UNIVERSITY^{OF} BIRMINGHAM University of Birmingham Research at Birmingham

Norm-based generalisation bounds for deep multiclass convolutional neural networks

Ledent, Antoine; Mustafa, Waleed; Lei, Yunwen; Kloft, Marius

License: None: All rights reserved

Document Version Peer reviewed version

Citation for published version (Harvard):

Ledent, A, Mustafa, W, Lei, Y & Kloft, M 2021, Norm-based generalisation bounds for deep multi-class convolutional neural networks. in *AAAI'21 Proceedings of the Thirty-fifth AAAI Conference on Artificial Intelligence*. Proceedings of the AAAI Conference on Artificial Intelligence. Proceedings of the AAAI Conference on Artificial Intelligence, no. 9, vol. 35, AAAI Press, pp. 8279-8287, 35th AAAI Conference on Artificial Intelligence, Vancouver, Canada, 2/02/21. https://ojs.aaai.org/index.php/AAAI/article/view/17007>

Link to publication on Research at Birmingham portal

General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

•Users may freely distribute the URL that is used to identify this publication.

•Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.

•User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?) •Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

Take down policy

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact UBIRA@lists.bham.ac.uk providing details and we will remove access to the work immediately and investigate.

Norm-based generalisation bounds for deep multi-class convolutional neural networks

Antoine Ledent¹, Waleed Mustafa¹, Yunwen Lei² and Marius Kloft¹

¹Department of Computer Science, TU Kaiserslautern, 67653 Kaiserslautern, Germany ²School of Computer Science, University of Birmingham, Birmingham B15 2TT, United Kingdom

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

Abstract

We show generalisation error bounds for deep learning with 1 two main improvements over the state of the art. (1) Our 2 bounds have no explicit dependence on the number of classes 3 except for logarithmic factors. This holds even when formu-4 lating the bounds in terms of the L^2 -norm of the weight 5 matrices, where previous bounds exhibit at least a square-6 7 root dependence on the number of classes. (2) We adapt the classic Rademacher analysis of DNNs to incorporate weight 8 sharing-a task of fundamental theoretical importance which 9 was previously attempted only under very restrictive assump-10 11 tions. In our results, each convolutional filter contributes only once to the bound, regardless of how many times it is applied. 12 Further improvements exploiting pooling and sparse connec-13 tions are provided. The presented bounds scale as the norms of 14 the parameter matrices, rather than the number of parameters. 15 In particular, contrary to bounds based on parameter count-16 17 ing, they are asymptotically tight (up to log factors) when the weights approach initialisation, making them suitable as a 18 basic ingredient in bounds sensitive to the optimisation pro-19 20 cedure. We also show how to adapt the recent technique of loss function augmentation to our situation to replace spectral 21 norms by empirical analogues whilst maintaining the advant-22 23 ages of our approach.

Introduction

24

Deep learning has enjoyed an enormous amount of success 25 in a variety of engineering applications in the last decade (Kr-26 izhevsky, Sutskever, and Hinton 2012; He et al. 2016; Karras, 27 Laine, and Aila 2018; Silver et al. 2018). However, providing 28 a satisfying explanation to its sometimes surprising gener-29 alisation capabilities remains an elusive goal (Zhang et al. 30 2017; Du et al. 2019; Asadi, Abbe, and Verdu 2018; Good-31 fellow, Shlens, and Szegedy 2015). The statistical learning 32 theory of deep learning approaches this question by providing 33 a theoretical analysis of the generalisation performance of 34 deep neural networks (DNNs) through better understanding 35 of the complexity of the function class corresponding to a 36 given architecture or training procedure. 37 This field of research has enjoyed a revival since 2017 38

This field of research has enjoyed a revival since 2017 with the advent of learning guarantees for DNNs expressed in terms of various norms of the weight matrices and classification margins (Neyshabur, Bhojanapalli, and Srebro 2018; Bartlett, Foster, and Telgarsky 2017; Zhang, Lei, and Dhillon

2018; Li et al. 2019; Allen-Zhu, Li, and Liang 2019). Many 43 improvements have surfaced to make bounds non-vacuous at 44 realistic scales, including better depth dependence, bounds 45 that apply to ResNets (He, Liu, and Tao 2019), and PAC-46 Bayesian bounds using network compression (Zhou et al. 47 2019), data-dependent Bayesian priors (Dziugaite and Roy 48 2018), fast rates (Suzuki 2018), and reduced dependence 49 on the product of spectral norms via data-dependent local-50 isation (Wei and Ma 2019; Nagarajan and Kolter 2019). A 51 particularly interesting new branch of research combines 52 53 norm-based generalisation bounds with the study of how the optimisation procedure (stochastic gradient descent) impli-54 citly restricts the function class (Cao and Gu 2019; Du et al. 55 2019; Arora et al. 2019; Zou et al. 2018; Jacot, Gabriel, and 56 Hongler 2018; Frankle and Carbin 2019). One idea at the 57 core of many of these works is that the weights stay relatively 58 close to initialisation throughout training, reinforcing lucky 59 guesses from the initialised network rather than constructing 60 a solution from scratch. Thus, in this branch of research, it 61 is critical that the bound is negligible when the network ap-62 proaches initialisation, i.e., the number of weights involved 63 is not as important as their size. This observation was first 64 made as early as in (Bartlett 1998). 65

Despite progress in so many new directions, we note that some basic questions of fundamental theoretical importance have remain unsolved. (1) How can we remove or decrease the dependence of bounds on the number of classes? (2) How can we account for weight sharing in convolutional neural networks (CNNs)? In the present paper, we contribute to an understanding of both questions.

Question (1) is of central importance in extreme classification (Prabhu and Varma 2014), where we deal with an extremely high number of classes (e.g. millions). (Bartlett, Foster, and Telgarsky 2017) showed a bound with no explicit class dependence (except for log terms). However, this bound is formulated in terms of the $L^{2,1}$ norms of the network's weight matrices. If we convert the occurring $L^{2,1}$ norms into the more commonly used L^2 norms, we obtain a square-root dependence on the number of classes.

Regarding (2), (Li et al. 2019) showed a bound that accounts for weight sharing. However, this bound is valid only under the assumption of orthonormality of the weight matrices. The assumption of unit norm weights—which is violated by typical convolutional architectures (GoogLeNet,

VGG, Inception, etc.)-makes it difficult to leverage the gen-87 eralisation gains from small weights, and it is a fortiori not 88 easy to see how the bounds could be expressed in terms of 89 distance to initialisation. 90

In this paper, we provide, up to only logarithmic terms, 91 a complete solution to both of the above questions. First, 92 our bound relies only on L^2 norms at the last layer, yet it 93 has no explicit (non-logarithmic) dependence on the number 94 of classes. In deep learning, no generalization bound other 95 than ours has ever achieved a lack of non-logarithmic class 96 dependency with L^2 norms. Second, our bound accounts for 97 weight sharing in the following way. The Frobenius norm 98 of the weight matrix of each convolutional filter contributes 99 only once to the bound, regardless of how many times it is ap-100 plied. Furthermore, our results have several more properties 101 of interest: (i) We exploit the L^{∞} -continuity of nonlinearities 102 such as pooling and ReLu to further significantly reduce the 103 explicit width dependence in the above bounds. (ii) We show 104 how to adapt the recent technique of loss function augmenta-105 tion to our setting to replace the dependence on the spectral 106 norms by an empirical Lipschitz constant with respect to 107 well chosen norms. (iii) Our bounds also have very little ex-108 plicit dependence on architectural choices and rely instead 109 110 on norms of the weight matrices expressed as distance to initialisation, affording a high degree of architecture robust-111 ness compared to parameter-space bounds. In particular, our 112 bounds are negligible as the weights approach initialisation. 113 In parallel to our efforts, (Long and Sedghi 2020) recently 114 made progress on question (2), providing a remedy to the 115 weight-sharing problem. Their work, which appeared at ICLR 116 2020, is independent of ours. This can be observed from 117 the fact that their work and ours were first preprinted on 118 arXiv on the very same day. Their approach is completely 119 different from ours, and both approaches have their merits 120 and disadvantages. We provide an extensive discussion and 121 comparison in Section H. 122

Related Work

123

We now discuss the related work on the statistical learning 124 theory (SLT) of DNNs. The SLT of neural networks can be 125 126 dated back to 1970s, based on the concepts of VC dimension, 153 fat-shattering dimension (Anthony and Bartlett 2002), and 127 154 Rademacher complexities (Bartlett and Mendelson 2002). 128 Here, we focus on recent work in the era of deep learning. 129

156 Let $(x_1, y_1), \ldots, (x_n, y_n)$ be training examples independ-130 ently drawn from a probability measure defined on the 157 131 sample space $\mathcal{Z} = \mathcal{X} \times \{1, \dots, K\}$, where $\mathcal{X} \subset \mathbb{R}^d$, d 158 132 is the input dimension, and K is the number of classes. 159 133 We consider DNNs parameterized by weight matrices $\mathcal{A} =$ 134 $\{A^1, \ldots, A^L\}$, so that the prediction function can be written $F_{\mathcal{A}}(x) = A^L \sigma_{L-1} (A^{L-1} \sigma_{L-2} (\cdots A^1 x))$, where L is the depth of the DNN, $A^l \in \mathbb{R}^{W_{l-1} \times W_l}$, $W_0 = d$, $W_L = K$, and $\sigma_i : \mathbb{R}^{W_i} \mapsto \mathbb{R}^{W_i}$ is the non linearity (including any pooling 135 136 137 138 and activation functions). 139

When providing PAC guarantees for DNNs, a critical quantity is the Rademacher complexity of the network obtained after appending any loss function. The first work in this area (Neyshabur, Tomioka, and Srebro 2015) therefore

focused on bounding the Rademacher complexity of networks satisfying certain norm conditions, where the last layer is one-dimensional. They apply the concentration lemma and a peeling technique to get a bound on the Rademacher complexity of the order $O(\frac{2^L}{\sqrt{n}}\prod_{i=1}^L ||A^i||_{\mathrm{Fr}})$, where $||A||_{\mathrm{Fr}}$ denotes the Frobenius norm of a matrix A. (Golowich, Rakhlin, and Shamir 2018) showed that this exponential dependency on the depth can be avoided by an elegant use of the contraction lemma to obtain bounds of the order $O((\sqrt{L}/\sqrt{n})\prod_{i=1}^{L} ||A^i||_{\mathrm{Fr}})$.¹ The most related work to ours is the spectrally-normalized margin bound by (Bartlett, Foster, and Telgarsky 2017) for multi-class classification. Writing $||A||_{\sigma}$ for the spectral norm, the result is of order $\tilde{O}(M/\gamma)$ with

$$M = \frac{1}{\sqrt{n}} \prod_{i=1}^{L} \|A^{i}\|_{\sigma} \left(\sum_{i=1}^{L} \frac{\|A^{i\top}\|_{2,1}^{\frac{2}{3}}}{\|A^{i}\|_{\sigma}^{\frac{2}{3}}} \right)^{\frac{1}{2}}, \qquad (1)$$

where $||A||_{p,q} = \left(\sum_{j} \left(\sum_{i} |A_{ij}|^{p}\right)^{\frac{q}{p}}\right)^{\frac{1}{q}}$ is the (p,q)-norm, and γ denotes the classification margin.

140

141

142

143

144

145

146

147

148

149

150

151

152

155

160

161

162

163

165

At the same time as the above result appeared, the authors in (Nevshabur, Bhojanapalli, and Srebro 2018) used a PAC Bayesian approach to prove an analogous result², where $W = \max\{W_0, W_1, \dots, W_L\}$ is the width:

$$\tilde{O}\left(\frac{L\sqrt{W}}{\gamma\sqrt{n}}\left(\prod_{i=1}^{L}\|A^{i}\|_{\sigma}\right)\left(\sum_{i=1}^{L}\frac{\|A^{i}\|_{\mathrm{Fr}}^{2}}{\|A^{i}\|_{\sigma}^{2}}\right)^{\frac{1}{2}}\right).$$
 (2)

These results provide solid theoretical guarantees for DNNs. However, they take very little architectural information into account. In particular, if the above bounds are applied to a CNN, when calculating the squared Frobenius norms $||A^i||_{\text{Fr}}^2$, the matrix A^i is the matrix representing the linear operation performed by the convolution, which implies that the weights of each filter will be summed as many times as it is applied. This effectively adds a dependence on the square root of the size of the corresponding activation map at each term of the sum. Furthermore, the L^2 version of the above bound (1) includes a dependence on the square root of the number of classes through the maximum width W of the network. This square-root dependence is not favorable when the number of classes is very large. Although many efforts have been performed to improve the class-size dependency in the shallow learning literature (Lauer 2018; Guermeur 2002, 2007; Koltchinskii and Panchenko 2002; Guermeur 2017; Musayeva, Lauer, and Guermeur 2019; Mohri, Rostamizadeh, and Talwalkar 2018; Lei et al. 2019), extensions of those results to deep learning are missing so far.

In late 2017 and 2018, there was a spur of research effort on the question of fine-tuning the analyses that provided the above bounds, with improved dependence on depth (Golowich, Rakhlin, and Shamir 2018), and some bounds for

¹Note that both of these works require the output node to be one dimensional and thus are not multiclass

^{2}Note that the result using formula (2) can also be derived by expressing (1) in terms of L^2 norms and using Jensen's inequality

recurrent neural networks (Chen, Li, and Zhao 2019; Zhang, 225 166 Lei, and Dhillon 2018)). Notably, in (Li et al. 2019), the au-167 thors provided an analogue of (1) for convolutional networks, 227 168 but only under some very specific assumptions, including 169 orthonormal filters. Those conditions are not satisfied by 170 228 the typical convolutional architectures (GoogLeNet, VGG, 171 229 Inception, etc.). 172

Independently of our work, (Long and Sedghi 2020, to ap-173 pear at ICLR 2020) address the weight-sharing problem using 174 a parameter-space approach. Their bounds scale roughly as 175 the square root of the number of parameters in the model. In 176 contrast to ours, their employed proof technique is more sim-177 ilar to (Li et al. 2019): it focuses on computing the Lipschitz 178 constant of the functions with respect to the parameters. The 179 result by (Long and Sedghi 2020) and ours, which we con-180 trast in detail in Section, both have their merits. In nutshell, 181 the bound by (Long and Sedghi 2020) remarkably comes 182 along without dependence on the product of spectral norms 183 (up to log terms), thus effectively removing the exponential 184 dependence on depth. Our result on the other hand comes 185 along without an explicit dependence on the number of para-186 meters, which can be very large in deep learning. As already 187 noted in (Bartlett 1998), this property is crucial when the 188 weights are small or close to the initialisation. 189

Lastly, we would like to point out that, over the course 190 of the past year, several techniques have been introduced 191 to replace the dependence on the product of spectral norms 192 by an empirical version of it, at the cost of either assuming 193 smoothness of the activation functions (Wei and Ma 2019) 194 or a factor of the inverse minimum preactivation (Nagarajan 195 and Kolter 2019). Slightly earlier, a similar bound to that 196 in (Long and Sedghi 2020) (with explicit dependence on 197 the number of parameters) had already been proved for an 198 unsupervised data compression task (which does not apply 199 to our supervised setting) in (Lee and Raginsky 2019). Re-200 cently, another paper addressing the weight sharing problem 201 appeared on arXiv (Lin and Zhang 2019). In this paper, which 234 202 was preprinted several months after (Long and Sedghi 2020) 235 203 and ours, the authors provided another solution to the weight 236 204 sharing problem, which incorporates elements from both our 237 205 approach and that of (Long and Sedghi 2020): they bound 206 238 the L^2 -covering numbers at each layer independently, but use 207 parameter counting at each layer, yielding *both* an unwanted 208 240 dependence on the number of parameters in each layer (from 241 209 the parameter counting) and a dependence on the spectral 242 210 norms from the chaining of the layers. 211

Further related work includes the following. (Du et al. 244 212 2018) showed size-free bounds for CNNs in terms of the num- 245 213 ber of parameters for two-layer networks. In (Sedghi, Gupta, 246 214 and Long 2019), the authors provided an ingenious way of 247 215 computing the spectral norms of convolutional layers, and 248 216 showed (interestingly) that regularising the network to make 249 217 them approach 1 for each layer is both feasible and beneficial 250 218 to accuracy. Several researchers have also provided interest-251 219 ing insights into DNNs from different perspectives, includ- 252 220 ing through model compression (Neyshabur, Bhojanapalli, 253 221 and Srebro 2018), capacity control by VC dimensions (Har- 254 222 vey, Liaw, and Mehrabian 2017), and the implicit restriction 255 223 on the function class imposed by the optimisation proced- 256 224

ure (Arora et al. 2018; Zhou et al. 2019; Nevshabur et al. 2019, to appear; Suzuki 2018; Du et al. 2019; Jacot, Gabriel, and Hongler 2018; Arora et al. 2019).

Contributions in a Nutshell

In this section, we state the simpler versions of our main results for specific examples of neural networks. The general results are described in in more technical detail in Section A.

Fully Connected Neural Networks

231

232

233

239

243

In the fully connected case, the bound is particularly simple:

Theorem 1 (Multi-class, fully connected). Assume that we are given some fixed reference matrices M^1, M^2, \ldots, M^L representing the initialised values of the weights of the network. Set $\widehat{R}_{\gamma}(F_{\mathcal{A}}) = (1/n)(\#(i : F(x_i)_{y_i} < \gamma +$ $\max_{j \neq y_i} F(x_i)_j)$, where # denotes the cardinality of a set. With probability at least $1 - \delta$, every network F_A with weight matrices $\mathcal{A} = (A_1, A_2, \dots, A_L)$ and every margin $\gamma > 0$ satisfy:

$$\mathbb{P}(\arg\max(F_{\mathcal{A}}(x)_j) \neq y) \le \widehat{R}_{\gamma}(F_{\mathcal{A}}) +$$
(3)

$$\widetilde{\mathcal{O}}\left(\frac{\max_{i=1}^{n} \|x_i\|_{\mathrm{Fr}} R_{\mathcal{A}}}{\gamma \sqrt{n}} \log(\bar{W}) + \sqrt{\frac{\log(1/\delta)}{n}}\right), \quad (4)$$

where $W = \overline{W} = \max_{i=1}^{L} W_i$ is the maximum width of the network, and

$$R_{\mathcal{A}} := \rho_L \max_i \|A_{i, \cdot}^L\|_{\mathrm{Fr}} \left(\prod_{i=1}^{L-1} \rho_i \|A^i\|_{\sigma}\right)$$
(5)

$$\left(\sum_{i=1}^{L-1} \frac{\left(\|A^{i} - M^{i}\|_{2,1}^{2/3}}{\|A^{i}\|_{\sigma}^{2/3}} + \frac{\|A^{L}\|_{\mathrm{Fr}}^{2/3}}{\max_{i} \|A_{i, \cdot}^{L}\|_{\mathrm{Fr}}^{2/3}}\right)^{\frac{3}{2}}.$$
 (6)

Note that the last term of the sum does not explicitly contain architectural information, and assuming bounded L^2 norms of the weights, the bound only implicitly depends on W_i for $i \leq L-1$ (through $||A^i - M^i||_{2,1} \leq$ $\sqrt{W_{i-1}} \|A^i - M^i\|_{2,1}$, but not on W_L (the number of classes). This means the above is a class-size free generalisation bound (up to a logarithmic factor) with L^2 norms of the last layer weight matrix. This improves on the earlier $L^{2,1}$ norm result in (Bartlett, Foster, and Telgarsky 2017). To see this, let us consider a standard situation where the rows of the matrix A^L have approximately the same L^2 norm, i.e., $||A_{i,\cdot}^L||_2 \approx a$. (In Section I in the Appendix, we show that similar conditions hold except on a subset of weight space of asymptotically vanishing measure and further discuss possible behaviour of the norms.) In this case, our bound involves $||A^L||_{\rm Fr} \simeq \sqrt{W_L}a$, which incurs a squareroot dependency on the number of classes. As a comparison, the bound in (Bartlett, Foster, and Telgarsky 2017) involves $||(A^L)^\top||_{2,1} \simeq W_L a$, which incurs a linear dependency on the number of classes. If we further impose an L_2 -constraint on the last layer as $||A^L||_{\text{Fr}} \leq a$ as in the SVM case for a constant a (Lei et al. 2019), then our bound would enjoy a logarithmic dependency while the bound in (Bartlett, Foster,

and Telgarsky 2017) enjoys a square-root dependency. This 257 cannot be improved without also changing the dependence 258 on n. Indeed, if it could, we would be able to get good guar-259 antees for classifiers working on fewer examples than classes. 260 Furthermore, in the above bound, the dependence on the 261 spectral norm of A^L in the other terms of the sum is reduced 262

to a dependence on $\max_i ||A_{i,\cdot}^L||_2^3$. Both improvements are 263

based on using the L^{∞} -continuity of margin-based losses. 264

Convolutional Neural Networks 265

Our main contribution relates to CNNs. To avoid blinding 266 the reader with notation, we present first simple versions of 267

our results. 268

Two-layers The topic of the present paper is often notation-269 280 ally cumbersome, which imposes an undue burden on readers. 270 281 Therefore, we first present a particular case of our bound for 271 282 a two-layer network composed of a convolutional layer and a 272 fully connected layer with a single input channel, with expli-273

*cit pre chosen norm constraints*⁴. Note that the restrictions ²⁸⁴ 274

are purely based on notational and reader convenience: more 285 275

276 general results are presented later and in the Appendix.

2-layer Notation: Consider a two layer network with a 287 convolutional layer and a fully connected layer. Write d, Cfor the dimensions of the input space and the number of 289 classes respectively. We write w for the *spacial* dimension 290 of the hidden layer after pooling⁵. Write \hat{A}^1, A^2 for the sets 291 of weights of the first and second layer, with the weights appearing only once in the convolutional case. A^2 is a mat-293 rix whilst A^1 can be arranged as a tensor or unfolded as a matrix. The matrix, which we denote by \tilde{A}^1 , representing 295 the convolution operation presents the weights of the matrix 296 A_1 repeated as many times as the filters are applied. For any 297 input $x \in \mathbb{R}^d$, we write $|x|_0$ for the maximum L^2 norm of any single convolutional patch of x. For instance, if x is an 299 image and the network applies 3×3 convolutions with stride 1, $|x|_0$ would be the maximum L^2 norm of any sub-image 300 of size $3 \times 3 \times m$ where m is the number of channels. The 301 network is represented by the function 302

$$F(x) = A^2 \sigma(\tilde{A}^1 x),$$

277 where σ denotes the non linearities (including both pooling) 305

and activation functions). As above, M^1, M^2 are the initial-278 ised weights. 279

Theorem 2. Let $a_1, a_2, a_*, b_0, b_1 > 0$. Suppose that the distribution over inputs is such that $|x|_0 \leq b$ a.s. With probabil- $\begin{array}{l} ity > 1 - \delta \text{ over the draw of the training set, for every network} \\ \mathcal{A} = (A^1, A^2) \text{ with weights satisfying } \|(A^1 - M^1)^\top\|_{2,1} \leq \end{array}$

⁴It is common practice to leave the post hoc step to the reader in this way. Cf.,e.g., (Long and Sedghi 2020))

⁵This is less than the number of convolutional patches in the 320 input and is not influenced by the number of filters applied. 321

 $a_1, ||A^2 - M^2||_{\mathrm{Fr}} \leq a_2 \text{ and } \sup_{c < C} ||A^2_{c, \bullet}||_2 \leq a_*, \text{ if }$ $\sup_{i < n} \|\tilde{A}^1 x_n\|_2 \le b_1$, then

$$\mathbb{P}\left(\arg\max_{j} (F_{\mathcal{A}}(x)_{j}) \neq y\right) \tag{7}$$

$$\leq \widehat{R}_{\gamma}(F_{\mathcal{A}}) + 3\sqrt{\frac{\log(\frac{2}{\delta})}{2n}} + \frac{\mathcal{C}}{\sqrt{n}}\mathcal{R}\left[\log_{2}(n^{2}\mathcal{D})\right]^{\frac{1}{2}}\log(n),$$

where C is an absolute constant,

$$\mathcal{R}^{2/3} = \left[b_0 a_1 \max\left(\frac{1}{b_1}, \frac{\sqrt{w}a_*}{\gamma}\right) \right]^{2/3} + \left[\frac{b_1 a_2}{\gamma}\right]^{2/3},$$
(8)

and the quantity in the log term is \mathcal{D} $\max(b_0a_1Wa_*/b_1, b_1a_2C/\gamma)$ where W is the number of hidden neurons before pooling.

Remarks:

286

294

303

304

306

307

308

309

310

311

312

315

316

317

318

319

- 1. Just as in the fully connected case, the implicit dependence on the number of classes is only through an L^2 norm of the full last layer matrix. b_1 is a an upper bound on the L^2 norms of hidden activations.
- 2. a_1 is the norm of the filter matrix A^1 , which counts each 288 filter only once regardless of how many times it is applied. This means our bound enjoys only logarithmic dependence on input size for a given stride, and takes weight sharing 292 into account.
 - 3. As explained in more detail at the end of Appendix H, there is also no explicit dependence on the size of the filters and the bound is stable through up-resolution. In fact, there is no explicit non logarithmic dependence on architectural parameters, and the bounds converges to 0 as a_1, a_2 tend to zero (in contrast to parameter space bounds such as (Long and Sedghi 2020)).
 - 4. a_* replaces the spectral norm of A^2 , and is only equal to the maximum L^2 norm of the second layer weight vectors corresponding to each class. This improvement, just as the improved dependence on the number of classes, comes from better exploiting the continuity of margin based losses with respect to the L^{∞} norm.
 - 5. The spectral norm of the first layer matrix \tilde{A}_1 is not neccessary and is absorbed into an empirical estimate of the hidden layer norms. The first term in the max relates to the estimation of the risk of a test point presenting with a hidden layer norm higher than (a multiple of) b_1 .
 - 6. b_0 refers to the maximum L^2 norm of a single convolutional patch over all inputs and patches. In particular, the bound exploits the sparsity of connections in CNNs.

Α result for the multi-layer case We assume we are given training and testing points $(x, y), (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ drawn iid from any probability distribution over $\mathbb{R}^d \times \{1, 2, \dots, C\}$. We suppose we have a convolutional architecture so that for each filter matrix $A^l \in \mathbb{R}^{m_l \times d_l}$ from layer l-1 to layer l, we can construct a larger matrix \tilde{A}^l representing the corresponding (linear) convolutional operation. The 0^{th} layer is the input,

³Replacing the $L^{2,1}$ norm by a L^2 norm without accumulating factors of the numbers of classes is the more substantial contribution. 313 On the other hand, the replacement of the spectral norm is down to 314 better Lipschitz management and has probably been achieved elsewhere. We know of one paper that provides a similar improvement under specific assumptions (the last layer weights being fixed and initialised as independent Bernouilli distributions) (Zou et al. 2018)

whist the L^{th} layer is the output/loss function. We write w_l 355 322 for the *spacial* width at layer l, W_l for the total width at layer 356 323 l (including channels), and W for $\max_l W_l$. For simplicity 324 357 of presentation, we assume that the activation functions are 325 358 composed only of ReLu and max pooling. 326

Theorem 3. With probability $\geq 1-\delta$, every network F_A with fliter matrices $\mathcal{A} = \{A^1, A^2, \dots, A^L\}$ and every margin $\gamma > 0$ satisfy:

$$\mathbb{P}\left(\arg\max_{j} (F_{\mathcal{A}}(x)_{j}) \neq y\right)$$

$$\leq \widehat{R}_{\gamma}(F_{\mathcal{A}}) + \widetilde{\mathcal{O}}\left(\frac{R_{\mathcal{A}}}{\sqrt{n}}\log(\bar{W}) + \sqrt{\frac{\log(1/\delta)}{n}}\right), \quad (9)$$

where \overline{W} is the maximum number of neurons in a single layer (before pooling) and

$$R_{\mathcal{A}}^{2/3} = \sum_{l=1}^{L} (T_l)^{2/3}$$

for where $T_l =$

$$B_{l-1}(X) \| (A^l - M^l)^\top \|_{2,1} \sqrt{w_l} \max_{U \le L} \frac{\prod_{u=l+1}^U \| \tilde{A}^u \|_{\sigma'}}{B_U(X)}$$

if
$$l \leq L - 1$$
 and for $l = L$, $T_l = \frac{B_{L-1}(X)}{\gamma} \|A^L - M^L\|_{\text{Fr}}.$

Here, w_l *is the width at layer l after pooling. By convention,* 380 327 $b_L = \gamma$, and for any layer l_1 , $B_{l_1}(X) := \max_i \left| F^{0 \to l_l}(x_i) \right|_{l_1}$ 328 381 denotes the maximum l^2 norm of any convolutional patch 329 382 of the layer l_1 activations, over all inputs. For $l \leq L - 1$, 330 383 $\|\tilde{A}_l\|_{\sigma'} \leq \|\tilde{A}_l\|$ denotes the maximum spectral norm of any 331 384 matrix obtained by deleting, for each pooling window, all but 332 385 one of the corresponding rows of \tilde{A} . In particular, for l = L, $\|\tilde{A}^L\|_{\sigma'} = \rho_L \max_i \|A_{i,*}^L\|_2$. Here $A_{i,*}^L$ denotes the *i*'th row of 333 386 334 387 A^L , and $\|\cdot\|_2$ denotes the Frobenius norm. 335 388

Similarly to the two layer case above, a notable property 389 336 of the above bounds is that the norm involved is that of the 390 337 matrix A^l (the filter) instead of \tilde{A}^l (the matrix representing 391 338 the full convolutional operation), which means we are only 392 339 adding the norms of each filter once, regardless of how many 393 340 patches it is applied to. As a comparison, although the gen-341 realization bound in (Bartlett, Foster, and Telgarsky 2017) 342 also applies to CNNs, the resulting bound would involve the 394 343 whole matrix A ignoring the structure of CNNs, yielding an 344 extra factor of O_{l-1} instead of $\sqrt{w_l}$, where O_l denotes the 345 number of convolutional patches in layer l: through exploit-346 ing weight sharing, we remove a factor of $\sqrt{O_{l-1}}$ in the l^{th} 347 term of the sum compared to a standard the result in (Bartlett, 348 Foster, and Telgarsky 2017), and we remove another factor 349 of $\sqrt{O_{l-1}/w_l}$ through exploitation of the L^{∞} continuity of 350 max pooling and our use of L^{∞} covering numbers. 351

A further significant improvement is in replacing the factor 352 $||X||_{2,2} \prod_{i=1}^{l-1} ||\tilde{A}_i||_{\sigma}$ from the classic bound by $B_{l-1}(X)$, 353 which is the maximum L^2 norm of a single convolutional 354

patch. This implicitly removes another factor of $\sqrt{O_{l-1}}$, this time from the local connection structure of convolutions.

We note that it is possible to obtain more simple bounds without a maximum in the definition of T_l by using the spectral norms to estimate the norms at the intermediary layers.

359

360

361 362

363

364

365

366

367

368

369

370

371

372

373

374

375

376 377 378

379

Empirical spectral norms; Lipschitz augmentation

A commonly mentioned weakness of norm-based bounds is the dependence on the product of spectral norms from above. In the case of fully connected networks, there has been a lot of progress last year on how to tackle this problem. In particular, it was shown in (Nagarajan and Kolter 2019) and in (Wei and Ma 2019) that the products of spectral norms can be replaced by empirical equivalents, at the cost of either a factor of the minimum preactivation in the Relu case (Nagarajan and Kolter 2019), or Lipschitz constant of the *derivative* of the activation functions if one makes stronger assumptions (Wei and Ma 2019). In the appendix, we adapt some of those techniques to our convolutional, ReLu situation and find that the quantity $\rho_l^{\mathcal{A}}$ can be replaced in our case by:

$$\rho_l^{\mathcal{A}} = \max\left(\max_i \max_{\tilde{l} \ge l} \frac{\rho_{l_1 \to l_2}^{\mathcal{A}, x_i}}{B_{l_2}(X)}, \max_i \max_{\tilde{l} \ge l} \frac{\theta_{l_1 \to l_2}^{\mathcal{A}, x_i}}{E_{l_2}(X)}\right)$$

where $E_l(X)$ denotes the minimum preactivation (or distance to the max/second max in max pooling) at layer lfor over every input, $\rho_{l_1 \to l_2}^{\mathcal{A}, x_i}$ (resp. $\theta_{l_1 \to l_2}^{\mathcal{A}, x_i}$) is the Lipschitz constant of gradient of $F^{l_1 \to l_2}$ with respect to the norms $|\cdot|_{\infty,l_1}$ and $|\cdot|_{l_2}$ (resp. $|\cdot|_{\infty,l_1}$ and $|\cdot|_{\infty}$). These quantities can easily be calculated explicitly: if $M = \nabla_{F^{0 \to l_1}(x_i)} F^{l_1 \to l_2}$ so that locally around $F^{0 \to l_1}(x_i)$, $F^{l_1 \to l_2}(x) = Mx$, then $\theta_{l_1 \to l_2}^{\mathcal{A}, x_i} = \|M^{\top}\|_{1,\infty}$ and $\rho_{l_1 \to l_2}^{\mathcal{A}, x_i} = \max_M \|M'\|_{1,2}$ where M' runs over all sub matrices of M obtained by keeping only the rows corresponding to a single convolutional patch of layer l_2 .

Note that an alternative approach is to obtain tighter bounds on the worst case Lipschitz constant. Theorem A.1 in the Appendix is a variation of Theorem 3 involving the explicit worst case Lipschitz constants across layers instead of spectral norms. These quantities can then be independently bounded, or made small via regularisation using recent techniques developed in, e.g., (Fazlyab et al. 2019; Latorre, Rolland, and Cevher 2020).

General proof strategy

Some key aspects of our proofs and general results rely on using the correct norms in activation spaces. On each activation space, we use the norm $|\cdot|_{\infty}$ to refer to the maximum absolute value of each neuron in the layer, the norm $|\cdot|_l$ to refer to the the maximum l^2 norm of a single convolutional patch (at layer l) and $|\cdot|_{\infty,l}$ for the maximum l^2 norm of a single pixel viewed as a vector over channels. Using these norms, we can for each pair of layers l_1, l_2 define the Lipschitz constant $\rho_{l_1 \to l_2}$ is the Lipschitz constant of the subnetwork $F^{l_1 \to l_2}$ with respect to the norms $|\cdot|_{\infty,l_1}$ and $|\cdot|_{l_2}$. Using those norms we can formulate a cleaner extention of Theorem 3 where the

quantity $R_{\mathcal{A}}$ can be replaced by

$$\left[\sum_{l=1}^{L-1} \left(B_{l-1}(X) \|A^{l} - M^{l}\|_{2,1} \max_{\tilde{l} > l} \frac{\rho_{l \to \tilde{l}}}{B_{\tilde{l}}(X)}\right)^{2/3} + \left(\frac{B_{L-1}(X)}{\gamma} \|A^{L} - M^{L}\|_{\mathrm{Fr}}\right)^{2/3}\right]^{3/2},$$

where for any layer l_1 , $B_{l_1}(X) := \max_i |F^{0 \to l_l}(x_i)|_{l_1}$ denotes 441 395 the maximum l^2 norm of any convolutional patch of the layer 442 396 l_1 activations, over all inputs. $B_L(X) = \gamma$. Our proofs derive 397 443 this result, and the Theorems above follow. cf. Theorem A.1⁶. 444 398 399 In the rest of this Section, we sketch the general strategy 445 of the proof, focusing on the (crucial) one-layer step. At this 446 400 point, we need to introduce notation w.r.t. the convolutional 401 447 channels: we will collect the data matrix of the previous layer 448 402 in the form of a tensor $X \in \mathbb{R}^{n \times U \times d}$ consisting of all the 449 403 convolutional patch stacked together: if we fix the first index 450 404 (sample i.d.) and the second index (patch i.d.), we obtain 405 451 a convolutional patch of the corresponding sample. For a 406 set of weights $A \in \mathbb{R}^{d \times m}$, the result of the convolutional 407 operation is written XA where is defined by $(XA)_{u,i,j} =$ 408 $\sum_{o=1}^{d} X_{u,i,o} A_{o,j} \text{ for all } u, i, j.$ 409 A key step in bounding the capacity of NN's is to bound the 410

covering numbers of individual layers. Recall the definition. 411 **Definition 1** (Covering number). Let $V \subset \mathbb{R}^n$ and $\|\cdot\|$ 412 be a norm in \mathbb{R}^n . The covering number w.r.t. $\|\cdot\|$, de-413 noted by $\mathcal{N}(V, \epsilon, \|\cdot\|)$, is the minimum cardinality m 414 of a collection of vectors $\mathbf{v}^1, \ldots, \mathbf{v}^m \in \mathbb{R}^n$ such that 415 452 $\sup_{\mathbf{v}\in V} \min_{j=1,\ldots,m} \|\mathbf{v}-\mathbf{v}^j\| \leq \epsilon. \text{ In particular, if } \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ is a function class and $X = (x_1, x_2, \ldots, x_n) \in \mathcal{X}^n$ are 416 453 417 454 data points, $\mathcal{N}(\mathcal{F}(X), \epsilon, (1/\sqrt{n}) \| \cdot \|_2)$ is the minimum car-418 dinality m of a collection of functions $\mathcal{F} \ni f^1, \ldots, f^m$: 419 456 $\mathcal{X} \to \mathbb{R}$ such that for any $f \in \mathcal{F}$, there exists $j \leq m$ 420 457 such that $\sum_{i=1}^{n} (1/n) \left| f^j(x_i) - f(x_i) \right|^2 \leq \epsilon^2$. Similarly, 421 $\mathcal{N}(\mathcal{F}(X), \overline{\epsilon}, \|\cdot\|_{\infty})$ is the minimum cardinality m of a 422 459 collection of functions $\mathcal{F} \ni f^1, \ldots, f^m : \mathcal{X} \to \mathbb{R}$ such 423 460 that for any $f \in \mathcal{F}$, there exists $j \leq m$ such that $i \leq m$ 424 $\left|f^{j}(x_{i}) - f(x_{i})\right| \leq \epsilon.$ n, 425

If we apply classical results on linear classifiers as is done 426 427 in (Bartlett, Foster, and Telgarsky 2017) (where results on 428 L^2 covering numbers are used) by viewing a convolutional layer as a linear map directly, we cannot take advantage of 429 weight sharing. In this work, we circumvent this difficulty 430 by applying results on the L^{∞} covering numbers of classes 431 of linear classifiers to a different problem where each "(con-432 volutional patch, sample, output channel)" combination is 433 mapped into a higher dimensional space to be viewed as a 434 single data point, as explained below. A further reduction in 435 explicit dependence on architectural parameters is achieved 436 by leveraging the L^{∞} -continuity of margin based loss func-437 tions, ReLu activation functions, and pooling. We will make 438 use of the following proposition from (Zhang 2002) (The-439 orem 4, page 537). 440

Proposition 4. Let $n, d \in \mathbb{N}$, a, b > 0. Suppose we are given n data points collected as the rows of a matrix $X \in \mathbb{R}^{n \times d}$, with $||X_{i,\bullet}||_2 \leq b, \forall i = 1, ..., n$. For $U_{a,b}(X) = \{X\alpha : ||\alpha||_2 \leq a, \alpha \in \mathbb{R}^d\}$, we have

$$\log \mathcal{N}\left(U_{a,b}(X),\epsilon, \|\cdot\|_{\infty}\right) \leq \frac{36a^2b^2}{\epsilon^2}\log_2\left(\frac{8abn}{\epsilon} + 6n + 1\right).$$

Note that this proposition is much stronger than Lemma 3.2 in (Bartlett, Foster, and Telgarsky 2017). In the latter, the cover can be chosen independently of the data set, and the metric used in the covering is an L^2 average over inputs. In Proposition 4, the metric used in the covering is a maximum over all inputs, and the data set must be chosen in advance, though the size of the cover only depends (logarithmically) on the sample size⁷.

Using the above proposition on the auxiliary problem based on (input, convolutional patch, ouput channel) triplets, we can prove the following bounds for the one layer case:

Proposition 5. Let positive reals (a, b, ϵ) and positive integer m be given. Let the tensor $X \in \mathbb{R}^{n \times U \times d}$ be given with $\forall i \in \{1, 2, ..., n\}, \forall u \in \{1, 2, ..., U\}, ||X_{i,u}, \cdot||_2 \leq b$. For any choice of reference matrix M, we have

$$\begin{split} \log \mathcal{N}\left(\{XA: A \in \mathbb{R}^{d \times m}, \|A - M\|_{\mathrm{Fr}} \leq a\}, \epsilon, \|\cdot\|_{\infty}\right) \\ \leq \frac{36a^2b^2}{\epsilon^2} \log_2\left[\left(\frac{8ab}{\epsilon} + 7\right)mnU\right], \end{split}$$

where the norm $\|\cdot\|_{\infty}$ is over the space $\mathbb{R}^{n \times U \times m}$.

Sketch of proof: By translation invariance, it is clear that we can suppose M = 0. We consider the problem of bounding the L^{∞} covering number of $\{(v_i^{\top}X^j)_{i \leq I, j \leq J} : \sum_{i \leq I} ||v_i||_2^2 \leq a^2\}$ (where $X^j \in \mathbb{R}^{d \times n}$ for all j) with only logarithmic dependence on n, I, J. Here, I plays the role of the number of output channels, while J plays the role of the number of convolutional patches. We now apply the above Proposition 4 on the $nIJ \times dI$ matrix constructed as follows:

(X^1	0 \mathbf{v}^1		$\begin{pmatrix} 0 \\ 0 \end{pmatrix}$	۲ ۱
	0	<i>X</i>	· · · · · · ·	0	
	$ \begin{array}{c} 0 \\ \mathbf{v}^2 \end{array} $	0	•••	X^1	
	$\begin{array}{c} \Lambda \\ 0 \end{array}$	X^2	•••		
				$\frac{\dots}{V^2}$	
	X^3	0	· · · · · · ·		
	$\frac{\dots}{\mathbf{Y}^J}$			· · · · 0	
	0	X^J		0	
			•••	$\frac{\dots}{\mathbf{V}^J}$	
<u>۱</u>	0	0	• • •		

with the corresponding vectors being constructed as $(v_1, v_2, \ldots, v_I) \in \mathbb{R}^{dI}$.

⁶Note that our assumption that the *worst case* Lipschitz constants are bounded removes some of the interactions between layers, yielding a simpler final formula compared to (Wei and Ma 2019; Nagarajan and Kolter 2019)

⁷We note that the proof is also much more obscure, although it is far more approachable to prove an analogous result with a squared log term instead, by going via the shattering dimension.

If we compose the linear map on $\mathbb{R}^{n \times d}$ represented by 503 463 $(v_1, v_2, \dots, v_I)^{\top}$ with k real-valued functions with L^{∞} 464 504 Lipschitz constant 1, the above argument yields comparable 465 505 bounds on the $\|\cdot\|_2$ covering number of the composition, los-466 506 ing a factor of \sqrt{k} only (for the last layer, k = 1, and for 467 507 convolutional layers, k is the number of neurons in the layer 468 508 left after pooling). 469 509

The proposition above is only enough to deal with the last layer, or a purely l^2 version of our bounds. To prove Theorem 3, which involves $\|\cdot\|_{2,1}$ norms, it is necessary to show the following extension: 512 513

Proposition 6. Let positive reals (a, b, ϵ) and positive integer m be given. Let the tensor $X \in \mathbb{R}^{n \times U \times d}$ be given with $\forall i \in \{1, 2, ..., n\}, \forall u \in \{1, 2, ..., U\}, ||X_{i,u, \bullet}||_2 \leq b$. For any fixed M:

514

515

516

517

518

519

520

521 522

523

524

538

539

540

541

542

$$\log \mathcal{N}\left(\{XA: A \in \mathbb{R}^{d \times m}, \|A - M\|_{2,1} \le a\}, \epsilon, \|\cdot\|_*\right)$$
$$\leq \frac{64a^2b^2}{\epsilon^2} \log_2\left[\left(\frac{8ab}{\epsilon} + 7\right)mnU\right],$$

474 where the norm $\|\cdot\|_*$ over the space $\mathbb{R}^{n \times U \times m}$ is defined by

475
$$||Y||_* = \max_{i \le n} \max_{j \le U} \left[\sum_{k=1}^m Y_{i,j,k}^2 \right]^2$$

Sketch of proof: We first assume fixed bounds on the L^2 525 476 526 norms $||A^{i, \cdot}||_2 = a_i$ of each filter, and use Proposition 5 with 477 527 m = 1 for each output channel with a different granularity 478 ϵ_i . We then optimize over the choice of ϵ_i , and make the 528 479 result apply to the case where only $a = \sum_{i} a_i \ge ||A||_{2,1}$ is 529 480 fixed in advance by l^1 covering the set of possible choices for 530 481 531 (a_1, a_2, \ldots, a_m) for each a, picking a cover for each such 482 532 choice and taking the union. We accumulate a factor of 2 483 533 because of this approach, but to our knowledge, it is not 484 534 possible to rescale the inputs by factors of $\sqrt{a_i}$ as was done 485 in (Bartlett, Foster, and Telgarsky 2017), as the input samples 535 486 536 in an L^{∞} covering number bound must be chosen in advance. 487 537

We can now **sketch the proof of the Theorem 7 :** we use the loss function

$$l(x_{i}, y_{i}) = \max \left[\lambda_{b_{1}} (\|\sigma(\tilde{A}^{1}x_{i})\|_{2} - b_{1}), \\ \lambda_{\gamma} (\max_{j \neq y} (A^{2}\sigma(\tilde{A}^{1}x_{i}))_{j} - (A^{2}\sigma(\tilde{A}^{1}x_{i}))_{y_{i}}) \right],$$

where for any $\theta > 0$ the ramp loss λ_{θ} is defined by 488 $\lambda_{\theta} = 1 + \min(\max(x, -\theta), 0)/\theta$. This loss incorporates 543 489 the following two failure scenarios: (1) the L^2 norm of the 490 544 hidden activations exceed a multiple of b_1 (2) The activations 545 491 behave normally but the network still outputs a wrong pre-546 492 diction. Since pooling is continuous w.r.t. the infty norm, the 547 493 above results for the one layer case applied to a layer yields 548 494 an ϵ cover of hidden layer w.r.t to the L^{∞} norm. The contri-549 495 butions to the error source (1) therefore follows directly from 550 496 the first layer case. The contribution of the 1st layer cover 551 497 error to (2) must be multiplied $1/\gamma$ and the Lipschitz constant 552 498 of A^2 with respect to the $||_{\infty,1}$ and L^{∞} norms respectively, 553 499 which we estimate by $\sqrt{w}a_*$ since the Euclidean norm of the 500 deviation from the cover at the hidden layer is bounded by 501 \sqrt{w} times the deviation in $||_{\infty,1}$ norm⁸. 502

Remarks and comparison to concurrent work

We have addressed the main problems of weight sharing and dependence on the number of classes. As mentioned earlier, (Long and Sedghi 2020) have recently studied the former problem independently of us. It is interesting to provide a comparison of their and our main results, which we do briefly here and in more detail in the Appendix.

The bound in (Long and Sedghi 2020) scales like $C\sqrt{\frac{\mathcal{W}(\sum_{l=1}^{L} s_l - \log(\gamma)) + \log(1/\delta)}{n}}$, where s_l is an upper bound on the spectral norm of the matrix corresponding to the l^{th} layer, γ is the margin, and \mathcal{W} is the number of parameters, taking weight sharing into account by counting each parameter of convolutional filters only once. The idea of the proof is to bound the Lipschitz constant of the map from the set of weights to the set of functions represented by the network, and use dimension-dependent results on covering numbers of finite dimensional function classes. Remarkably, this doesn't require chaining the layers, which results in a lack of a non logarithmic dependence on the product of spectral norms. Note that the term $\sum_{l=1}^{L} s_l$ comes from a log term via the inequality $\prod (1 + s_i) \leq \exp(\sum s_i)$.

On the other hand, the bound scales at least as the square root of the number of parameters, even if the weights are arbitrarily close to initialisation. In contrast, our bound (3) scales like $O(\sqrt{1/n})$ up to log terms when the weights approach initialisation. Furthermore, if we fix an explicit upper bound on the relevant norms (cf. Theorem C.2)⁹, **the bound then converges to zero** as the bounds on the norms go to zero. In a refined treatment via the NTK literature (cf. (Arora et al. 2019)), explicit bounds would be provided for those quantities via other tools. In addition, a small modification of the proofs can make the constant towards which the post hoc bounds converges at initialisation arbitrarily small at the cost of slightly worse log terms away from initialisation¹⁰.

Finally, note that the main advantages and disadvantages of our bounds compared to (Long and Sedghi 2020) are connected through a tradeoff in the proof where one can decide which quantities go inside or outside the log. In particular, it is not possible to combine the advantages of both. We refer the reader to Appendix H for a more detailed explanation.

Conclusion

We have proved norm-based generalisation bounds for deep neural networks with significantly reduced dependence on certain parameters and architectural choices. On the issue of class dependency, we have completely bridged the gap between the states of the art in shallow methods and in deep learning. Furthermore, we have, simultaneously with (Long and Sedghi 2020), provided the first satisfactory answer to the weight sharing problem in the Rademacher analysis of neural networks. Contrary to independent work, our bounds are norm-based and are negligible at initialisation.

⁸Recall this $||_{\infty,1}$ norm is a supremum over the spacial locations of the L^2 norms over the channel directions.

⁹The bounds in (Long and Sedghi 2020) and other works deal only with this case, leaving the post hoc case to the reader

¹⁰The post hoc step from Thm C.2 to (e.g.) Thm 3 is a discrete equivalent to setting priors on the maximum norms $||(A_l - M_l)^\top||_{2,1}$.

References

607

613

626

635

640

641

642

645

646

653

- ⁵⁵⁵ Allen-Zhu, Z.; Li, Y.; and Liang, Y. 2019. Learning and ⁶⁰⁸ ⁵⁵⁶ Generalization in Overparameterized Neural Networks, Go-
- ⁵⁵⁷ ing Beyond Two Layers. In Wallach, H.; Larochelle, H.; ₆₁₀
- Beygelzimer, A.; d Alché-Buc, F.; Fox, E.; and Garnett, R., 611
- eds., Advances in Neural Information Processing Systems 32, 612
- 560 6155–6166. Curran Associates, Inc.

 561
 Anthony, M.; and Bartlett, P. 2002. Neural Network Learning:
 614

 562
 Theoretical Foundations. ISBN 978-0-521-57353-5. doi:
 615

 563
 10.1017/CBO9780511624216.
 616

- ⁵⁶⁴ Arora, S.; Du, S. S.; Hu, W.; Li, Z.; and Wang, R. 2019. ⁶¹⁷
- 565 Fine-Grained Analysis of Optimization and Generalization 618
- for Overparameterized Two-Layer Neural Networks. *arXiv* 619
 e-prints arXiv:1901.08584. 620
- 568 Arora, S.; Ge, R.; Neyshabur, B.; and Zhang, Y. 2018. 621
- 569 Stronger Generalization Bounds for Deep Nets via a Com- 622
- 570 pression Approach. In Dy, J.; and Krause, A., eds., Pro- 623
- 571 ceedings of the 35th International Conference on Machine 624
- 572 Learning, volume 80 of Proceedings of Machine Learning
- 573 *Research*, 254–263. Stockholm, Sweden: PMLR.
- 574 Asadi, A.; Abbe, E.; and Verdu, S. 2018. Chaining Mutual In- 627
- 575 formation and Tightening Generalization Bounds. In Bengio, 628
- 576 S.; Wallach, H.; Larochelle, H.; Grauman, K.; Cesa-Bianchi,
- 577 N.; and Garnett, R., eds., Advances in Neural Information
- 578 Processing Systems 31, 7234–7243. Curran Associates, Inc.
- ⁵⁷⁹ Bartlett, P. L. 1998. The sample complexity of pattern classi-
- fication with neural networks: the size of the weights is more important than the size of the network. *IEEE Transactions on*
- ⁵⁸² *Information Theory* 44(2): 525–536. doi:10.1109/18.661502.

Bartlett, P. L.; Foster, D. J.; and Telgarsky, M. J. 2017. ⁶³⁶
 Spectrally-normalized margin bounds for neural networks. ⁶³⁷
 6240–6249. Curran Associates, Inc. ₆₃₈

- 586 Bartlett, P. L.; and Mendelson, S. 2002. Rademacher and ⁶³⁹
- 587 Gaussian complexities: Risk bounds and structural results.
- *Journal of Machine Learning Research* 3(Nov): 463–482.
- 589 Cao, Y.; and Gu, Q. 2019. Generalization Bounds of
- 590 Stochastic Gradient Descent for Wide and Deep Neural Net-643 works. *arXiv e-prints* arXiv:1905.13210.
- Chen, M.; Li, X.; and Zhao, T. 2019. On GeneralizationBounds of a Family of Recurrent Neural Networks.
- ⁵⁹⁴ Du, S. S.; Wang, Y.; Zhai, X.; Balakrishnan, S.; ⁶⁴⁷
- ⁵⁹⁵ Salakhutdinov, R. R.; and Singh, A. 2018. How Many ⁶⁴⁸
- 596 Samples are Needed to Estimate a Convolutional Neural 649
- ⁵⁹⁷ Network? In Bengio, S.; Wallach, H.; Larochelle, H.; Grau-⁶⁵⁰
- man, K.; Cesa-Bianchi, N.; and Garnett, R., eds., Advances in 651
 Neural Information Processing Systems 31, 373–383. Curran 652
- 600 Associates, Inc.
- 601 Du, S. S.; Zhai, X.; Poczos, B.; and Singh, A. 2019. Gradi- 654
- ent Descent Provably Optimizes Over-parameterized Neural
 Networks. In *International Conference on Learning Repres*-656
- 604 *entations*. 657
- Dziugaite, G.; and Roy, D. 2018. Data-dependent PAC-Bayes 658 priors via differential privacy . 659

Fazlyab, M.; Robey, A.; Hassani, H.; Morari, M.; and Pappas,G. J. 2019. Efficient and Accurate Estimation of LipschitzConstants for Deep Neural Networks. *CoRR* abs/1906.04893.

Frankle, J.; and Carbin, M. 2019. The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks. In *International Conference on Learning Representations*.

Golowich, N.; Rakhlin, A.; and Shamir, O. 2018. Size-Independent Sample Complexity of Neural Networks. In Bubeck, S.; Perchet, V.; and Rigollet, P., eds., *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, 297–299. PMLR.

Goodfellow, I.; Shlens, J.; and Szegedy, C. 2015. Explaining and Harnessing Adversarial Examples. In *International Conference on Learning Representations*.

Guermeur, Y. 2002. Combining Discriminant Models with New Multi-Class SVMs. *Pattern Analysis & Applications* 5(2): 168–179. ISSN 1433-7541. doi:10.1007/s100440200015.

Guermeur, Y. 2007. VC Theory of Large Margin Multi-Category Classifiers. *Journal of Machine Learning Research* 8: 2551–2594.

Guermeur, Y. 2017. Lp-norm Sauer–Shelah lemma for margin multi-category classifiers. *Journal of Computer and System Sciences* 89: 450 – 473. ISSN 0022-0000. doi: https://doi.org/10.1016/j.jcss.2017.06.003.

Harvey, N.; Liaw, C.; and Mehrabian, A. 2017. Nearly-tight VC-dimension bounds for piecewise linear neural networks. In Kale, S.; and Shamir, O., eds., *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*, 1064–1068. Amsterdam, Netherlands: PMLR.

He, F.; Liu, T.; and Tao, D. 2019. Why ResNet Works? Residuals Generalize. *arXiv e-prints* arXiv:1904.01367.

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep Residual Learning for Image Recognition. In *CVPR*, 770–778. IEEE Computer Society.

Jacot, A.; Gabriel, F.; and Hongler, C. 2018. Neural Tangent Kernel: Convergence and Generalization in Neural Networks. *CoRR* abs/1806.07572.

Karras, T.; Laine, S.; and Aila, T. 2018. A Style-Based Generator Architecture for Generative Adversarial Networks. *CoRR* abs/1812.04948.

Koltchinskii, V.; and Panchenko, D. 2002. Empirical Margin Distributions and Bounding the Generalization Error of Combined Classifiers. *Ann. Statist.* 30(1): 1–50. doi: 10.1214/aos/1015362183.

Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012. ImageNet Classification with Deep Convolutional Neural Networks. In Pereira, F.; Burges, C. J. C.; Bottou, L.; and Weinberger, K. Q., eds., *Advances in Neural Information Processing Systems 25*, 1097–1105. Curran Associates, Inc.

Latorre, F.; Rolland, P.; and Cevher, V. 2020. Lipschitz constant estimation of Neural Networks via sparse polynomial

554

- optimization. In International Conference on Learning Rep- 715 660
- resentations. URL https://openreview.net/forum?id=rJe4_ 716 661 xSFDB. 662 717
- 718 Lauer, F. 2018. Error bounds with almost radical dependence 663
- on the number of components for multi-category classifica-719 664
- tion, vector quantization and switching regression. In Con- 720 665
- férence sur l'Apprentissage automatique (CAp) French Con-721 666
- ference on Machine Learning (FCML), Proc. of the French 667
- 668 Conference on Machine Learning (CAp/FCML). Rouen, 723 France. 669 724
- Lee, J.; and Raginsky, M. 2019. Learning Finite-Dimensional 670 725
- Coding Schemes with Nonlinear Reconstruction Maps. SIAM 671 726
- Journal on Mathematics of Data Science 1: 617–642. doi: 727 672
- 10.1137/18M1234461. 673
- Lei, Y.; Dogan, Ü.; Zhou, D.; and Kloft, M. 2019. Data-729 674
- Dependent Generalization Bounds for Multi-Class Classific- 730 675
- ation. IEEE Trans. Information Theory 65(5): 2995-3021. 731 676 doi:10.1109/TIT.2019.2893916. 677
- Li, X.; Lu, J.; Wang, Z.; Haupt, J.; and Zhao, T. 2019. On 678
- Tighter Generalization Bounds for Deep Neural Networks: 679 734
- 680 CNNs, ResNets, and Beyond.
- Lin, S.; and Zhang, J. 2019. Generalization Bounds for 736 681 Convolutional Neural Networks. 682
- Long, P. M.; and Sedghi, H. 2020. Size-free generalization 683 bounds for convolutional neural networks. In International 684
- Conference on Learning Representations. 685
- 741 Mohri, M.; Rostamizadeh, A.; and Talwalkar, A. 2018. Found-686
- 742 ations of Machine Learning. Adaptive Computation and Ma-687 chine Learning. Cambridge, MA: MIT Press, 2 edition. ISBN 743 688 978-0-262-03940-6. 744 689
- 745 Musayeva, K.; Lauer, F.; and Guermeur, Y. 2019. 690 746 Rademacher complexity and generalization performance of 691
- multi-category margin classifiers. Neurocomputing 342: 6 747 692
- 15. ISSN 0925-2312. doi:https://doi.org/10.1016/j.neucom. 748 693
- 2018.11.096. Advances in artificial neural networks, machine 749 694
- learning and computational intelligence. 695
- Nagarajan, V.; and Kolter, J. Z. 2019. Deterministic PAC- 751 696 Bayesian generalization bounds for deep networks via gener- 752 697 753
- alizing noise-resilience. CoRR abs/1905.13344. 698
- 754 Neyshabur, B.; Bhojanapalli, S.; and Srebro, N. 2018. A 699
- 755 PAC-Bayesian Approach to Spectrally-Normalized Margin 700
- 701 Bounds for Neural Networks. In International Conference 756
- on Learning Representations. openreview.net. 702
- Neyshabur, B.; Li, Z.; Bhojanapalli, S.; LeCun, Y.; and 703
- Srebro, N. 2019, to appear. The role of over-parametrization 704
- in generalization of neural networks. In International Con-705
- ference on Learning Representations. 706
- Neyshabur, B.; Tomioka, R.; and Srebro, N. 2015. Norm-707
- Based Capacity Control in Neural Networks. In GrACEnwald, 708
- P.; Hazan, E.; and Kale, S., eds., Proceedings of The 28th 709
- *Conference on Learning Theory*, volume 40 of *Proceedings* 710 711 of Machine Learning Research, 1376–1401. Paris, France:
- PMLR. 712
- Prabhu, Y.; and Varma, M. 2014. FastXML: A Fast, Accurate 713
- and Stable Tree-classifier for Extreme Multi-label Learning. 714

- In Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14, 263-272. New York, NY, USA: ACM. ISBN 978-1-4503-2956-9. doi:10.1145/2623330.2623651.
- Sedghi, H.; Gupta, V.; and Long, P. M. 2019. The Singular Values of Convolutional Layers. In International Conference on Learning Representations.

722

728

732

733

735

737

738

739

740

750

757

758

- Silver, D.: Hubert, T.: Schrittwieser, J.: Antonoglou, I.: Lai, M.; Guez, A.; Lanctot, M.; Sifre, L.; Kumaran, D.; Graepel, T.; Lillicrap, T.; Simonyan, K.; and Hassabis, D. 2018. A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play. Science 362(6419): 1140-1144. ISSN 0036-8075. doi:10.1126/science.aar6404.
- Suzuki, T. 2018. Fast generalization error bound of deep learning from a kernel perspective. In Storkey, A.; and Perez-Cruz, F., eds., Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics, volume 84 of Proceedings of Machine Learning Research, 1397-1406. Playa Blanca, Lanzarote, Canary Islands: PMLR.
- Wei, C.; and Ma, T. 2019. Data-dependent Sample Complexity of Deep Neural Networks via Lipschitz Augmentation. In Wallach, H.; Larochelle, H.; Beygelzimer, A.; dAlché-Buc, F.; Fox, E.; and Garnett, R., eds., Advances in Neural Information Processing Systems 32, 9725–9736. Curran Associates, Inc.
- Zhang, C.; Bengio, S.; Hardt, M.; Recht, B.; and Vinyals, O. 2017. Understanding deep learning requires rethinking generalization.
- Zhang, J.; Lei, Q.; and Dhillon, I. S. 2018. Stabilizing Gradients for Deep Neural Networks via Efficient SVD Parameterization. In ICML, volume 80 of Proceedings of Machine Learning Research, 5801–5809. PMLR.
- Zhang, T. 2002. Covering Number Bounds of Certain Regularized Linear Function Classes. J. Mach. Learn. Res. 2: 527-550. ISSN 1532-4435. doi:10.1162/153244302760200713. URL https://doi.org/10.1162/153244302760200713.
- Zhou, W.; Veitch, V.; Austern, M.; Adams, R. P.; and Orbanz, P. 2019. Non-vacuous Generalization Bounds at the ImageNet Scale: a PAC-Bayesian Compression Approach. In International Conference on Learning Representations. openreview.net.
- Zou, D.; Cao, Y.; Zhou, D.; and Gu, Q. 2018. Stochastic Gradient Descent Optimizes Over-parameterized Deep ReLU Networks. CoRR abs/1811.08888.