $\Delta L_q[D]$ to be the training loss decrease caused by updating the q -th parameter. We can
 122 show that when λ is very small $\Delta L_q[D] = \lambda k_q^2(D) + O(\lambda^2)$. Therefore when $\lambda \rightarrow 0$

$$\mathcal{R}(\mathcal{Z}) = 1 - \frac{1}{n} \sum_{q=1}^Q W_q \frac{1}{r_q + \frac{1}{n}}, \quad \text{where } W_q = \frac{E_{D \sim \mathcal{Z}^n}(\Delta L_q[D])}{E_{D \sim \mathcal{Z}^n}(\Delta L[D])} \quad \text{with } \sum_{q=1}^Q W_q = 1 \quad (20)$$

123 Equation (20) shows that the GSNR (equation 1) plays a crucial role in the model's generalization
 124 ability. In the non-overfitting limit approximation and with a small enough learning rate, the one-step
 125 generalization ratio in gradient descent equals 1 minus the weighted average of $\frac{1}{r_q + \frac{1}{n}}$ over all model
 126 parameters divided by n . The weight is proportional to the expectation of the training loss decrease
 127 caused by updating that parameter. This result implies that larger GSNR of model parameters during
 128 training leads to smaller generalization gap growth and better generalization performance of the final
 129 model. Also note when $n \rightarrow \infty$, we have $\mathcal{R} \rightarrow 1$ and hence good generalization performance.

130 2.4 EXPERIMENTAL VERIFICATION OF THE RELATION BETWEEN GSNR AND OSGR

131 The relation between GSNR and OSGR, i.e. equation 17 or 20, is solid and can be accurately verified
 132 using any datasets if: (1) The dataset include enough number of samples to construct enough number

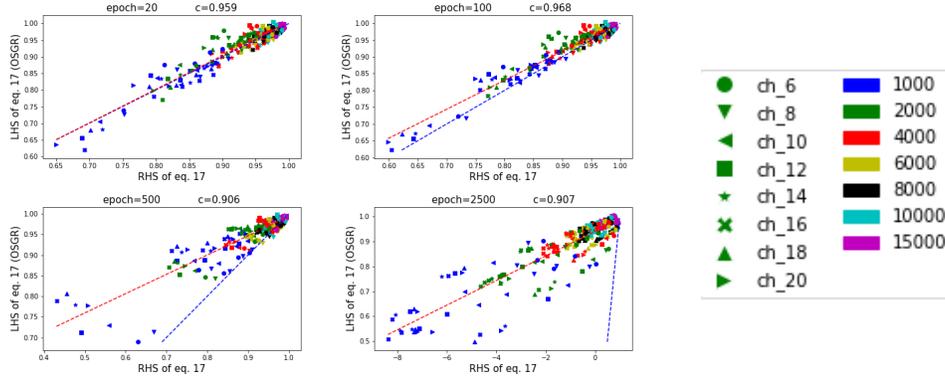


Figure 1: Left hand (LHS) and right side (RHS) of equation 17. Points are drawn under different experiment setting. Left figure: LHS vs RHS relation at epoch 20, 100, 500, 2500. Each point is drawn by LHS and RHS computed at the given epoch under different model structure (channel number) or training data size; red dotted line is the line of best fit computed by Least Square Method; blue dotted line is the line of reference representing LHS = RHS; the value of c in each title represents the Pearson correlation coefficient between LHS and RHS computed by points in figure. Right figure: The legend, different symbols and colors stand for different channel number and training data size respectively. Different random noise is not presented.

of training datasets and a large enough test dataset so that the estimation of ρ_q^2 , $E_{D \sim \mathcal{Z}^n}(k_q^2(D))$ and OSGR can be accurate. (2) The learning rate is small enough. (3) In the early training stage of gradient descent.

To verify equation 17, here we show how to estimate the left and right hand sides of it. Suppose we have M number of training datasets with n samples each and a large test dataset with n' samples. We initialize a model and train it separately on the M training datasets and test it with the same test dataset. For the t -th training iteration, we denote the training loss and test loss of the model trained on the m -th training dataset as L_{tm} and L'_{tm} , respectively. Then the left hand side, i.e. OSGR, of the t -th iteration can be estimated as

$$\mathcal{R}_t(\mathcal{Z}) \approx \frac{\sum_{m=1}^M L'_{t+1m} - L'_{tm}}{\sum_{m=1}^M L_{t+1m} - L_{tm}} \quad (21)$$

For the model trained on the m -th training dataset, we can compute the t -th step gradient mean and sample-wise gradient variance of q -th parameter on the corresponding training dataset, denoted as \hat{k}_{qmt} and $\hat{\rho}_{qmt}^2$, then the right hand side of equation 17 can be estimated

$$E_{D \sim \mathcal{Z}^n}(k_{qt}^2(D)) \approx \frac{1}{M} \sum_{m=1}^M \hat{k}_{qmt}^2, \quad \rho_{qt}^2 \approx \frac{1}{M} \sum_{m=1}^M \hat{\rho}_{qmt}^2 \quad (22)$$

Note that using equation 22, the right hand side can be computed merely with the samples in the training datasets.

To verify equation 17, we carry out the experiment on MNIST training dataset with simple CNNs which consists of 2 Conv-Relu-MaxPooling blocks and 2 fully-connected layers. First, to estimate equation 22 with $M = 10$, we randomly sample 10 training datasets with n samples each and a test set with 10000 samples. Then to cover different conditions, we (1) set $n \in \{1000, 2000, 4000, 6000, 8000, 10000, 15000\}$, respectively; (2) inject noise in the datasets by randomly changing the labels with proportion $p_{random} \in \{0.0, 0.1, 0.2, 0.3, 0.5\}$; (3) experiment on various model structures which are different only at the number of channels in layers. For model details, please see Appendix A. Moreover, we use the gradient descent training (Not SGD), with a small learning rate of 0.001. The left and right hand sides of 17 at different epochs are shown in Figure 1, where each point represents one specific combination of above settings.

At the beginning of training, the data points closely distributed along the dashed line of $y = x$ which shows that equation 17 fits quite well under a variety of different settings. As training proceeds, equation 17 gradually losses its accuracy because the non-overfitting limit approximation no longer holds, but strong positive correlation between the left and right hand sides of equation 17 remains even when the training converges (at epoch 2500).

162 Through analytically derivations and experimental verifications, we prove that GSNR is strongly
 163 related to OSGR which indicates the generalization ability, thus demonstrate that the larger GSNR
 164 during training leads to better generalization performance.

165 3 TRAINING DYNAMICS OF DNNs NATURALLY LEADS TO LARGE GSNR

166 In this section, we analyze and explain one observed phenomenon: the parameters' GSNR of DNN
 167 models rises in the early stages of training, while the GSNR of shallow models such as logistic
 168 regression or support vector machines declines during the entire training process. This difference of
 169 behaviors provides GSNR large practical values during DNN training, which in turn is associated
 170 with good generalization. We analyze the dynamics behind this phenomenon both experimentally
 171 and theoretically, which deepens our understanding of the good generalization ability of DNNs.

172 3.1 GSNR BEHAVIOR OF DNNs TRAINING

173 For shallow models, the GSNR of parameters decreases in the whole training process because gra-
 174 dients become small as learning proceeds when the optimizer finds the local minimum. But for
 175 DNN it is not the case. We trained DNNs on the CIFAR datasets and computed the GSNR aver-
 176 aged over all model parameters. Because $E_{D \sim \mathcal{Z}^n}(k_q^2(D)) = \frac{1}{n}\rho_q^2 + \mu_q^2$ and we assume n is large,
 177 $E_{D \sim \mathcal{Z}^n}(k_q^2(D)) \approx \mu_q^2$. In the case of only one large training datasets, we estimate GSNR of t -th
 178 iteration by

$$\hat{r}_{qt}(\theta) \approx \hat{k}_{qt}^2(\theta) / \hat{\rho}_{qt}^2(\theta) \quad (23)$$

179 As shown in Figure 2, the GSNR starts out low with randomly initialized parameters. As learning
 180 progresses, the GSNR increases in the early training stage and stays at a high level in the whole
 181 learning process. We also computed the proportion of the samples that have the same gradient
 182 sign (positive or negative) for each parameter, denoted as p_{same_sign} . In Figure 2c, we plot the
 183 mean timeseries of this proportion for all the parameters. This value increases from about 50% (half
 184 positive half negative due to random initialization) at beginning to about 56% finally, which indicates
 185 that for most parameters, the gradient signs on different samples become to have a certain degree
 186 of consistency as the training proceeds. We hypothesis this is because meaningful features begin to
 187 emerge in the learning process and the gradients of the weights on these features tend to have the
 188 same sign among different samples.

189 Previous research (Zhang et al., 2016) showed that DNNs achieved zero training loss by memorizing
 190 training samples even if the labels were randomized. We also plot the average GSNR for model
 191 trained using data with randomized labels in Figure 2 and find that the GSNR stays at a low level
 192 throughout the training process. Although the training loss of both the original and randomized
 193 labels go to zero (not shown), the GSNR curves clearly distinguish between these two cases and
 194 reveal the lack of meaningful patterns in the latter one. We believe this is the reason why DNNs
 195 trained on real and random data lead to completely different generalization behavior.

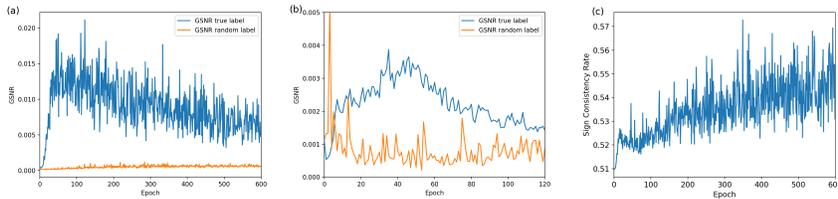


Figure 2: (a): GSNR curves generated by a simple network based on real and random data. An obvious upward process in the early training stage was observed for real data.(b): Same figure but for ResNet18 (c): average of p_{same_sign} for the same model in (a)

196 3.2 TRAINING DYNAMICS BEHIND THE GSNR BEHAVIOR

197 In this section we show that the feature learning ability of DNNs is the key reason why the GSNR
 198 curve behavior of DNNs is different from that of shallow models during the gradient descent training
 199 process. To demonstrate this clearly, a simple two-layer perceptron regression model is constructed.

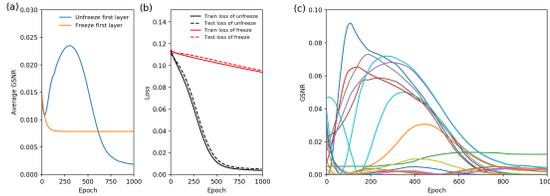


Figure 3: Average GSNR (a) and loss (b) curves for frozen and non-frozen the first layer. GSNR curves (c) of individual parameters for the non-frozen case

200 A synthetic dataset is generated as following. Each data point is constructed i.i.d. using $y = x_0x_1 + \epsilon$,
 201 where x_0 and x_1 are drawn from uniform distribution $[-1, 1]$ and ϵ is drawn from uniform distribu-
 202 tion $[-0.01, 0.01]$. The training set and test set sizes are 200 and 10000, respectively. We use a very
 203 simple 2-layer MLP structure with 2 input, 20 hidden neurons and 1 output.

204 We randomly initialized the model parameters and trained the model on the synthetic training dataset.
 205 As a control setup we also tried to freeze model weights in the first layer to prevent it from learning
 206 features. Note that a two layer MLP with the first layer frozen is equivalent to a linear regression
 207 model. That is, regression weights are learned on the second layer using fixed features extracted by
 208 the first layer. We plot the average GSNR of the second layer parameters for both the frozen and
 209 non-frozen cases. Figure 3 shows that in the non-frozen case, the average GSNR over parameters
 210 of the second layer has a significant upward process, whereas in the frozen case the average GSNR
 211 decreases in the beginning and remains at a low level during whole training process.

212 GSNR curve of individual parameters of the second layer of the non-frozen case are shown in Figure
 213 3. It can be observed that the GSNR for some parameters have a significant upward process. We
 214 computed the Pearson correlation between the features extracted from the first layer and y corre-
 215 sponding to these parameters at the beginning of the training process and the maximum point of the
 216 GSNR curve. We can see that the learning process changes these features from a random initialized
 217 one to a "good" feature with stronger correlation with y as shown in Table 1. This shows that in the
 218 training of DNN, the feature learning process is closely related to the GSNR increasing process.

219 3.3 ANALYTICAL ANALYSIS OF TRAINING DYNAMICS BEHIND DNNs' GSNR BEHAVIOR

220 In this section, we will analytically investigate the training dynamics behind the GSNR curve behav-
 221 ior of DNNs. In the case of fully connected network structure, we can analytically show that the
 222 numerator of GSNR, the squared gradient mean of model parameters, tends to increase in the early
 223 training stage through feature learning.

224 In a fully connected network, whose parameters are denoted as $\theta \equiv$
 225 $\{\omega_1, \mathbf{b}_1, \omega_2, \mathbf{b}_2, \dots, \omega_{l_{max}}, \mathbf{b}_{l_{max}}\}$, ω_1, \mathbf{b}_1 is the weight and bias of the first layer, and so on.
 226 We denote the activations of the l -th layer as $[a_l]_{c_1}(\phi)$, $c_1 = \{1, 2, \dots, C_1\}$, where C_1 is the number
 227 of activations of this layer. ϕ is the collection of model parameters over all the layers before this
 228 layer, i.e. $\phi = \{\omega_1, \mathbf{b}_1, \omega_2, \mathbf{b}_2, \dots, \omega_{l-1}, \mathbf{b}_{l-1}\}$. In the forward pass of the i -th sample, $\mathbf{a}_l(\phi)$ will
 229 be multiplied by the weight of this layer ω_l , which can be expressed as matrix multiplication.

$$[o_l]_{ic_2} = \sum_{c_1} [\omega_l]_{c_2c_1} [a_l]_{ic_1}(\phi) \quad (24)$$

230 where \mathbf{o}_l is the output of the matrix multiplication of the l -th layer and $c_2 = \{1, 2, \dots, C_2\}$ where C_2
 231 is the number of activations of the next layer. We use \mathbf{k}_l to denote the gradient mean of weights of
 232 the l -th layer ω_l , i.e. $\mathbf{k}_l = \frac{1}{n} \sum_{i=1}^n \frac{L_i}{\partial \omega_l}$, where L_i is the loss of the i -th sample.

233 Here we analytically show that the feature learning ability of DNNs plays a crucial role in the GSNR
 234 increasing process. To be more precise, we show that the learning of features $\mathbf{a}_l(\phi)$, i.e. the learning
 235 of parameters ϕ , tends to increase the absolute value of \mathbf{k}_l . Let's consider, the one-step change of
 236 gradient mean $\Delta[k_l] = ([k_l])_{t+1} - ([k_l])_t$ with the learning rate $\lambda \rightarrow 0$. In one training step, θ will
 237 be updated by $\Delta\theta = \theta_{t+1} - \theta_t = -\lambda \mathbf{k}(\theta)$. Use linear approximation with $\lambda \rightarrow 0$, we have

$$\Delta[k_l]_{c_1c_2} \approx \sum_{q=1}^Q \frac{\partial [k_l]_{c_1c_2}}{\partial \theta_q} \Delta\theta_q = \sum_{q=1}^{Q_l} \frac{\partial [k_l]_{c_1c_2}}{\partial \phi_q} \Delta\phi_q + \sum_{q=Q_l+1}^Q \frac{\partial [k_l]_{c_1c_2}}{\partial \theta_q} \Delta\theta_q \quad (25)$$

238 where Q is the total number of model parameters and Q_l is the number of model parameters of all
 239 the layers before the l -th layer. We will focus on the first term of equation 25, i.e. the one-step

240 change of k_l caused by learning ϕ . Substituting $k_l = \frac{1}{n} \sum_{i=1}^n \frac{L_i}{\omega_l}$ and $\Delta\phi_q = (-\lambda \frac{1}{n} \sum_{i=1}^n \frac{\partial L_i}{\partial \phi_q})$ in
 241 to equation 25, we have

$$\Delta[k_l]_{c_1 c_2} = -\frac{\lambda}{n^2} \sum_{q=1}^{Q_l} [\omega_l]_{c_1 c_2} \left(\sum_{i=1}^n \frac{\partial L_i}{\partial [\omega_l]_{c_1 c_2}} \frac{\partial [a_l]_{i c_1}}{\partial \phi_q} \right)^2 + \text{other terms} \quad (26)$$

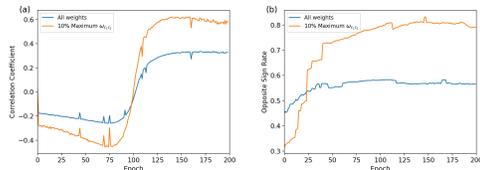
242 The detailed derivation of equation 26 can be found in Appendix B. We can see the first term (which
 243 is a summation of Q_l terms) in equation 26 has opposite sign with $[\omega_l]_{c_1 c_2}$. This term will make
 244 $\Delta[k_l]_{c_1 c_2}$ negatively correlated with $[\omega_l]_{c_1 c_2}$. We plot the correlation between $\Delta[k_l]_{c_1 c_2}$ with $[\omega_l]_{c_1 c_2}$
 245 for a model trained on MNIST for 200 epochs as shown in Figure 4a. In the early training stage,
 246 they are indeed negatively correlated. For top-10% weights with larger absolute values, the negative
 247 correlation is more significant.

248 Here we show that this negative correlation between $\Delta[k_l]_{c_1 c_2}$ and $[\omega_l]_{c_1 c_2}$ will tend to increase
 249 the absolute value of $[k_l]$ through an interesting mechanism. Consider the weights $[\omega_l]_{c_1 c_2}$ with
 250 $\{[\omega_l]_{c_1 c_2} > 0, [k_l]_{c_1 c_2} < 0\}$. Learning ϕ would tend to decrease $[k_l]_{c_1 c_2}$ and thus increase its
 251 absolute value because the first term in equation 26 is negative in this case. Learning $[\omega_l]_{c_1 c_2}$ would
 252 increase $[\omega_l]_{c_1 c_2}$ and its absolute value because $\Delta[\omega_l]_{c_1 c_2} = -\lambda[k_l]_{c_1 c_2}$ is positive in this case. This
 253 will form a positive feedback process, in which the numerator of GSNR, $([k_l]_{c_1 c_2})^2$, would increase
 254 and so is the GSNR. Similar things happen for the weights with $\{[\omega_l]_{c_1 c_2} < 0, [k_l]_{c_1 c_2} > 0\}$.

255 Then what about the weights with $\{[\omega_l]_{c_1 c_2} [k_l]_{c_1 c_2} > 0\}$? Here we show that the weights with
 256 $\{[\omega_l]_{c_1 c_2} [k_l]_{c_1 c_2} > 0\}$ tends to change into weights with $\{[\omega_l]_{c_1 c_2} [k_l]_{c_1 c_2} < 0\}$ during training
 257 process. Consider the case where $\{[\omega_l]_{c_1 c_2} > 0, [k_l]_{c_1 c_2} > 0\}$, the first term in equation 26 will be
 258 negative, learning ϕ tends to decrease $[k_l]_{c_1 c_2}$ or even change the sign of $[k_l]_{c_1 c_2}$. Another possibility
 259 is that learning $[\omega_l]_{c_1 c_2}$ tends to changes the sign of $[\omega_l]_{c_1 c_2}$ because $\Delta[\omega_l]_{c_1 c_2} = -\lambda[k_l]_{c_1 c_2}$ is
 260 negative in this case. Both case will change the weights with $\{[\omega_l]_{c_1 c_2} [k_l]_{c_1 c_2} > 0\}$ to weights with
 261 $\{[\omega_l]_{c_1 c_2} [k_l]_{c_1 c_2} < 0\}$. Same things happen for the weights with $\{[\omega_l]_{c_1 c_2} < 0, [k_l]_{c_1 c_2} < 0\}$.

262 Therefore $\{[\omega_l]_{c_1 c_2} [k_l]_{c_1 c_2} < 0\}$ is a more stable state than $\{[\omega_l]_{c_1 c_2} [k_l]_{c_1 c_2} > 0\}$ in the early
 263 training process. For a simple model trained on Mnist, We plot the proportion of weights that satisfy
 264 $\{[\omega_l]_{c_1 c_2} [k_l]_{c_1 c_2} < 0\}$ in Figure 4b and find that there are more weights with $\{[\omega_l]_{c_1 c_2} [k_l]_{c_1 c_2} < 0\}$
 265 than the opposite. Because weights with small absolute value easily change sign during training, we
 266 also plot this proportion for the top-10% weights with larger absolute values. It can be seen that for
 267 the weights with large absolute values, nearly 80% of these weights have opposite sign with their
 268 gradient mean as training proceeds. And for these weights, the numerator of GSNR, $([k_l]_{c_1 c_2})^2$,
 269 tends to increase through a positive feedback process in the early training stage of DNN.

270 Through discussions above, we further demonstrate that the learning of corresponding features tends
 271 to increase the GSNR of a weight based on analytical analysis. This GSNR increasing process leads
 272 to larger GSNR in the whole gradient descent training process of DNN, which in turn provides good
 273 generalization performance of DNN.



274 Figure 4: Dataset: Mnist. Left: Correlation between $\Delta[k_l]_{c_1 c_2}$ with $[\omega_l]_{c_1 c_2}$. Right : Ratio of weights that have opposite sign with its gradient mean.

Table 1: Pearson correlation of the features with y .

q	Beginning of training	Maximum of GSNR curve
1	-0.11	0.47
2	-0.33	0.53
3	-0.21	-0.27
4	0.07	0.40
5	0.11	0.44

275 4 SUMMARY

276 In this paper we investigated the relation between generalization of DNNs and the GSNR of the
 277 model parameters in gradient descent. We also analyzed the GSNR behavior and the mechanism
 278 behind it in the DNNs training process. Through our analysis, we hope to shed more light on the
 279 mechanisms behind DNNs impressive generalization ability.

280 REFERENCES

- 281 Madhu S Advani and Andrew M Saxe. High-dimensional dynamics of generalization error in neural
282 networks. *arXiv preprint arXiv:1710.03667*, 2017.
- 283 Mendelson S Bartlett P L. Rademacher and gaussian complexities: Risk bounds and structural
284 results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- 285 Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize
286 for deep nets. In *Proceedings of the 34th International Conference on Machine Learning-Volume*
287 *70*, pp. 1019–1028. JMLR. org, 2017.
- 288 Gintare Karolina Dziugaite and Daniel M Roy. Computing nonvacuous generalization bounds for
289 deep (stochastic) neural networks with many more parameters than training data. *arXiv preprint*
290 *arXiv:1703.11008*, 2017.
- 291 Daniel Jakubovitz etc. Generalization error in deep learning, 2018.
- 292 Alex Graves. Agenerating sequences with recurrent neural networks, 2013. arXiv:1308.0850v5.
- 293 Moritz Hardt, Benjamin Recht, and Yoram Singer. Train faster, generalize better: Stability of
294 stochastic gradient descent. *arXiv preprint arXiv:1509.01240*, 2015.
- 295 Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the generaliza-
296 tion gap in large batch training of neural networks. *Advances in Neural Information Processing*
297 *Systems*, pp. 1731–1741, 2017.
- 298 Kenji Kawaguchi, Leslie Pack Kaelbling, and Yoshua Bengio. Generalization in deep learning.
299 *arXiv preprint arXiv:1710.05468*, 2017.
- 300 Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Pe-
301 ter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv*
302 *preprint arXiv:1609.04836*, 2016.
- 303 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
304 *arXiv:1412.6980*, 2014.
- 305 Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring general-
306 ization in deep learning. In *Advances in Neural Information Processing Systems*, pp. 5947–5956,
307 2017.
- 308 Novak R, Bahri Y, Abolafia D A, Pennington J, and Sohl-Dickstein J. Sensitivity and generalization
309 in neural networks: An empirical study. *arXiv:1802.08760*, 2018.
- 310 Tom Rainforth, Adam R Kosiorek, Tuan Anh Le, Chris J Maddison, Maximilian Igl, Frank Wood,
311 and Yee Whye Teh. Tighter variational bounds are not necessarily better. *arXiv preprint*
312 *arXiv:1802.04537*, 2018.
- 313 Arora S., Ge R., Neyshabur B., and Y. Zhang. Stronger generalization bounds for deep nets via a
314 compression approach, 2018. arXiv:1802.05296.
- 315 Vladimir N Vapnik and A Ja Chervonenkis. The necessary and sufficient conditions for consistency
316 of the method of empirical risk. *Pattern Recognition and Image Analysis*, 1(3):284–305, 1991.
- 317 Matthew D. Zeiler. Adadelta: An adaptive learning rate method, 2012. arXiv:1212.5701.
- 318 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
319 deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.

320 A APPENDIX A

321 Model Structure in Section 2.4. As shown in Table 2, all models in the experiment consist of 2
 322 Conv-Relu-MaxPooling blocks and 2 fully-connected layers, but they are different at the number
 323 of channels. We choose the number of channel p as $p \in \{6, 8, 10, 12, 14, 16, 18, 20\}$.

Table 2: Model structure in Section 2.4. p is the number of channel in experiment setting and $q = \text{int}(2.5 * p)$

Layer	input channel number	output channel number
conv + relu + maxpooling	1	p
conv + relu + maxpooling	p	q
flatten	-	-
fc + relu	$16 * q$	$10 * q$
fc + relu	$10 * q$	10
softmax	-	-

324 B APPENDIX B

The step-by-step derivation of equation 26

$$\Delta[k_l]_{c_1 c_2} = \sum_{q=1}^{Q_l} \frac{\partial[k_l]_{c_1 c_2}}{\partial\phi_q} \Delta\phi_q + \text{other terms} \quad (27)$$

$$= \sum_{q=1}^{Q_l} \frac{\partial(\frac{1}{n} \sum_{i=1}^n \frac{\partial L_i}{\partial[\omega_l]_{c_1 c_2}})}{\partial\phi_q} (-\lambda \frac{1}{n} \sum_{i=1}^n \frac{\partial L_i}{\partial\phi_q}) + \text{other terms} \quad (28)$$

$$= \sum_{q=1}^{Q_l} \frac{\partial(\frac{1}{n} \sum_{i=1}^n \frac{\partial L_i}{\partial[o_l]_{ic_2}} \frac{\partial[o_l]_{ic_2}}{\partial[\omega_l]_{c_1 c_2}})}{\partial\phi_q} (-\lambda \sum_{i=1}^n \sum_{c'_1 c'_2} \frac{\partial L_i}{\partial[o_l]_{ic'_2}} \frac{\partial[o_l]_{ic'_2}}{\partial[a_l]_{ic'_1}} \frac{\partial[a_l]_{ic'_1}}{\partial\phi_q}) + \text{other terms} \quad (29)$$

$$= -\frac{\lambda}{n^2} \sum_{q=1}^{Q_l} \frac{\partial(\sum_{i=1}^n \frac{\partial L_i}{\partial[o_l]_{ic_2}} [a_l]_{ic_1})}{\partial\phi_q} (\sum_{i=1}^n \sum_{c'_1 c'_2} \frac{\partial L_i}{\partial[o_l]_{ic'_2}} [\omega_l]_{c'_1 c'_2} \frac{\partial[a_l]_{ic'_1}}{\partial\phi_q}) + \text{other terms} \quad (30)$$

$$= -\frac{\lambda}{n^2} \sum_{q=1}^{Q_l} \sum_{i=1}^n (\frac{\partial L_i}{\partial[o_l]_{ic_2}} \frac{\partial[a_l]_{ic_1}}{\partial\phi_q} + \frac{\partial^2 L_i}{\partial[o_l]_{ic_2} \partial\phi_q} [a_l]_{ic_1}) (\sum_{c'_1 c'_2} [\omega_l]_{c'_1 c'_2} \sum_{i=1}^n \frac{\partial L_i}{\partial[o_l]_{ic'_2}} \frac{\partial[a_l]_{ic'_1}}{\partial\phi_q}) + \text{other terms} \quad (31)$$

325 In the above derivation, we use relation $\frac{\partial[o_l]_{ic'_2}}{\partial[a_l]_{ic'_1}} = [\omega_l]_{c'_1 c'_2}$ and $\frac{\partial[o_l]_{ic_2}}{\partial[\omega_l]_{c_1 c_2}} = [a_l]_{ic_1}$ which can both be
 326 derived from equation 24. Consider the first term of equation 31 and when $c'_1 = c_1, c'_2 = c_2$ in the
 327 c'_1, c'_2 summation, we have

$$\Delta[k_l]_{c_1 c_2} = -\frac{\lambda}{n^2} \sum_{q=1}^{Q_l} [\omega_l]_{c_1 c_2} (\sum_{i=1}^n \frac{\partial L_i}{\partial[o_l]_{ic_2}} \frac{\partial[a_l]_{ic_1}}{\partial\phi_q})^2 + \text{other terms} \quad (32)$$

328 Note that the term related to $\frac{\partial^2 L_i}{\partial[o_l]_{ic_2} \partial\phi_q} [a_l]_{ic_1}$ and the terms when $c'_1 \neq c_1$ or $c'_2 \neq c_2$ in equation
 329 31 is added into *other terms* of equation 32.