UNDERSTANDING NONLINEAR IMPLICIT BIAS VIA RE GION COUNTS IN INPUT SPACE

Anonymous authors

Paper under double-blind review

ABSTRACT

One explanation for the strong generalization ability of neural networks is implicit bias. Yet, the definition and mechanism of implicit bias in non-linear contexts remains little understood. In this work, we propose to characterize implicit bias by the count of connected regions in the input space with the same predicted label. Compared with parameter-dependent metrics (e.g., norm or normalized margin), region count can be better adapted to nonlinear, overparameterized models, because it is determined by the function mapping and is invariant to reparametrization. Empirically, we found that small region counts align with geometrically simple decision boundaries and correlate well with good generalization performance. We also observe that good hyper-parameter choices such as larger learning rates and smaller batch sizes can induce small region counts. We further establish the theoretical connections and explain how larger learning rate can induce small region counts in neural networks.

023 024 025

026

027

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

One mystery in deep neural networks lies in their ability to generalize, despite having significantly more learnable parameters than the number of training examples (Zhang et al., 2017a). The choice of network architectures, including factors such as nonlinearity, depth, and width, along with training procedures like initialization, optimization algorithms, and loss functions, can result in vastly diverse generalization performances (Sutskever et al., 2013; Smith et al., 2017; Wilson et al., 2017; Li et al., 2019). The varied generalization abilities exhibited by neural networks are often explained by many researchers through the theory of *implicit bias* (Brutzkus et al., 2017; Soudry et al., 2018). Implicit bias refers to inherent tendencies of how the network learns and generalizes from the training data, even without explicit regularizations or constraints.

The implicit bias of linear neural networks has been extensively studied. One of the classical setting is linear classification with logistic loss. Brutzkus et al. (2017); Soudry et al. (2018); Arora et al. (2019) show that the parameter converges to the direction that maximizes the L_2 margin. For regression problems, it is proved that gradient descent or stochastic gradient descent converges to a parameter that is closest to the initialization in terms of L_2 norm (Gunasekar et al., 2018a). The results from the linear regression model can be extended to deep linear neural networks by generalizing the definition of min-norm and max-margin solutions (Ji & Telgarsky, 2018a; Vaskevicius et al., 2019; Woodworth et al., 2020).

Compared to linear models, defining implicit biases in non-linear networks poses significant challenges. One line of work studies homogeneous networks and demonstrates that gradient flow solutions converge to a KKT point of the max-margin problem (Lyu & Li, 2019; Ji & Telgarsky, 2020; Wang et al., 2021; Jacot et al., 2022). Further research extends this analysis, showing that gradient flow converges to a max-margin solution under various norms (Ongie et al., 2019; Chizat & Bach, 2020).
Other studies focus on describing the implicit bias of neural networks using sharpness, such as (Foret et al., 2020; Montúfar et al., 2022; Andriushchenko et al., 2023).

We note that previous definitions of implicit bias in neural networks mostly focus on certain metrics
 of network parameters. Such approaches enable explicit analyses of training trajectories, but face
 new challenges when applied to nonlinear networks: reparametrization of the network may preserve



Figure 1: A schematic illustration of main results in this paper. Left: The region counts in 2-dimension input space. Each distinct region represents an area where the neural network makes the same prediction for all points within that region. Middle: A strong correlation between region counts and the generalization gap. Right: Larger learning rate or smaller batch size induces smaller region counts.

the function mapping but gives completely different parameters, and consequently, different implicit
 biases. We will discuss this point in detail in Section 3.

072 Motivated by the above studies, we instead focus on leveraging the decision boundaries in the input 073 space to characterize implicit bias. Our research identifies a metric called region count, which is 074 defined by the average number of regions in the predictor's decision boundary (See Figure 1). We 075 select a low-dimensional space, use the neural network to predict the labels of all points within this 076 space, and define the number of connected components with the same label as the region count in 077 that subspace. This definition differs from previous studies on the number of linear regions (Hanin & Rolnick, 2019a;b; Safran et al., 2022) which is defined as the set of inputs that correspond to the 078 079 same activation pattern in the network. We find that this region count has a strong correlation with the generalization gap, defined as the difference in percentage between the training and test errors. The experiments in Figure 1 suggests that models with fewer region counts tend to generalize better. 081

Furthermore, we show that neural networks trained with large learning rate or small batch size, which are typically deemed as beneficial for generalization, are biased towards solutions that have small region counts. Therefore, region count empirically serves as an accurate generalization metric as well as an implicit bias indicator. We also provide theoretical analyses to explain this phenomenon. We prove that for two-layer ReLU neural networks, gradient descent with large learning rate induces a small region count, which accords well with our empirical findings.

- ⁸⁸ The main contributions of this paper are listed as follows:
 - 1. We introduce a novel measure of implicit bias via the region count in the input space. Through extensive experiments, we verify a strong correlation between region count and the generalization gap. This correlation remains robust across different learning methods, datasets, training parameters, and counting methods.
 - 2. We assess the factors that induce small region count, discovering that training with larger learning rates and smaller batch sizes typically results in fewer regions. This provides a possible cause for the implicit bias in neural networks.
 - 3. We conduct theoretical analyses on region counts, and show that for two-layer ReLU neural networks, gradient descent with large learning rate induces a small region count.
- 099 100

090

091

092

093 094

095

096

098

066

067

068

069

- 2 RELATED WORKS
- 101 102

Implicit Bias of Linear Neural Network The implicit bias in linear neural networks are thoroughly
 investigated in recent works. For linear logistic regression on linearly separable data, full-batch
 gradient descent converges in the direction of the maximum margin solution (Soudry et al., 2018).
 This foundational work has various follow-ups, including extensions to non-linearly-separable data (Ji
 & Telgarsky, 2018b; 2019), stochastic gradient descent (Nacson et al., 2019), and other loss functions and optimizers (Gunasekar et al., 2018a).

These findings in linear logistic regression are generalized to deep linear networks. For fullyconnected neural networks with linear separable data, Ji & Telgarsky (2018a) show that the direction of weight also converges to L_2 max-margin solution. For linear diagonal networks, the gradient flow maximizes the margin with respect to a specific quasi-norm that is related to the depth of network (Gunasekar et al., 2018b; Woodworth et al., 2020; Pesme et al., 2021), leading to a bias towards sparse linear predictors as the depth goes to infinity. This sparsity bias also exists in linear convolutional networks (Gunasekar et al., 2018b; Yun et al., 2020).

115

116 Implicit Bias of Non-linear Neural Network The non-linearity of modern non-linear neural 117 networks pose challenges to studying its implicit bias. Initial works in this area (Lyu & Li, 2019; Ji & 118 Telgarsky, 2020) focus on homogeneous networks. These studies show that with exponentially-tailed 119 classification losses, both gradient flow and gradient descent converge directionally to a KKT point 120 in a maximum-margin problem. Further studies, for instance in (Wang et al., 2021), consider a more general setup that includes different optimizers and prove that both Adam and RMSProp are capable 121 of maximizing the margin in neural networks while Adagrad is not. Ongie et al. (2019); Chizat 122 & Bach (2020) showcased a bias towards maximizing the margin in a variation norm for infinite-123 width two-layer homogeneous networks. Lyu et al. (2021); Jacot et al. (2022) identified margin 124 maximization in two-layer Leaky-ReLU networks trained with linearly separable and symmetric 125 data. More recent investigations into non-linear neural networks, such as (Jacot, 2022), focus on 126 the homogeneity of the non-linear layer, demonstrating an implicit bias characterized by a novel 127 non-linear rank.

128 129

Region Counts of Neural Network Many previous works focus on calculating the linear regions 130 of neural networks (Hanin & Rolnick, 2019a;b). A linear region is a set of inputs that share the same 131 activation pattern in the network. Safran et al. (2022) proves that for a two-layer ReLU network 132 with width r, gradient flow will converge directionally to a network characterized by no more than 133 O(r) linear regions. Serra et al. (2018) and Cai et al. (2023) explain that maximizing the number 134 of linear regions can lead to better-performing networks, and they explore how network structures 135 can be designed to achieve more linear regions. The number of linear regions is independent of 136 the network's label output, focusing more on its representational capacity rather than generalization 137 ability. In contrast, we define decision regions as connected areas in the input space that correspond 138 to the same label, and our work explores the relationship between the number of regions and the generalization gap. 139

140 The paper (Nguyen et al., 2018) is the most relevant to our work, as it similarly defines decision 141 regions corresponding to label predictions of neural networks. They define decision regions in 142 the entire input space, whereas our definition and counting method focus on a subspace. Their 143 conclusion states that, under certain conditions, each label has only one connected decision region in 144 the entire space. However, two points being connected in the entire space does not imply they are connected in a subspace. For instance, two points may be connected in a three-dimensional space, 145 but a two-dimensional cross-section may not provide a connecting path. Thus, we observe multiple 146 regions in 1D or 2D subspaces and find a strong correlation with generalization ability. 147

148 149

3 MOTIVATION

150 151

Norm-based and margin-based characterizations belong to the most popular measures of implicit
 bias. Various definitions for norm and margin exist. For simplicity, we consider the following two
 definitions.

Example 1 (Norm and Margin). Let $W = \{W_1, \dots, W_l\}$ denote the post-training weight parameters of an *l*-layer neural network $f_W(x) = W_l \sigma(W_{l-1} \cdots W_2 \sigma(W_1 x))$, with $\sigma(\cdot)$ as the ReLU activation function. Denote the weight initialization as $W^0 = \{W_1^0, \dots, W_l^0\}$. Consider the Frobenious norm between network weights and initialization:

- 159
- 160

$$d(W) = \sqrt{\sum_{i=1}^{l} \|W_i - W_i^0\|_F^2},$$



Figure 2: Norm-based and margin-based measures may not be predictive of generalization gaps. We train ResNet18 on the CIFAR-10 dataset using various hyperparameters. These implicit bias measures can be ineffective for general non-linear neural networks.

 $\gamma(W) = \mathbb{E}_{(x,y) \in D_{train}} \left[f(x)_y - \max_{i \neq y} f(x)_i \right].$

and the output-space margin

172

173

174

175

178 179

 $d(\cdot)$ and $\gamma(\cdot)$ are commonly used indicator for implicit bias of linear models (Soudry et al., 2018; 181 Ji & Telgarsky, 2018b). However, both of them are not invariant to network reparameterization. We can construct a different set of network weight parameters by scaling the parameters 182 as $\hat{W} = \{2W_1, \frac{1}{2}W_2, \cdots, W_l\}$, such that $f_W = f_{\hat{W}}$ but $d(W) \neq d(\hat{W})$ in general. Similarly, 183 we can scale the last layer weights and get $\tilde{W} = \{W_1, \dots, 2W_l\}$, such that $\gamma(\tilde{W}) \neq \gamma(W)$, but $\operatorname{argmax}_i f_W(x) = \operatorname{argmax}_i f_{\tilde{W}}(x)$. This reparameterization trick also works for more complicated 185 norm-based and margin-based generalization metric in (Jiang et al., 2019), or the sharpness-based metrics (Andriushchenko et al., 2023). 187

188 We numerically investigate whether they are effective measures, by training a ResNet18 on Cifar10 189 dataset, using different hyperparameters as in Table 1. The results are presented in Figure 2, which 190 indicates that these measures have a low correlation with the generalization gap in the deep learning 191 regime. One could choose other definitions of norms to achieve stronger correlations, but such choices are often problem-specific and require domain expertise, as discussed in (Jiang et al., 2019). 192

193 The definition of margin may also be improved to the input-space margin, *i.e.*, the ℓ_2 distance of 194 input data x to decision boundary defined by the classifier, which is able to characterize the quality 195 and robustness of the classifier. This metric is invariant to reparameterization and therefore more 196 intrinsic to the underlying classifier. However, due to the highly nonconvex loss landscape, the 197 input-space margin is NP-complete to compute and even hard to approximate (Katz et al., 2017; Weng et al., 2018). Therefore, quantitatively analyzing the decision boundary of a neural network and characterizing its implicit bias remains a challenge. 199

200 Our motivation can be summarized by a simple idea: although the margin in the input space is hard 201 to compute, we can quantify the regions split by the decision boundary. This measure is invariant 202 to model reparametrization and can also capture the complexity of the decision boundary. This 203 motivates us to consider the region counts as an implicit bias metric.

204 205

206

4 PRELIMINARY

207 Although region count is a natural measure for the complexity of a predictor, and it depends only on 208 the decision function rather than the model parameterization, its formal definition and computability 209 remains unclear. In this section, we first provide the definition and low-dimensional approximation of 210 region counts. We then empirically verify that region counts correlate with generalization gap.

211

212 4.1 DEFINITION OF REGION COUNTS

213

Let d denote the training data dimension and $f : \mathbb{R}^d \to \{1, 2, \dots, N\}$ denote a neural network for 214 a classification task with N classes. For a subset $U \subset \mathbb{R}^d$, we can define the connectedness of its 215 element as follows:

216 **Definition 1** (Connectedness). We say the data points $x_1, x_2 \in U$ are (path) connected with respect 217 to a neural network f if they satisfy: 218

• $f(x_1) = f(x_2) = c$,

219 220

221

222 223

224

225

226 227

228

229 230

231

232

233 234

235

236 237

238 239

243

244

245 246 247

263

264

268

269

• There exist a continuous mapping $\gamma: [0,1] \to U, \gamma(0) = x_1, \gamma(1) = x_2$, and for any $t \in [0, 1], f(\gamma(t)) = c.$

Then we define the connected region in this subset:

Definition 2 (Maximally Connected Region). We say $V \subset U$ is a maximally connected region in $U \in \mathbb{R}^d$ with respect to a neural network f if it satisfies the following property:

- For any $x, y \in V$, they are connected.
- For any $x \in V$, $y \in U \setminus V$, they are not connected.

Finally, we formally define the region count as follows:

Definition 3 (Region Count). For a subset $U \subseteq \mathbb{R}^d$, we define its region count R_U as the number of maximally connected regions in U with respect to a neural network f:

$$R_U = \operatorname{card}\{V \subset U | V \text{ is a maximally connected region}\},\$$

where card is the cardinality of a set.

4.2 ESTIMATING REGION COUNTS

Calculating the region count in the original high-dimensional input space can be computationally 240 intractable. Therefore, we propose a computationally efficient surrogate by calculating the region 241 counts on low dimensional subspace spanned by training data points. 242

Definition 4 (Region Count in *d*-Dimensional Subspace). We randomly sample d + 1 datapoints in the training set D_{train} to generate a convex region in \mathbb{R}^d subspace. The d-dimensional region count R_d is defined as the expectation of number of maximally connected regions:

$$R_d = E_{x_1, x_2, \dots, x_{d+1} \sim D_{train}} [R_{\text{Conv}\{x_1, x_2, \dots, x_{d+1}\}}],$$

where $x_1, x_2, \ldots, x_{d+1}$ are sampled from the training dataset, and Conv $\{x_1, x_2, \ldots, x_{d+1}\}$ is the 248 convex hull formed by these d + 1 points. 249

250 This paper primarily focuses on low dimension 251 spaces, which is illustrated as below. In practice, 252 we randomly sample training data points for mul-253 tiple times and take the average region counts. In 254 Section 7, we show that the choice of subspace 255 dimension d does not significantly affect the results. The details on how to count the regions and 256 generate the polytopes are provided in Appendix B. 257

258 Example 2 (Region counts in 1D and 2D subspace). 259 For region count in 1-dimensional subspace, we 260 randomly sample two data points, denoted as x_1 and x_2 , from the training set, and calculate the 261 region count on the line segment connecting them: 262

$$\{\alpha x_1 + (1 - \alpha)x_2, 0 \le \alpha \le 1\}.$$



For the 2-dimensional case, we randomly sample 265 three data points, x_1 , x_2 , and x_3 , from the training 266

set, and calculate the region count in the convex hull spanned by them: 267

 $\{\alpha x_1 + \beta x_2 + (1 - \alpha - \beta) x_3, \alpha \ge 0, \beta \ge 0, \alpha + \beta \le 1\}.$

We provide an illustration in Figure 3.



Figure 3: Illustrations of region counts in 1D and 2D subspace. We use different colors to represent different outputs of the neural network.

²⁷⁰ 5 REGION COUNTS CORRELATE WITH GENERALIZATION GAPS

In this section, we present our major empirical findings, which reveal a strong correlation between
 region counts and the generalization error of neural networks.

We conduct image classification experiments on the CIFAR-10 dataset, using different architectures, including ResNet18 (He et al., 2016), EfficientNetB0 (Tan & Le, 2019), and SeNet18 (Hu et al., 2018).
Results on other architectures are deferred to ablation studies. We vary the hyperparameters for training, such as learning rate, batch size and weight decay coefficient, whose numbers are reported in Table 1. The region count is calculated using randomly generated 1D hyperplanes, as described in Example 2. We run each experiment 100 times and report the average number.

We plot the region count and generalization gap of different setups in Figure 4, and calculate the correlation between them. For each network architecture, we observe a strong correlation as high as 0.98. The overall correlation for all the three networks still reaches 0.93. This reveals a remarkably high correlation between region counts and generalization gap.

We also explore whether such a strong correlation exists for traditional machine learning algorithms. We conduct experiments with decision trees and random forests on the same classification tasks, with hyperparameters specified in Table 1. We observe a similar linear trend between region count and generalization gap, with a correlation of 0.96 in decision trees and 0.98 in random forests. The overall correlation coefficient is 0.90. Therefore, region count serves as a good indicator for generalization performance across various setups.

Table 1: The hyperparameters for experiments. Left: We vary the learning rate, batch size, and weight decay
 for training a neural network, to modulate the model's generalization ability. Right: We adjust the training
 parameters for traditional machine learning models, such as decision tree and random forest.

Hyperparameters	Value	Hyperparameters	Value
Learning rate	0.1, 0.01, 0.001	Depth	$3, 4 \cdots, 17$
Batch size	256, 512, 1024	Criterions	gini, entropy
Weight decay	$10^{-5}, 10^{-6}, 10^{-7}$	Splitter	best, random



Figure 4: Strong correlation between region counts and generalization gap. Left: We conduct experiments using three neural networks on the CIFAR-10 dataset, with various hyperparameters. There is a strong correlation between region counts and the generalization gap, with a correlation coefficient of 0.98 for each network and 0.93 across all networks. Right: We conduct experiments using Decision Tree and Random Forest on the CIFAR-10 dataset. The result also reveals a strong correlation between region counts and the generalization gap.

317 318 319

303

305

306

307

308

309

310

311

312

313

314

315

316

6 REGION COUNTS QUANTIFY IMPLICIT BIAS

320 321

In this section, we further investigate the implicit bias of neural networks via region counts. We show both empirically and theoretically that neural networks trained with appropriate hyperparameters tend to have smaller region counts, thus achieving better generalization performance.



Figure 5: Large learning rate and small batch size reduce region counts. We train three networks on the CIFAR-10 dataset, varying the batch sizes and learning rates. Our findings reveal that a smaller batch size or a higher learning rate results in smaller region counts, allowing the network to learn a simpler decision boundary and generalize better.

6.1 THE BIAS FROM TRAINING HYPERPARAMETERS

Training neural networks requires careful selection of many hyperparameters, such as learning rates, batch sizes, optimizers, epochs and so on. Here, we primarily focus on learning rate and batch size, and study their impact on the region count.

Learning Rates. We provide the relationship between the learning rate and the region count in
 Figure 5. Our findings indicate that a larger learning rate tends to simplify the decision boundary and
 results in a smaller region count in the hyperplane. This accords well with real practices, where large
 learning rates of 0.1 or 0.01 are often favored for better generalization.

Batch Sizes. Similarly, the training batch size can impact the number of regions. As shown in
Figure 5, smaller batch sizes lead to a model with fewer regions. This result reveals the advantage of
small-batch training, which leads to better generalization accuracy.

Previous studies (Keskar et al., 2016; Jastrzkebski et al., 2017; Hoffer et al., 2017; Novak et al., 2018)
find that certain hyperparameters, such as a large learning rate and a small batch size, can improve
the generalization of the neural network. Our observations provide a possible explanation: these
good hyperparameter choices lead to a reduced region count. Such simplicity bias can decrease the
generalization gap of neural networks.

360 361

375

337

338

339 340 341

342 343

344

345

346

351

6.2 THEORETICAL EXPLANATIONS

362
 363 Next, we present a theoretical analysis to explain why some hyperparameter choices, such as large
 364 learning rate, can lead to small region counts.

Consider a two layer ReLU neural network $f_W(x) = \sum_{i=1}^p a_i \sigma(w_i^\top x)$. The second layer weights a_i are initialized uniformly from $\{1, -1\}$ and fixed throughout training. Let $\mathcal{D} = \{(x_i, y_i)\}_{1 \le i \le N}$ denote the training set. Consider training f_W on \mathcal{D} using gradient descent (GD) with learning rate η . We choose the quadratic loss $l(W, x, y) = \frac{1}{2}(y - f_W(x))^2$ and denote $L(W) = \frac{1}{N} \sum_{i=1}^N l(W, x_i, y_i)$. Denote the GD trajectory as $\{W_i\}_{i\ge 0}$. For two input data x_a, x_b , Let $R(x_a, x_b, W)$ denote the region count on the line segment connecting them, and $N(x_a, W)$ denote the number of activated neurons with input x_a , i.e., the number of i such that $w_i^\top x_a > 0$.

We make the following assumption on the data distribution.

Assumption 1. The training dataset $\mathcal{D} = \{(x_i, y_i)\}_{1 \le i \le N}$ satisfies the following two properties:

- 1. $||x_i|| \ge r$ for all $1 \le i \le N$,
- 376 377 2. With probability one, any $W \in \{W_i\}_{i \ge 0}$ satisfies $w_i^\top x_j \neq 0$ for all $1 \le i \le q, 1 \le j \le N$, where the randomness comes from weight initialization.

The validity of Assumption 1 comes from the fact that the bifurcation zone (Bertoin et al., 2021) of ReLU neural networks, which contains its non-differentiable points, has Lebesgue measure zero (Bolte & Pauwels, 2020; 2021; Bianchi et al., 2022). Therefore, if the distribution of weights are absolutely continuous with respect to the Lebesgue measure, the bifurcation zone can be avoided with probability one. We conjecture that it can be proved rigorously, but leave it as an assumption since the proof diverges from the main content in this paper.

The next assumption characterizes the sharpness along the training trajectory. This is actually from the well-known edge of stability phenomenon (Cohen et al., 2020; Damian et al., 2022; Arora et al., 2022; Ahn et al., 2024), which states that the sharpness of neural networks, characterized by the ℓ_2 norm of the Hessian matrix, hovers around $\frac{2}{\eta}$.

Assumption 2 (Edge of Stability). There exist a $T \in \mathbb{N}$, such that for $t \ge T$, with we have

$$\lambda_{\max}(\nabla_W^2 L(W_t)) = \Theta\left(\frac{1}{\eta}\right)$$

where λ_{\max} denotes the maximum eigenvalue of a matrix.

We are now ready to present the main theorem, which establishes a relationship between region count and learning rate.

Theorem 1. Under Assumption 1 and 2, we have that for neural net weights W_t at training step $t \ge T$ with probability one, the average region count $R(X, X', W_t)$ for random training data point X, X' can be bounded as:

$$\mathbb{E}_{X,X'}[R(X,X',W_t)] = \frac{1}{N^2} \sum_{i=1}^{N} \sum_{j=1}^{N} R(x_i,x_j,W_t) \le O\left(\frac{N}{r^2\eta}\right).$$

401 402 403

389 390 391

394

395

396

397

398

399 400

The theorem demonstrates that with a larger learning rate, gradient descent has the implicit bias to yield solutions with smaller region counts. This aligns well with the previous observations.

We defer the proof of this theorem to Appendix D, and sketch the proof as follows. The proof begins with bounding the region count using the activation pattern of ReLU neurons, as stated in the following lemma.

Lemma 2. The region counts between a pair of data points is upper-bounded by the number of active neurons. For two inputs x_a, x_b , we have $R(x_a, x_b, W) \le N(x_a, W) + N(x_b, W) + 2$.

⁴¹² Then we prove that the activation pattern gives a bound on the smoothness of the training loss.

413 **Lemma 3.** The sharpness of a neural network is lower-bounded by the number of active neurons: 414 $\lambda_{\max} \left(\nabla_W^2 L(W) \right) \geq \frac{r^2}{N^2} \sum_{i=1}^N N(x_i, W).$

416 Note that this lemma brings in an additional N in the denominator, which leads to a N-dependent
417 bound in Theorem 1. We conjecture that the N-dependency can be optimized under further structural
418 assumptions on the data distribution, and leave it for further investigations. Equipped with these two
419 lemmas, Theorem 1 is a consequence of the sharpness condition in Assumption 2.

420 421

7 ABLATION STUDIES

422

427

This section presents an ablation study to validate the robustness and consistency of our findings. We
systematically vary key aspects of our experimental setup, including the network architecture, dataset,
optimizer, and the method of computing the plane, and evaluate their impact on our main results of
the correlation between region count and the generalization gap.

More Architectures, Datasets and Hyperplane Dimensions. We first examine the influence of neural network architectures and datasets on our results. We provide additional results on various neural network architectures such as ResNet34 (He et al., 2016), VGG19 (Simonyan & Zisserman, 2014), MobileNetV2 (Sandler et al., 2018), ShuffleNetV2 (Ma et al., 2018), RegNet200MF (Radosavovic et al., 2020), and SimpleDLA (Yu et al., 2018). We also use various datasets such as

CIFAR-100 (Krizhevsky et al., 2009) and ImageNet (Deng et al., 2009). Region counts and general ization gaps are evaluated across various learning rates, batch sizes, and weight decay parameters as
 listed in Table 1.

We also explore the effects of different methods for generating the hyperplane in the input space.
In our previous experiments, we generate the 1-dimensional plane using random pairs of samples from the training set and calculate the region count on them. Here we explore region counts in higher dimensional planes, that are spanned by 2 to 5 data randomly-selected points from the training set, using the CIFAR-10 dataset.

The experiment results of the correlation are presented in Table 2. We also provide correlation plots for each network in Appendix A.2. We observe that the strong correlation between region count and the generalization gap remains consistent in various setups. The consistency indicates that our findings reveal a fundamental characteristic of non-linear neural networks.

We also provide evaluations by varying the optimizer and hyperplane generation algorithms. The results are deferred to Appendix C.

447

448 449

Table 2: Experimental consistency across networks, datasets, and counting methods. We conduct experiments on various types of networks across multiple datasets. We also alter the method of calculating the region counts. The results of the correlation indicate that our findings are consistent across different setups.

Natwork	Dataset			Counting Dimension			
INCLWOIK	CIFAR-10	CIFAR-100	ImageNet	2	3	4	5
ResNet18	0.98	0.96	0.91	0.96	0.97	0.97	0.96
ResNet34	0.98	0.98	0.82	0.98	0.98	0.98	0.99
VGG19	0.94	0.85	0.78	0.88	0.86	0.84	0.86
MobileNet	0.95	0.95	0.92	0.99	0.99	0.99	0.99
SENet18	0.98	0.85	0.80	0.97	0.97	0.97	0.93
ShuffleNetV2	0.95	0.92	0.92	0.94	0.95	0.95	0.93
EfficientNetB0	0.98	0.84	0.93	0.99	0.99	0.99	0.98
RegNetX_200MF	0.98	0.87	0.97	0.98	0.99	0.99	0.98
SimpleDLA	0.98	0.94	0.84	0.99	0.99	0.98	0.99

462 463

464 465

Data Augmentations. Mixup (Zhang et al., 2017b) is a data augmentation technique that creates training samples by linearly interpolating pairs of input data and their corresponding labels. We train a ResNet-18 model using mixup, with other hyperparameters in Table 1. The plot in Figure 6 illustrates that Mixup induces smoother decision boundaries with smaller region count and has a better generalization performance.

Random crop and random horizontal flip is another way to enhance the diversity of the training dataset. We apply random crop of size 32×32 with padding 4 and random horizontal flip with a probability of 0.5 as data augmentations. As depicted in Figure 7, we observe that compared with mixup, random crop and random flip result in a more evident vertical shift in the performance curve.

Both Figure 6 and Figure 7 show that employing these techniques does not alter the correlation between region counts and generalization gaps.

477

Evolution of Region Counts during Training. Next we study how the region count, generalization
gap and their correlation evolve during training. Following the setup in Section 5, we train a ResNet18
model on CIFAR-10 dataset, and report the region count and generalization gap during the training
process. The statistics are averaged over different hyperparameter choices as in Table 1.

The results are provided in Table 3. We recorded the average values of region count and generalization gap for these data points in the second and third columns to show their changes during the training process. We observe that the correlation is very low at initialization, but steadily increases during training. This suggests that the metric of region count is not a property of the neural network initialization, but rather inherently involved with the neural network training algorithm.



Figure 6: The impact of mixup. This figure shows that mixup improves the model's generalization ability and reduces the number of regions in the hyperplane.



Figure 7: The impact of random crop and random flip. Unlike mixup, data augmentation results in a vertical shift in the performance curve, accompanied by a decrease in the number of regions and a more significant enhancement in test accuracy.

Table 3: The evolution of region count, generalization gap and their correlation. The correlation is very low at initialization, but steadily increases during training.

Training Epoch	Region Counts	Generalization Gap	Correlation
0	1.13	N/A	N/A
20	3.14	11.2	-0.53
40	3.02	7.3	-0.29
60	3.10	18.2	0.35
80	3.22	26.7	0.77
100	3.25	30.4	0.98
130	3.28	31.2	0.97
160	3.27	31.7	0.98
200	3.26	31.6	0.98

8 CONCLUSIONS AND FUTURE DIRECTIONS

This paper introduces a novel approach to characterizing the implicit bias of neural networks. We study
the region counts in the input space and identify its strong correlation with generalization gap in nonlinear neural networks. These findings are consistent across various network architectures, datasets,
optimizers. Our analysis offers a new perspective to quantify and understand the generalization
property and implicit bias of neural networks.

Our paper suggests several promising directions for future research. Firstly, our analyses of why large
 learning rate induces small region counts mainly focus on a simplified setup. The analyses for more
 general settings remain open. Secondly, extending the definition of region count to non-classification
 tasks, such as natural language generation, would be a worthwhile direction. Lastly, region count
 can be leveraged to design new architectures or regularization, that can potentially improve the
 generalization performance of neural networks.

540 REFERENCES

548

576 577

578

579

- Kwangjun Ahn, Sébastien Bubeck, Sinho Chewi, Yin Tat Lee, Felipe Suarez, and Yi Zhang. Learning
 threshold neurons via edge of stability. *Advances in Neural Information Processing Systems*, 36, 2024.
- Maksym Andriushchenko, Francesco Croce, Maximilian Müller, Matthias Hein, and Nicolas Flammarion. A modern look at the relationship between sharpness and generalization. *arXiv preprint* arXiv:2302.07011, 2023.
- Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- Sanjeev Arora, Zhiyuan Li, and Abhishek Panigrahi. Understanding gradient descent on the edge of stability in deep learning. In *International Conference on Machine Learning*, pp. 948–1024.
 PMLR, 2022.
- David Bertoin, Jérôme Bolte, Sébastien Gerchinovitz, and Edouard Pauwels. Numerical influence of
 relu'(0) on backpropagation. *Advances in Neural Information Processing Systems*, 34:468–479, 2021.
- Pascal Bianchi, Walid Hachem, and Sholom Schechtman. Convergence of constant step stochastic
 gradient descent for non-smooth non-convex functions. *Set-Valued and Variational Analysis*, 30 (3):1117–1147, 2022.
- Jérôme Bolte and Edouard Pauwels. A mathematical model for automatic differentiation in machine
 learning. Advances in Neural Information Processing Systems, 33:10809–10819, 2020.
- Jérôme Bolte and Edouard Pauwels. Conservative set valued fields, automatic differentiation, stochastic gradient methods and deep learning. *Mathematical Programming*, 188:19–51, 2021.
- Alon Brutzkus, Amir Globerson, Eran Malach, and Shai Shalev-Shwartz. SGD learns over parameterized networks that provably generalize on linearly separable data. *arXiv preprint arXiv:1710.10174*, 2017.
- Junyang Cai, Khai-Nguyen Nguyen, Nishant Shrestha, Aidan Good, Ruisen Tu, Xin Yu, Shandian Zhe, and Thiago Serra. Getting away with more network pruning: From sparsity to geometry and linear regions. In *International Conference on Integration of Constraint Programming, Artificial Intelligence, and Operations Research*, pp. 200–218. Springer, 2023.
- Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks
 trained with the logistic loss. In *Conference on Learning Theory*, pp. 1305–1338. PMLR, 2020.
 - Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2020.
- Alex Damian, Eshaan Nichani, and Jason D Lee. Self-stabilization: The implicit bias of gradient descent at the edge of stability. *arXiv preprint arXiv:2209.15594*, 2022.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.
- Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. In *International Conference on Machine Learning*, pp. 1832–1841.
 PMLR, 2018a.
- 593 Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. *Advances in neural information processing systems*, 31, 2018b.

602

608

614

625

626

627

594	Boris Hanin and David Rolnick. Complexity of linear regions in deep networks. In International
595	Conference on Machine Learning, pp. 2596–2604. PMLR, 2019a.
596	

- Boris Hanin and David Rolnick. Deep relu networks have surprisingly few activation patterns.
 Advances in neural information processing systems, 32, 2019b.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the generaliza tion gap in large batch training of neural networks. *Advances in neural information processing systems*, 30, 2017.
- Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7132–7141, 2018.
- Arthur Jacot. Implicit bias of large depth networks: a notion of rank for nonlinear functions. *arXiv* preprint arXiv:2209.15055, 2022.
- Arthur Jacot, Eugene Golikov, Clément Hongler, and Franck Gabriel. Feature learning in *l*_2-regularized dnns: Attraction/repulsion and sparsity. *Advances in Neural Information Processing Systems*, 35:6763–6774, 2022.
- Stanislaw Jastrzkebski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey. Three factors influencing minima in SGD. *arXiv preprint arXiv:1711.04623*, 2017.
- Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. *arXiv preprint arXiv:1810.02032*, 2018a.
- Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression. *arXiv preprint arXiv:1803.07300*, 2018b.
- Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In
 Conference on Learning Theory, pp. 1772–1798. PMLR, 2019.
 - Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. Advances in Neural Information Processing Systems, 33:17176–17186, 2020.
- Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic
 generalization measures and where to find them. *arXiv preprint arXiv:1912.02178*, 2019.
- Guy Katz, Clark Barrett, David L Dill, Kyle Julian, and Mykel J Kochenderfer. Reluplex: An efficient smt solver for verifying deep neural networks. In *Computer Aided Verification: 29th International Conference, CAV 2017, Heidelberg, Germany, July 24-28, 2017, Proceedings, Part I 30*, pp. 97–117. Springer, 2017.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter
 Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Yuanzhi Li, Colin Wei, and Tengyu Ma. Towards explaining the regularization effect of initial large
 learning rate in training neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks.
 arXiv preprint arXiv:1906.05890, 2019.
- Kaifeng Lyu, Zhiyuan Li, Runzhe Wang, and Sanjeev Arora. Gradient descent on two-layer nets: Margin maximization and simplicity bias. *Advances in Neural Information Processing Systems*, 34:12978–12991, 2021.

658

665

648	Ningning Ma Xiangyu Zhang Hai-Tao Zheng and Jian Sun, Shufflenet v2: Practical guidelines for
649	efficient cnn architecture design. In <i>Proceedings of the European conference on computer vision</i>
650	(ECCV) np 116–131 2018
651	(2007), pp. 110-101, 2010.

- Guido Montúfar, Yue Ren, and Leon Zhang. Sharp bounds for the number of regions of maxout networks and vertices of minkowski sums. *SIAM Journal on Applied Algebra and Geometry*, 6(4): 618–649, 2022.
- Mor Shpigel Nacson, Nathan Srebro, and Daniel Soudry. Stochastic gradient descent on separable
 data: Exact convergence with a fixed learning rate. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 3051–3059. PMLR, 2019.
- Quynh Nguyen, Mahesh Chandra Mukkamala, and Matthias Hein. Neural networks should be wide
 enough to learn disconnected decision regions. In *International conference on machine learning*,
 pp. 3740–3749. PMLR, 2018.
- Roman Novak, Yasaman Bahri, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl Dickstein. Sensitivity and generalization in neural networks: an empirical study. *arXiv preprint arXiv:1802.08760*, 2018.
- Greg Ongie, Rebecca Willett, Daniel Soudry, and Nathan Srebro. A function space view of bounded norm infinite width relu nets: The multivariate case. *arXiv preprint arXiv:1910.01635*, 2019.
- Scott Pesme, Loucas Pillaud-Vivien, and Nicolas Flammarion. Implicit bias of sgd for diagonal linear networks: a provable benefit of stochasticity. *Advances in Neural Information Processing Systems*, 34:29218–29230, 2021.
- Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing
 network design spaces. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10428–10436, 2020.
- Itay Safran, Gal Vardi, and Jason D Lee. On the effective number of linear regions in shallow univariate relu networks: Convergence guarantees and implicit bias. *Advances in Neural Information Processing Systems*, 35:32667–32679, 2022.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mo bilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4510–4520, 2018.
- Thiago Serra, Christian Tjandraatmadja, and Srikumar Ramalingam. Bounding and counting linear regions of deep neural networks. In *International conference on machine learning*, pp. 4558–4566. PMLR, 2018.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image
 recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Samuel L Smith, Pieter-Jan Kindermans, Chris Ying, and Quoc V Le. Don't decay the learning rate, increase the batch size. *arXiv preprint arXiv:1711.00489*, 2017.
- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit
 bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):
 2822–2878, 2018.
- Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pp. 1139–1147.
 PMLR, 2013.
- Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks.
 In *International conference on machine learning*, pp. 6105–6114. PMLR, 2019.
- 701 Tomas Vaskevicius, Varun Kanade, and Patrick Rebeschini. Implicit regularization for optimal sparse recovery. *Advances in Neural Information Processing Systems*, 32, 2019.

702 703 704	Bohan Wang, Qi Meng, Wei Chen, and Tie-Yan Liu. The implicit bias for adaptive optimization algorithms on homogeneous neural networks. In <i>International Conference on Machine Learning</i> , pp. 10849–10858. PMLR, 2021.
705 706 707 708	Lily Weng, Huan Zhang, Hongge Chen, Zhao Song, Cho-Jui Hsieh, Luca Daniel, Duane Boning, and Inderjit Dhillon. Towards fast computation of certified robustness for relu networks. In <i>International Conference on Machine Learning</i> , pp. 5276–5285. PMLR, 2018.
709 710 711	Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. <i>Advances in neural information processing systems</i> , 30, 2017.
712 713 714 715	Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In <i>Conference on Learning Theory</i> , pp. 3635–3673. PMLR, 2020.
716 717 718	Fisher Yu, Dequan Wang, Evan Shelhamer, and Trevor Darrell. Deep layer aggregation. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 2403–2412, 2018.
719 720 721	Chulhee Yun, Shankar Krishnan, and Hossein Mobahi. A unifying view on implicit bias in training linear neural networks. <i>arXiv preprint arXiv:2010.02501</i> , 2020.
722 723	Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization, 2017a.
724 725 726	Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. <i>arXiv preprint arXiv:1710.09412</i> , 2017b.
727	
728	
729	
730	
731	
732	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
740	
750	
751	
752	
753	
754	
755	

A EXPERIMENT DETAILS

In this section, we provide the detailed experiment settings.

760

A.1 DETAILS ON ARCHITECTURES AND DATASETS

We conduct experiments on different neural network architectures, including ResNet18 and
ResNet34 (He et al., 2016), EfficientNetB0 (Tan & Le, 2019), SENet18 (Hu et al., 2018), VGG19 (Simonyan & Zisserman, 2014), MobileNetV2 (Sandler et al., 2018), ShuffleNetV2 (Ma et al., 2018),
RegNet200MF (Radosavovic et al., 2020), SimpleDLA (Yu et al., 2018). We conduct all experiments
using NVIDIA RTX 6000 graphics card.

767 We use CIFAR-10/100 (Krizhevsky et al., 2009) and Imagenet-1k (Deng et al., 2009) as datasets. For 768 CIFAR-10 and CIFAR-100 dataset, each network was trained for 200 epochs using the Stochastic 769 Gradient Descent (SGD) algorithm with cosine learning rate schedule. We choose 27 combinations 770 of hyperparameters in Table 1, and for each hyperparameter we use 3 random seeds and report the 771 average metrics. For the Imagenet-1k dataset, each network was trained for 50 epochs with random data crop and random flip. We use the same optimizer and 27 combinations of hyperparameters as 772 in CIFAR-10 and CIFAR-100 experiments. It is worth noting that we make minor adjustments on 773 hyperparameters for certain networks to ensure stable training. For example, in the case of VGG19, 774 the training is unable to converge when the learning rate is set to 0.1; therefore, we adjust it to 0.05. 775

776 777

784

785

A.2 CORRELATION PLOTS

We show the correlation plot of average regions and test accuracy in Figure 8. The figure consists of results from training different networks on CIFAR-10 dataset with SGD for 200 epochs, using the hyperparameters specified in Table 1. The results show that different networks have different number of regions, ranging from 2 to 20. However, the correlation of test accuracy and average number of regions are consistently high in all the networks.

B HOW TO CALCULATE THE NUMBER OF REGIONS

In this section, we study different methods to calculate the region count, and discuss their impact on the results. Since it is impossible in practice to calculate the predictions of an infinite number of data points on the hyperplane, we select grid points from the hyperplane to calculate the region count.

Assume we have divided a region of the hyperplane into several equidistant small squares. We can use an algorithm similar to breadth-first search to calculate the number of connected components within these small squares, thereby determining the number of regions. Here, we use a 2-dimensional hyperplane as an example (the 1-dimensional case can be considered a degenerate version of this algorithm). The algorithm for calculating the number of regions in this setup is given in Algorithm 1.

Therefore, it is necessary to determine the granularity of splits for the plane. We experimented with different setups of splitting parameters, and the results averaged by 100 independent trials are presented in Table 4. From the results, using 200 grid points in the 1D case and 30x30 grid points in the 2D case is an optimal choice. Splitting the plane into fewer points results in an inadequate approximation of regions, while increasing the number of points does not significantly enhance accuracy but incurs greater computational costs. Therefore, in our experiments, we split the plane into 200 grid points for the 1D case and 30x30 grid points for the 2D case.

Subsequently, we study the number of random samples in calculating the average number of regions.
We experiment with different numbers of hyperplanes, and the results are presented in Table 5. From the results, we know that using 100 samples to calculate the average provides a reliable answer with relatively low computational costs. Therefore, in the experiments we randomly generate 100 lines or planes and calculate the average number of regions.

807 In our paper we use the convex hull of two points $\{\alpha x_1 + (1 - \alpha x_2)\}$ to calculate the region counts 808 in 1D case with $\alpha \in [0, 1]$. We also conduct ablation studies with varied coordinate ranges α . We 809 train ResNet18 on CIFAR10 using hyperparameter in Table 1 in our manuscript, where we vary the range of α and analyze the correlation. The results are shown in Table 6. These studies confirm



Figure 8: The correlation plot of all networks between average regions and test accuracy for CIFAR-10 dataset with optimizer SGD. From the graph we know that the correlations are all very high for different non-linear networks. Different structures of neural networks incur different scope of the average regions.

Algorithm 1 Calculate the Number of Region

843 **Input:** Prediction Matrix P with dimension(w, h). 844 **Output:** Number of connected regions N. 845 1: Initialize a mark matrix M to a zero matrix, with the same dimensions as P846 2: Initialize count of connected regions $N \leftarrow 0$ 847 3: for $i \leftarrow 0$ to w - 1 do 848 for $j \leftarrow 0$ to h - 1 do 4: 849 5: if M[i][j] is already marked then continue 850 6: 7: end if 851 Mark position (i, j) in M as visited 8: 852 Perform Breadth-First Search (BFS) starting from position (i, j)9: 853 10: In BFS, enqueue all neighboring points that have the same value as (i, j) in P and mark 854 them as visited in M855 11: Continue BFS until the queue is empty 856 Increment count of connected regions $N \leftarrow N + 1$ 12: 857 13: end for 858 14: end for 859 15: **return** N 860 861

862

838

839

that expanding the range does not influence the strong correlation between region counts and test accuracy.

865		2	•		
866	Splitting Numbers	Region Counts		Splitting Numbers	Region Counts
867	50	2.74		10×10	2.76
868	50	2.74		10×10	2.70
869	100	2.76		20×20	2.76
870	200	2.78		30×30	2.78
871	300	2.78		40×40	2.78
872	500	2.70		10/(10	2.70
873	500	2.78		50×50	2.78

 Table 4: The mean value of region counts with different splitting numbers in 1d (left) and 2d (right) planes.

Table 5: The mean value of region counts with different number of random samples.

Number of Samples	Region Counts
10	2.24
50	2.56
100	2.78
300	2.80
500	2.79

Table 6: The impact of interpolation range on region counts.

The range of α	Region Counts	Correlation
[0, 1]	3.56	0.98
[-1, 2]	4.47	0.96
[-2, 3]	5.86	0.92
[-3, 4]	6.39	0.93

918 C MORE ABLATION STUDIES

Gradient Optimizers. We calculate the region count of models trained by different optimizers, including SGD, Adam, and Adagrad. The correlation between region count and the generalization gap is consistent for them, as detailed in Table 7.

Table 7: The impact of optimizers on the correlation between region counts and generalization gap.

Optimizer Network	SGD	Adam	Adagrad
ResNet18	0.98	0.92	0.96
ResNet34	0.98	0.92	0.91
VGG19	0.94	0.92	0.87
MobileNet	0.95	0.95	0.99
SENet18	0.98	0.78	0.91
ShuffleNetV2	0.95	0.83	0.99
EfficientNetB0	0.98	0.97	0.99
RegNetX_200MF	0.98	0.95	0.99
SimpleDLA	0.98	0.95	0.88

942 Hyperplane Generation Methods. We explore the effects of different methods for generating the
hyperplane in the input space. In the main experiments, we generate a 1-dimensional plane using
random pairs of samples from the training set and calculated the number of distinct regions between
them. In this section, we apply various techniques for plane generation: selecting two data points
from the test set, choosing one data point from the training set and extending it in a random direction
by a fixed length. We calculate the number of regions for each of these setups. The results in Table 8
are consistent across different hyperplane computational approaches.

Table 8: The impact of calculation methods on the correlation between region counts and generalization gap.

Counting Network	Test	Train	Random
ResNet18	0.98	0.98	0.98
ResNet34	0.98	0.96	0.94
VGG19	0.94	0.89	0.78
MobileNet	0.95	0.94	0.88
SENet18	0.98	0.96	0.99
ShuffleNetV2	0.95	0.95	0.92
EfficientNetB0	0.98	0.98	0.92
RegNetX_200MF	0.98	0.97	0.92
SimpleDLA	0.98	0.97	0.96

972 Proof D 973

977

979

984

989

998

999

1002

1003

974 This section contains the proof of the theorem in this paper. 975

We first prove two lemmas. 976

Lemma 2. The region counts between a pair of data points is upper-bounded by the number of active neurons. For two inputs x_a, x_b , we have $R(x_a, x_b, W) \leq N(x_a, W) + N(x_b, W) + 2$. 978

Proof. If $R(x_a, x_b, W) \leq 2$, then the equation naturally holds. Next we consider $R(x_a, x_b, W) > 2$. 980 From the definition of region count, one can find $R := R(x_a, x_b, W)$ points on the line segment 981 between x_a and x_b , such that the neural network gives different predictions. Denote these points as 982 $\tilde{x}_1, \cdots, \tilde{x}_R$. We have 983

$$f_W(\tilde{x}_i)f_W(\tilde{x}_{i+1}) < 0, 0 \le i \le R-1.$$

985 Consider $\tilde{x}_i, \tilde{x}_{i+1}, \tilde{x}_{i+2}$. Since the neural network gives alternating predictions on these three points, it is nonlinear and has activation sign changes on the line segment connecting them. Therefore, we 986 can find a $1 \le n(i) \le p$, such that $(w_{n(i)}^{\top} \tilde{x}_i)(w_{n(i)}^{\top} \tilde{x}_{i+2}) < 0$. 987

We prove it by contradiction. IF for all $1 \le n(i) \le p$, such that $(w_{n(i)}^{\top} \tilde{x}_i)(w_{n(i)}^{\top} \tilde{x}_{i+2}) \ge 0$. Suppose $\tilde{x}_{i+1} = \lambda \tilde{x}_i + (1-\lambda)\tilde{x}_{i+2}$, then we have

$$f_{W}(\tilde{x}_{i+1}) = \sum_{i=1}^{p} a_{i}\sigma(w_{i}^{\top}x_{i+1}) = \sum_{i=1}^{p} a_{i}\sigma(w_{i}^{\top}x_{i+1})$$
$$= \sum_{i=1}^{p} a_{i}[\lambda\sigma(w_{i}^{\top}x_{i}) + (1-\lambda)\sigma(w_{i}^{\top}x_{i+2})]$$
$$= \lambda f_{W}(\tilde{x}_{i}) + (1-\lambda)f_{W}(\tilde{x}_{i+2}).$$

Therefore $f_W(\tilde{x}_{i+1})$ has the same sign of $f_W(\tilde{x}_i)$ and $f_W(\tilde{x}_{i+2})$, contradict with the condition that they have alternative signs. So we can find a $1 \le n(i) \le p$, such that $(w_{n(i)}^{\top} \tilde{x}_i)(w_{n(i)}^{\top} \tilde{x}_{i+2}) < 0$.

1000 Since $w_{n(i)}^{+}x$ is linear in x, this implies that 1001

$$(w_{n(i)}^{\top}x_a)(w_{n(i)}^{\top}x_b) < 0$$

We also prove it by contradiction. If they have the same sign then the convex combination of them 1004 have the same sign so $(w_{n(i)}^{\top}\tilde{x}_i)(w_{n(i)}^{\top}\tilde{x}_{i+2}) \geq 0$. 1005

We have the following two observations about n(i). Firstly, we can choose an n(i) such that $a_{n(i)}w_{n(i)}^{\dagger}\tilde{x}_{i+2}$ and $f_W(\tilde{x}_{i+2})$ have the same sign, since there exists at least one such neuron that 1008 contributes to the sign change of f_W . This implies that $n(i) \neq n(i+1)$, since $f_W(\tilde{x}_i)$ have alternating signs. Secondly, since $w_{n(i)}^{+}x$ is a linear function in x, it can only changes sign for at most one time. 1010 This implies that $n(i) \neq n(j)$ if $j - i \geq 2$. Putting them together, we know that $n(i) \neq n(j)$ for 1011 $i \neq j$. 1012

Recall that for each $1 \le i \le R-2$, we have $(w_{n(i)}^{\top} x_a)(w_{n(i)}^{\top} x_b) < 0$. Therefore, there exists R-21013 neurons that are activated for either x_a or x_b . This gives $N(x_a, W) + N(x_b, W) \ge R - 2$, which 1014 completes the proof. 1015

1016 1017

1020

Lemma 4. The sharpness of a neural network is lower-bounded by the number of active neurons: $\lambda_{\max}\left(\nabla_W^2 L(W)\right) \geq \frac{r^2}{N^2} \sum_{i=1}^N N(x_i, W).$ 1018 1019

1021 *Proof.* The Hessian of l(W, x, y) can be expressed as

- $\nabla_W^2 l(W, x, y) = \begin{bmatrix} v_1 v_1^\top & \cdots & v_1 v_p^\top \\ \vdots & & \vdots \\ \vdots & & \vdots \end{bmatrix},$ 1023
- 1025

$$\begin{bmatrix} v_p v_1^\top & \cdots & v_p v_p^\top \end{bmatrix}$$

where $v_i = a_i \sigma'(w_i^{\top} x) x$. Suppose $V = [v_1^{\top}, \dots, v_p^{\top}]$. As the nonzero eigenvalue of $V^{\top} V$ and VV^{\top} is the same, this implies that

 $\lambda_{\max}(\nabla_W^2 l(W, x, y)) = \lambda_{\max}(VV^{\top}) = \sum_{i=1}^p \|v_i\|_2^2 = \sum_{i=1}^p \sigma'(w_i^{\top} x) \|x\|^2 \ge \sum_{i=1}^p \sigma'(w_i^{\top} x)r^2.$

From the definition of $\nabla^2_W L(W)$ and the positive definiteness of Hessian matrices, we know that

$$\lambda_{\max}\left(\nabla_W^2 L(W)\right) = \frac{1}{N} \lambda_{\max}\left(\sum_{i=1}^N \nabla_W^2 l(W, x_i, y_i)\right) \ge \frac{1}{N^2} \sum_{i=1}^N \lambda_{\max}\left(\nabla_W^2 l(W, x_i, y_i)\right).$$

Plug in the previous calculation and use the definition of N(x), we have

$$\lambda_{\max}\left(\nabla_{W}^{2}L(W)\right) \geq \frac{r^{2}}{N^{2}} \sum_{i=1}^{N} \sum_{j=1}^{p} \sigma'(w_{i}^{\top}x_{j}) = \frac{r^{2}}{N^{2}} \sum_{i=1}^{N} N(x_{i}, W).$$

Proof of Theorem 1. Theorem 1 is a direct consequence of Assumption 2 and the following two lemmas.

$$\mathbb{E}_{X,X'}[R(X,X',W_t)] = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N R(x_1,x_2,W_t)$$

$$\leq \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N (N(x_i,W_t) + N(x_j,W_t) + 2)$$

$$= \frac{2}{N} \sum_{i=1}^N (N(x_i,W_t) + 1)$$

$$\leq \frac{2N}{r^2} \lambda_{\max}(\nabla_W^2 L(W_t)) + 2$$

$$= O\left(\frac{N}{r^2\eta}\right)$$