
The Volume of Non-Restricted Boltzmann Machines and Their Double Descent Model Complexity

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The double descent risk phenomenon has received much interest in the machine
2 learning and statistics community. Motivated through Rissanen’s minimum descrip-
3 tion length (MDL) principle, and Amari’s information geometry, we investigate
4 how a double descent-like behavior may manifest by considering the $\log V$ mod-
5 eling term - which is the logarithm of the model volume. In particular, the $\log V$
6 term will be studied for the general class of fully-observed statistical lattice models,
7 of which Boltzmann machines form a subset. Ultimately, it is found that for such
8 models the $\log V$ term can decrease with increasing model dimensionality, at a
9 rate which appears to overwhelm the classically understood $\mathcal{O}(D)$ complexity
10 terms of AIC and BIC. Our analysis aims to deepen the understanding of how the
11 double descent behavior may arise in deep lattice structures, and by extension, why
12 generalization error may not necessarily continue to grow with increasing model
13 dimensionality.

14 1 Introduction

15 Model selection is a problem which underpins the field of machine learning. Key to its formulation
16 is the notion of learning an appropriate predictor, $h^* : \mathbb{R}^D \rightarrow \mathbb{R}$ from an underlying model class,
17 \mathcal{H} , based on N input training examples $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, with each $(\mathbf{x}_i, y_i) \in \mathbb{R}^D \times \mathbb{R}$. Typically, the
18 predictor, h^* , is chosen so as to minimize some risk functional; that is, $h^* = \arg \min_{h \in \mathcal{H}} R(h)$
19 with $R(h) = \mathbb{E}_{p(x,y)}[L(h(x), y)]$, where $L : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is the risk functional, and $p(x, y)$ denotes
20 the probability density function (pdf) over the data. Fundamentally, the aim of such an approach
21 is to ensure that h^* provides good generalization capability, so that after training it minimizes the
22 *out-of-sample* test error [12]. This is historically estimated via the Akaike information criterion (AIC)
23 [3], the Bayesian information criterion (BIC) [29], or through cross validation [12]. AIC and BIC are
24 derived based on asymptotic assumptions in the sample size N , and work similarly. Moreover, both
25 criteria suggest that out-of-sample error increases as $\mathcal{O}(D)^1$, suggesting that an over-parameterized
26 model should generalize poorly, which is an idea consistent with traditional empirical evidence, via
27 the U-shaped *train-test* curves [12].

28 Recently, however, particular classes of highly parameterized models such as deep neural networks,
29 and random forests have been shown to generalize extremely well, working in contrary to the implied
30 $\mathcal{O}(D)$ model complexity effects. In fact, strong empirical evidence has been presented by Belkin et
31 al. [10], where it was shown that a *double descent risk* phenomenon may be observed for a variety of
32 models which transition into highly parameterized regimes. This phenomenon is shown in Figure
33 A2a in Appendix, with many additional experiments made clear in [25]. In an effort to explain such
34 trends Belkin et al. [11] have tried to infer some similarities between ReLU networks and traditional

¹Technically AIC has a $2D$ model complexity term, and BIC has a $D \log N$ model complexity term. We will refer to the implied effects of both model complexities simply as the “ $\mathcal{O}(D)$ terms”.

35 kernel models, and Geiger et al. [15] have tried to connect the double descent cusp-like behaviour
 36 with diverging norms, through a neural tangent kernel framework. In addition, double descent risk
 37 has been explored in a variety of simpler (and shallow) model classes [17, 6, 13], with various risk
 38 asymptotics established [9, 20]. Lastly and rather interestingly, it has been found that there exist
 39 certain parallels between double descent risk, and the notion of the *jamming transition* which occurs
 40 in physical materials which undergo a phase transition [31, 16]. In this paper we look into the problem
 41 of how increasing the underlying dimensionality of statistical lattice models (of which Boltzmann
 42 machines form a subset) could imply increasing generalizability. Such ideas are motivated from the
 43 recent empirical findings of double descent risk. This will be achieved via the notion of a *model*
 44 *volume*, which carries interpretations from information theory, and Occam’s razor. In particular it
 45 will be shown the it is possible to have a decreasing model volume in such models as D increases,
 46 which tends to offset the $\mathcal{O}(D)$ complexity terms of AIC and BIC.

47 1.1 Model Selection and Occam’s Razor

48 In the late 90s and early 2000s, extensions to the base AIC and BIC formulations were developed
 49 by Rissanen [28] and Balasubramanian [8], which include additional model-specific terms. From
 50 the perspective of coding theory, Rissanen developed a notion of *stochastic model complexity*,
 51 which builds upon Shannon’s information criteria used for lossless encoding [30]. Upon this notion,
 52 Rissanen formalized an extension of binary Shannon entropy to continuous function classes, via the
 53 discretization of the model manifold over approximately equivalent model classes. This approach
 54 establishes an intuition behind “model distinguishability”, which is also echoed by Balasubramanian.
 55 In particular, under Rissanen’s construction information is encoded in *nats* (as opposed to bits) and it
 56 is formally recognized as the *Minimum Description Length* (MDL). This is shown in Equation (1),

$$\begin{array}{c}
 \text{AIC / BIC - like term} \qquad \qquad \qquad \text{Log - Model Volume} \\
 \hline
 -\log(p(\mathbf{x})) = -\log(\hat{\mathcal{L}}) + \frac{D}{2} \log\left(\frac{N}{2\pi e}\right) + \log \int_{\Theta} \sqrt{\det(\mathcal{I}(\boldsymbol{\theta}))} d\boldsymbol{\theta} + o(1), \quad (1)
 \end{array}$$

57 where $\mathbf{x} = \{\mathbf{x}_i\}_{i=1}^N$ denotes a random vector of N data samples, $\hat{\mathcal{L}} \in \mathbb{R}$, is the likelihood function
 58 evaluated at its optimal parameter setting, with Θ being the space of possible parameter settings,
 59 and $\mathcal{I}(\boldsymbol{\theta}) \in \mathbb{R}^{D \times D}$ denotes the Fisher information matrix (FIM), which is traditionally used as
 60 a lower bound on the variance of unbiased estimators, and in Jeffrey’s prior [12]. In a parallel
 61 fashion, Balasubramanian approached the problem of model selection, albeit from a more Bayesian
 62 perspective, which yielded a very similar formulation to Rissanen’s MDL [7]. To achieve this he takes
 63 an alternate route based on specifying a Jeffrey’s prior over the underlying parameter space, which is
 64 noted to act as an appropriate measure for the density of distinguishable distributions [8]. Ultimately,
 65 in both Rissanen’s and Balasubramanian’s model selection criteria, there is a term which acts like the
 66 $\mathcal{O}(D)$ model complexity used in AIC and BIC, and an additional term: $\log \int_{\Theta} \sqrt{\det(\mathcal{I}(\boldsymbol{\theta}))} d\boldsymbol{\theta}$, which
 67 they collectively referred to as the model distinguishability-like term [19]. Moreover, since these
 68 methods are built upon the log-marginal: $-\log(p(\mathbf{x}))$, they also appeal to a Bayesian Occam’s razor-
 69 like principle. Ultimately, the ensuing study will aim to explore the term: $\log \int_{\Theta} \sqrt{\det(\mathcal{I}(\boldsymbol{\theta}))} d\boldsymbol{\theta}$, for
 70 statistical lattice models. In particular, it will be made clear that this term describes the underlying
 71 log-model volume, which we denote by $\log V$, and that this term can in fact *decrease* with increasing
 72 D in statistical lattice models. This behaviour is significant as it suggests that certain model classes
 73 can have the power to generalize well when transitioning into the over-parameterized regimes.

74 Before proceeding to analyze this volume term, it is important to clarify certain points of Equation
 75 (1). In particular, it should be noted that the $o(1)$ term is defined with respect to N , and thus it is *not*
 76 *necessarily* $o(1)$ with respect to D . In fact, it is possible for there to exist additional terms which
 77 behave as a function of D . Indeed if one approaches the derivation of an MDL-like criteria through
 78 Balasubramanian’s razor, the additional modeling term: $\log \sqrt{\det(\mathcal{I}(\hat{\boldsymbol{\theta}}))/\det(\tilde{\mathcal{I}}(\hat{\boldsymbol{\theta}}))}$ appears, where
 79 $\tilde{\mathcal{I}}$ denotes the empirical FIM, and $\hat{\boldsymbol{\theta}}$ denotes the value of $\boldsymbol{\theta}$ at the maximum likelihood location. This is
 80 often interpreted as an indicator of model *robustness* with respect to the location $\hat{\boldsymbol{\theta}}$ on the manifold [8].
 81 It is also possible to derive explicit notions of geometric curvature (in particular, the Ricci-curvature
 82 tensor) in the MDL expression if one so wishes [23] - but such expressions are not central to the
 83 present study. This is because the authors are more concerned with the holistic (i.e. geometrically
 84 intrinsic) viewpoint of the underlying model space, but these model curvature terms typically require
 85 a specific $\hat{\boldsymbol{\theta}}$ value. Moreover, the geometric model volume provides an intuitive means in which to

86 establish notions of *model distinguishability* [7, 28, 24]. In other words, we are concerned with terms
 87 which encode intrinsic *geometric complexity*, whereas $\log \sqrt{\det(\mathcal{I}(\hat{\theta}))/\det(\tilde{\mathcal{I}}(\hat{\theta}))}$ is often said to
 88 model *relative complexity* [24].

89 1.2 Information Geometry

90 Information geometry concerns the application of differential geometry to statistical models. In
 91 particular, it considers a statistical manifold, $\mathcal{M} = \{p(x; \theta)\}$, over a θ co-ordinate system. A
 92 Riemannian metric, \mathcal{G} , can be placed on \mathcal{M} , where $\mathcal{G} : \mathcal{T}_p(\mathcal{M}) \times \mathcal{T}_p(\mathcal{M}) \rightarrow \mathbb{R}_{\geq 0}$ for each point
 93 $p \in \mathcal{M}$, with $\mathcal{T}_p(\mathcal{M})$ defined as the local tangent space at point p on the manifold. Principally, \mathcal{G}
 94 is a generalization of the inner product on Euclidean spaces to Riemannian manifolds. In addition,
 95 Amari defines a dually coupled affine co-ordinate system on statistical manifolds. Dually coupled co-
 96 ordinates arise naturally from the *dually flat* property which is intrinsic to many information manifolds.
 97 These co-ordinates are known as the θ (e-flat) and η (m-flat) co-ordinates for the exponential family
 98 in particular, and are related through the Legendre transformation $\eta = \nabla\psi(\theta)$ and $\theta = \nabla\varphi(\eta)$
 99 via two convex functions $\psi, \varphi : \mathbb{R}^D \rightarrow \mathbb{R}$ [5]. The θ and η co-ordinates for exponential models
 100 correspond to the natural and expectation parameters, respectively. Furthermore, the FIM defines
 101 a natural Riemannian metric tensor: $\mathcal{G}_{ij} = \mathbb{E}[\partial_i \log(p(x; \theta)) \partial_j \log(p(x; \theta))] = \mathcal{I}_{ij}$ [26, 5]. Thus
 102 the motivation for using information geometry is clear as Rissanen’s MDL, and Balasubramanian’s
 103 Occam Razor, depend on the FIM, which is, geometrically speaking, the metric tensor. Consequently,
 104 the term $\log \int_{\Theta} \sqrt{\det(\mathcal{I}(\theta))} d\theta$ has a clear definition in differential geometry as being the *log-volume*
 105 of the underlying information manifold; that is, the square root of the determinant of the Fisher
 106 information matrix is the manifold volume [5, 21].

107 2 Statistical Lattice Models and Model Volumes

108 Statistical lattice models are popular, traditional machine learning models which include *Boltzmann*
 109 *machines* (or Ising models) [1], log-linear models [4], and the matrix balancing problem [33]. In this
 110 section, we will work the hierarchical encoding of a probability distribution via a *lattice* structure
 111 [22, 32], which will be shown to naturally lead into the η co-ordinates (*m-flat*) from information
 112 geometry (see §1.2), and analyze the learning of distributions over a lattice structured domain.

113 Formally, a *partially ordered set* (poset) is a tuple, $(\mathcal{P}, \leq_{\mathcal{P}})$, where \mathcal{P} is a set of elements, and $\leq_{\mathcal{P}}$
 114 denotes an ordering structure, such that (1) $\forall p^2 \in \mathcal{P}, p \leq_{\mathcal{P}} p$ (reflexivity), (2) If $p \leq_{\mathcal{P}} q$, and
 115 $q \leq_{\mathcal{P}} p$, then $p = q$ (antisymmetry), and (3) If $p \leq_{\mathcal{P}} q$, and $q \leq_{\mathcal{P}} r$, then $p \leq_{\mathcal{P}} r$. Note that
 116 not every element may be directly comparable to every other element in the set (which would be
 117 known as a *total* ordering). In addition, a poset $(\mathcal{P}, \leq_{\mathcal{P}})$ is called a *lattice* if every pair of elements
 118 $p, q \in \mathcal{P}$ has the least upper bound $p \vee q$ and the greatest lower bound $p \wedge q$ [14]. We assume that
 119 \mathcal{P} is finite. In working with posets it is common to consider the zeta function, $\zeta : \mathcal{P} \times \mathcal{P} \rightarrow \{0, 1\}$
 120 such that $\zeta(p, q) = \mathbf{1}_{p \leq q}$ [18]. The lattice structure always gives us the θ and η co-ordinates of a
 121 statistical manifold: $\log \mathbb{P}(p) = \sum_{q \in \mathcal{P}} \zeta(q, p) \theta_q = \sum_{q \leq p} \theta_q$ and $\eta_p = \sum_{q \in \mathcal{P}} \zeta(p, q) \mathbb{P}(q) =$
 122 $\sum_{q \geq p} \mathbb{P}(q)$ [33]. For example, for Boltzmann machines with d binary variables, the lattice space
 123 $\mathcal{P} = \{0, 1\}^n$, where $p = (p_1, \dots, p_n) \leq_{\mathcal{P}} q = (q_1, \dots, q_n)$ if $p_i \leq q_i$ for all $i \in \{1, \dots, n\}$. The
 124 size $D = |\mathcal{P}| = 2^n$ in this case. The metric tensor for the information manifold of the proposed
 125 lattice structure is shown in Theorem 1, which was previously derived by Sugiyama et al. [33]. We
 126 assume that $\mathcal{P} = \{1, \dots, |\mathcal{P}|\}$ such that 1 corresponds to the least element without loss of generality.

127 **Theorem 1 (Lattice Metric Tensor [33]).** $\mathcal{G}_{ij} = \sum_{p \in \mathcal{P}} \zeta(i, p) \zeta(j, p) \mathbb{P}(p) - \eta_i \eta_j = \eta_{i \vee j} - \eta_i \eta_j$.

128 In Theorem 1, we can replace $\sum_{p \in \mathcal{P}} \zeta(i, p) \zeta(j, p) \mathbb{P}(p)$ with $\eta_{i \vee j}$ as we assume that \mathcal{P} is a lattice and
 129 $\eta_{i \vee j}$ always exists, which says that we only consider those poset structures in which the η co-ordinate
 130 values are shared (nested) between η_i and η_j for the off-diagonal terms in the metric tensor. In
 131 Theorem 1 it is clear that this metric tensor is expressed in terms of the η co-ordinates. The equivalent
 132 metric tensor in terms of the θ co-ordinates is available, but it is much more difficult to work with
 133 (requires the Möbius function instead of the zeta function). Moreover, in this definition of the metric
 134 tensor, the first row and column are always zeros, resulting in again, a singular geometry. Luckily this

²Since elements from a poset are denoted by p , we denote a probability measure as \mathbb{P} when referencing poset systems.

135 time however, since $-\theta_1$ corresponds to the partition function, it is generally removed in practice, [34]
 136 so that we effectively work with co-ordinates $\boldsymbol{\eta}' = (\eta_2, \dots, \eta_D)$, resulting in $\mathcal{G}' \succ 0$.

137 Based on this set-up, it is possible to derive the upper and lower bounds for $\log V$. For lattice models,
 138 the η co-ordinates lie compactly on a simplex within the unit D -hypercube, that is $\Omega = [0, 1]^D$,
 139 which makes evaluations much simpler [32]. In fact, it is possible to perform the re-parameterization:
 140 $\boldsymbol{\delta} = f(\boldsymbol{\eta})$, which allows $\boldsymbol{\delta}$ to be sampled from a Dirichlet distribution. This makes the η co-ordinate
 141 more intuitive to work with, and provides us with a tractable way to evaluate the volume integral via
 142 sampling. We provide details of this re-parameterization in Appendix A.1. The $\log V$ bounds which
 143 result are shown in Theorem 2, with the proof clarified in Appendix A.2.

144 **Theorem 2 (Lattice Log Volume Bounds).** $\log V$ is bound as in Equation (2), where $\mathcal{G} = \mathcal{M}^\top \mathcal{M}$,
 145 $\boldsymbol{\delta} = f(\boldsymbol{\eta})$ is a re-parameterization, and $\Gamma(D) = (D - 1)!$ is the standard Gamma function.

$$\mathbb{E} \left[\overbrace{\sum_{i=1}^D \log(\mathcal{M}_{ii}(\boldsymbol{\delta}))}^{\text{“Richness”}} \right] + \log \left(\overbrace{\frac{1}{\Gamma(D)}}^{\text{“Distinguishability”}} \right) \leq \log V \leq \log \left(\mathbb{E} \left[\overbrace{\sqrt{\prod_{i=1}^D \mathcal{G}_{ii}(\boldsymbol{\delta})}}^{\text{“Richness”}} \right] \right) + \log \left(\overbrace{\frac{1}{\Gamma(D)}}^{\text{“Distinguishability”}} \right) \quad (2)$$

146 As Theorem 2 makes clear, the $\log V$ term can be decomposed into two components which we define
 147 as: (i) *Model richness*: which is driven by the elements found in the metric tensor, and (ii) *Model*
 148 *distinguishability*: which refers to the volume of a *probability simplex*. Intuitively, (i) For (higher-
 149 order) Boltzmann machines, multi-way interactions can be encoded in a desired lattice structure, and
 150 this in turn will drive the construction of the metric tensor via the η co-ordinate system (Theorem
 151 1). Thus, since the metric tensor depends strongly upon the chosen lattice, and since it is used in
 152 defining angles and geodesics over a manifold, it would appear that this expression encodes a notion
 153 of *model richness* and or expressibility over the manifold. As for point (ii), since the η co-ordinate
 154 system is constrained to lie on a simplex geometry (as made explicit by the poset structuring) the
 155 volume is: $1/\Gamma(D) = 1/(D - 1)!$. Evidently, as dimensionality increases the simplex volume
 156 decreases. A nice combinatorial intuition of this is that for D randomly sampled numbers $\{n_i\}_i^D$, the
 157 probability of obtaining a permutation which is precisely the total ordering of these D points (i.e.
 158 $n_1 < n_2 < \dots < n_D$) is $1/D!$. Thus, even if the AIC $\mathcal{O}(D)$ model complexity term grows without
 159 bound, the simplex constraint over the D parameters serves to act as a strong counter-balance. In
 160 fact, this counter-balancing imbues the MDL expression with a *double descent*-like behaviour. As D
 161 increases, the value of the MDL expression first increases, and then decreases. Since MDL relates
 162 strongly to the notion of model generalizability the double descent phenomenon which arises here
 163 paints an intriguing picture for the double descent risk phenomenon that has empirically arisen in
 164 the deep learning field. A visual example of this on toy values is made clear in Figure A2b. It is
 165 important to note that our analysis assumes access to a *fully observable* lattice model (and thus is
 166 not immediately applicable to restricted Boltzmann machines without careful consideration). This is
 167 because in latent hierarchical settings the underlying geometry can become geometrically singular,
 168 which complicates the global model volume calculation with respect to the FIM [36? , 35].

169 In regards to (ii) model distinguishability, it should be noted that the term *distinguishability* is
 170 somewhat overloaded here, since in classic MDL literature it has been traditionally reserved for the
 171 entire expression: $\int_{\Theta} \sqrt{\det(\mathcal{I}(\boldsymbol{\theta}))} d\boldsymbol{\theta}$ [8, 28]. Our proposal is to elaborate it slightly by splitting this
 172 term into two terms, where one term works to clarify the importance of the model architecture, and
 173 the other term clarifies the constraints that exist in the underlying parameter space.

174 Inspecting these terms in the limit of the over-parameterized regime results in the following (see
 175 Appendix A.3 for its proof).

176 **Remark 1 (Limiting Lattice Volume).** $\lim_{D \rightarrow \infty} V = 0$.

177 Thus, the volume of the proposed statistical lattice models tends towards zero for large D , and this
 178 limit can be said to converge factorially due to the simplex volume. Finally, it is interesting (and
 179 tempting) to relate the volume calculations established here to the model volumes of tree structures,
 180 particularly since these models have similar graphical topologies to statistical lattices. Doing so, one
 181 is able to establish Remark 2, in which it can be seen that two topologically similar models do indeed
 182 consider similar volume expressions, and that in both expressions $\lim_{D \rightarrow \infty} V = 0$ holds.

183 **Remark 2 (Rissanen’s Tree Volume [28]).** *Binary decision trees for encoding D -classes have*
 184 *log-volume: $\log_2 V = D/2 \log_2 \pi - \log_2 \Gamma(D/2)$. Ignoring higher order terms, statistical lattice*
 185 *structures have their log-volume upper-bounded as: $\log_e V \leq D/2 \log_e \pi - \log_e \Gamma(D/2)$.*

References

- 186
- 187 [1] David H. Ackley, Geoffrey E. Hinton, and Terrence J. Sejnowski. A learning algorithm for
188 Boltzmann machines. *Cognitive Science*, 9(1):147–169, 1985.
- 189 [2] Alan Agresti. *Categorical Data Analysis*. Wiley, 3 edition, 2012.
- 190 [3] Hirotogu Akaike. Information theory and an extension of the maximum likelihood principle.
191 In B. N. Petrov and F. Caski, editors, *Proceedings of the 2nd International Symposium on*
192 *Information Theory*, pages 267–281, 1973.
- 193 [4] Shun-ichi Amari. Information geometry on hierarchy of probability distributions. *IEEE*
194 *Transactions on Information Theory*, 47(5):1701–1711, 2001.
- 195 [5] Shun-ichi Amari and Hiroshi Nagaoka. *Methods of information geometry*, volume 191. Ameri-
196 can Mathematical Soc., 2007.
- 197 [6] Jimmy Ba, Murat Erdogdu, Taiji Suzuki, Denny Wu, and Tianzong Zhang. Generalization of
198 two-layer neural networks: An asymptotic viewpoint. In *International Conference on Learning*
199 *Representations*, 2020.
- 200 [7] Vijay Balasubramanian. A geometric formulation of occam’s razor for inference of parametric
201 distributions. *arXiv:adap-org/9601001*, 1996.
- 202 [8] Vijay Balasubramanian. MDL, Bayesian inference, and the geometry of the space of probability
203 distributions. *Advances in Minimum Description Length: Theory and Applications*, pages
204 81–98, 2005.
- 205 [9] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in
206 linear regression. *arXiv:1906.11300*, 2019.
- 207 [10] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-
208 learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy*
209 *of Sciences*, 116(32):15849–15854, 2019.
- 210 [11] Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to
211 understand kernel learning. *arXiv:1802.01396*, 2018.
- 212 [12] Christopher M Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- 213 [13] Stéphane d’Ascoli, Maria Refinetti, Giulio Biroli, and Florent Krzakala. Double trouble in
214 double descent: Bias and variance(s) in the lazy regime. *arXiv:2003.01054*, 2020.
- 215 [14] Brian A. Davey and Hilary A. Priestley. *Introduction to Lattices and Order*. Cambridge
216 University Press, 2 edition, 2002.
- 217 [15] Mario Geiger, Arthur Jacot, Stefano Spigler, Franck Gabriel, Levent Sagun, Stéphane d’Ascoli,
218 Giulio Biroli, Clément Hongler, and Matthieu Wyart. Scaling description of generalization
219 with number of parameters in deep learning. *Journal of Statistical Mechanics: Theory and*
220 *Experiment*, 2020(2):023401, 2020.
- 221 [16] Mario Geiger, Stefano Spigler, Stéphane d’Ascoli, Levent Sagun, Marco Baity-Jesi, Giulio
222 Biroli, and Matthieu Wyart. Jamming transition as a paradigm to understand the loss landscape
223 of deep neural networks. *Physical Review E*, 100(1):012115, 2019.
- 224 [17] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized
225 two-layers neural networks in high dimension. *arXiv:1904.12191*, 2019.
- 226 [18] Gerhard Gierz, Karl H. Hofmann, Klaus Keimel, Jimmie D. Lawson, Michael Mislove, and
227 Dana S. Scott. *Continuous Lattices and Domains*. Cambridge University Press, 2003.
- 228 [19] Peter D Grünwald. *The minimum description length principle*. MIT press, 2007.
- 229 [20] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-
230 dimensional ridgeless least squares interpolation. *arXiv:1903.08560*, 2019.

- 231 [21] Harold Jeffreys. An invariant form for the prior probability in estimation problems. *Proceedings*
232 *of the Royal Society of London. Series A. Mathematical and Physical Sciences*, 186(1007):453–
233 461, 1946.
- 234 [22] Simon Luo and Mahito Sugiyama. Bias-variance trade-off in hierarchical probabilistic models
235 using higher-order feature interactions. In *Proceedings of the AAAI Conference on Artificial*
236 *Intelligence*, volume 33, pages 4488–4495, 2019.
- 237 [23] Bruno Mera, Paulo Mateus, and Alexandra M Carvalho. On the minmax regret for statistical
238 manifolds: the role of curvature. *arXiv preprint arXiv:2007.02904*, 2020.
- 239 [24] In Jae Myung, Vijay Balasubramanian, and Mark A Pitt. Counting probability distributions:
240 Differential geometry and model selection. *Proceedings of the National Academy of Sciences*,
241 97(21):11170–11175, 2000.
- 242 [25] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever.
243 Deep double descent: Where bigger models and more data hurt. In *International Conference on*
244 *Learning Representations*, 2020.
- 245 [26] C Radhakrishna Rao. Information and the accuracy attainable in the estimation of statistical
246 parameters. In *Breakthroughs in Statistics*, pages 235–247. Springer, 1992.
- 247 [27] Jorma J Rissanen. Fisher information and stochastic complexity. *IEEE Transactions on*
248 *Information Theory*, 42(1):40–47, 1996.
- 249 [28] Jorma J Rissanen. Stochastic complexity in learning. *Journal of Computer and System Sciences*,
250 55(1):89–95, 1997.
- 251 [29] Gideon Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464,
252 1978.
- 253 [30] Claude E Shannon. A mathematical theory of communication. *Bell System Technical Journal*,
254 27(3):379–423, 1948.
- 255 [31] Stefano Spigler, Mario Geiger, Stéphane d’Ascoli, Levent Sagun, Giulio Biroli, and Matthieu
256 Wyart. A jamming transition from under-to over-parametrization affects loss landscape and
257 generalization. *arXiv:1810.09665*, 2018.
- 258 [32] Mahito Sugiyama, Hiroyuki Nakahara, and Koji Tsuda. Information decomposition on struc-
259 tured space. In *2016 IEEE International Symposium on Information Theory*, pages 575–579,
260 2016.
- 261 [33] Mahito Sugiyama, Hiroyuki Nakahara, and Koji Tsuda. Tensor balancing on statistical manifold.
262 In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages
263 3270–3279, 2017.
- 264 [34] Mahito Sugiyama, Hiroyuki Nakahara, and Koji Tsuda. Legendre decomposition for tensors. In
265 *Advances in Neural Information Processing Systems*, pages 8811–8821, 2018.
- 266 [35] Ke Sun and Frank Nielsen. Lightlike neuromanifolds, occam’s razor and deep learning.
267 *arXiv:1905.11027*, 2019.
- 268 [36] Sumio Watanabe. *Algebraic Geometry and Statistical Learning Theory*, volume 25. Cambridge
269 University Press, 2009.

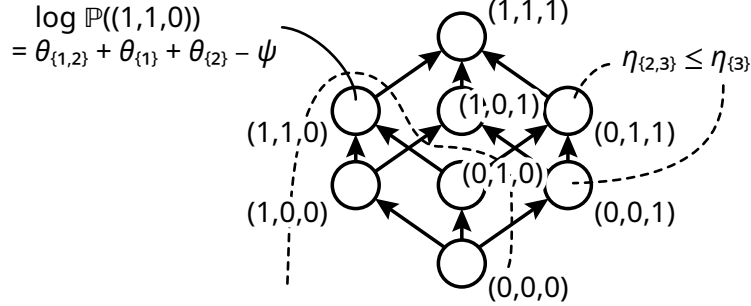


Figure A1: An example of a lattice structure for the binary domain $\{0, 1\}^3$. Each arrow denotes the partial order between elements in the lattice.

270 Appendix

271 A Statistical Lattice Models

272 A.1 Reparameterizing the Dual Geometry of Lattices

273 Lattices are useful structures, as they allow one to efficiently encode information hierarchically. A
 274 geometry over lattice-type structures based on modelling higher order feature interactions via log
 275 probabilities has been derived in the work of Sugiyama et al. [32]. The well-known log-linear model
 276 for binary variables in question is formulated as,

$$\log \mathbb{P}(\mathbf{x}) = \sum_i \theta_{\{i\}} x_i + \sum_{i < j} \theta_{\{i,j\}} x_i x_j + \sum_{i < j < k} \theta_{\{i,j,k\}} x_i x_j x_k + \dots + \theta_{\{1,\dots,n\}} x_1 \dots x_n - \psi,$$

277 where $\mathbf{x} \in \{0, 1\}^n$, each $\theta \in \mathbb{R}$ denotes the connection strength of a particular higher order
 278 interaction, each $x_i \in \{0, 1\}$ denotes a binary valued variable which activates a particular connection
 279 strength, $\psi \in \mathbb{R}$ denotes the normalization constant for the probability model [4, 2]. Under this
 280 structure, $\mathbb{P}(\mathbf{x})$ is a member of the exponential family of distributions. If we define a particular
 281 instance of the partial ordering as $\mathbf{x} = (x_1, \dots, x_n) \leq \mathbf{y} = (y_1, \dots, y_n)$, where $x_i \leq y_i$ for all
 282 $i \in \{1, \dots, n\}$, and denote by $\Sigma(\mathbf{x})$ as the set of indices of “1” in \mathbf{x} , then when can instantiate a
 283 lattice, and can condense the representation of the above log-linear model as:

$$\log \mathbb{P}(\mathbf{x}) = \sum_{\mathbf{s}} \delta(\mathbf{s}, \mathbf{x}) \theta_{\Sigma(\mathbf{s})} = \sum_{\mathbf{s} \leq \mathbf{x}} \theta_{\Sigma(\mathbf{s})}, \quad \text{where } \psi = -\theta_{\emptyset}.$$

284 Hence the lattice is a natural representation of this hierarchical structure over the sample space of
 285 $\{0, 1\}^n$. Sugiyama et al. [32] studied geometric structure of statistical lattice models and showed that
 286 distributions over not only $\{0, 1\}^n$ but any lattices belong to the exponential family. Note that posets
 287 are originally used in [32], which is a more general structure than lattices. Although we treat only
 288 lattices in this paper, most interesting statistical models (such as Boltzmann machines) are lattices.
 289 We thus proceed in this direction as lattice structures entail a simple co-ordinate representation of the
 290 metric tensor as we have described in Theorem 1. An example of a lattice structure for $\{0, 1\}^3$ is
 291 illustrated in Figure A1.

292 As Amari notes [5], the exponential family of distributions induces a statistical manifold which
 293 possesses an interesting dualistic structure. That is, two co-ordinate systems can be dually connected
 294 and allow one to generalize notions such as the Pythagoras theorem in Euclidean manifolds, to more
 295 general statistical manifolds. In particular, for the exponential family the first of these co-ordinate
 296 systems is given by $\boldsymbol{\theta} = (\theta_1, \dots, \theta_D)$ (as defined in the specified log-linear model of this subsection),
 297 and the second is given by $\boldsymbol{\eta} = (\eta_1, \dots, \eta_D)$. Note that in the case of a binary log-linear model, we
 298 have $D = 2^n$ and

$$\begin{aligned} \eta_{\{i\}} &= \mathbb{E}[x_i] = \Pr(x_i = 1) \\ \eta_{\{i,j\}} &= \mathbb{E}[x_i x_j] = \Pr(x_i = 1, x_j = 1) \\ \eta_{\{1,\dots,n\}} &= \mathbb{E}[x_1 \dots x_n] = \Pr(x_1 = 1, \dots, x_n = 1), \end{aligned}$$

299 and $(\boldsymbol{\theta}, \boldsymbol{\eta})$ are explicitly dually connected via the Legendre transformation [5, 32]. Note that since
 300 the η co-ordinate system is defined probabilistically, and is thus constrained to be in $[0, 1]^D$, which
 301 is convenient, and simplifies many calculations. Moreover, note that the η co-ordinate system is
 302 built hierarchically, in the sense that it adheres to a partial ordering similar to the following *example*
 303 structure:

$$\eta_1 \geq \eta_2, \eta_3 \quad \eta_2 \geq \eta_4, \quad \eta_3 \geq \eta_5, \quad \dots$$

304 This ordering is model specific and always uniquely determined from the lattice structure, and
 305 it is thus difficult to perform integrations over such arbitrary orderings. However, owing to the
 306 probabilistic nature of the η co-ordinate system, it is possible to impose the following recursive
 307 re-parameterisation:

$$\begin{aligned} \eta_1 &= \delta_1 \\ \eta_2 &= \eta_1 + \delta_2 \\ \eta_3 &= \eta_1 + \delta_3 \\ &\vdots \\ \eta_D &= \eta_{D-1} + \delta_D \end{aligned}$$

308 where each $\delta_i \in [0, 1]$ for $1 \leq i \leq D$, and $\sum_{i=1}^D \delta_i = 1$. Geometrically speaking, $\boldsymbol{\delta} = \{\delta_i\}_{i=1}^D$
 309 represents points in a D -simplex. We can then proceed to formally encode the lattice ordering
 310 constraints via an additional zeta matrix, resulting in $\boldsymbol{\eta} = \mathcal{Z}\boldsymbol{\delta}$, where each $\mathcal{Z}_{ij} \in \{0, 1\}$ is the value
 311 of zeta function $\zeta(q_i, q_j) = \mathbf{1}_{q_i \leq q_j}$ for the corresponding elements q_i and q_j in the lattice. In other
 312 words, points from the D simplex, $\boldsymbol{\delta}$, can be transformed into $\boldsymbol{\eta}$ co-ordinates for the poset manifold
 313 through a linear mapping. In order to re-express the volume integral via the $\boldsymbol{\delta}$ co-ordinates, it
 314 is necessary to calculate the determinant of the Jacobian transformation matrix between the co-ordinate
 315 systems. However this is trivially one, since \mathcal{Z} is by construction upper triangular, resulting in
 316 $\det\left(\frac{\partial \boldsymbol{\eta}}{\partial \boldsymbol{\delta}}\right) = 1$. Therefore it is possible to calculate the log-volume integral as,

$$\log\left(\int_{\Delta_D} \sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))} \cdot \det\left(\frac{\partial \boldsymbol{\eta}}{\partial \boldsymbol{\delta}}\right)^2 d\boldsymbol{\delta}\right) = \log\left(\int_{\Delta_D} \sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))} d\boldsymbol{\delta}\right), \quad (3)$$

317 where the coordinate transformation was performed using the square of the determinant, as the
 318 metric tensor is rank (0,2) - that is, it is a doubly covariant object, and we define the D -simplex as
 319 Δ_D . In this form it is natural to re-express the volume integral via the expectation operator, where
 320 the expectation is taken with respect to a Dirichlet distribution over $\boldsymbol{\delta}$, as the Dirichlet distribution
 321 represents a pdf over the probability simplex. In other words we consider,

$$\begin{aligned} \log\left(\int_{\Delta_D} \sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))} d\boldsymbol{\delta}\right) &= \log\left(\int_{\Delta_D} \sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))} \cdot \frac{w(\boldsymbol{\delta})}{w(\boldsymbol{\delta})} d\boldsymbol{\delta}\right) \\ &= \log\left(\mathbb{E}\left[\frac{\sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))}}{w(\boldsymbol{\delta})}\right]\right), \end{aligned} \quad (4)$$

322 where $w(\boldsymbol{\delta}) = \frac{\prod_{i=1}^D \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^D \alpha_i)} \prod_{i=1}^D x_i^{\alpha_i-1} := \text{Dir}(\boldsymbol{\delta}; \boldsymbol{\alpha})$, with $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_D)$, and $\Gamma : \mathbb{R} \rightarrow \mathbb{R}$
 323 being the standard Gamma function. Here, the choice of $\boldsymbol{\alpha}$ controls the manner in which sampling
 324 is performed over Δ_D . We opt for a uniform exploration over the D -simplex, which equates to
 325 requiring that $\alpha_d = 1$ for all $d \in D$. Doing so means that we get, $w(\boldsymbol{\delta}) = \frac{\prod_{i=1}^D \Gamma(1)}{\Gamma(\sum_{i=1}^D 1)} = \frac{1}{\Gamma(D)}$, where
 326 $\Gamma(D) = (D-1)!$. Ultimately, Equation (4) becomes,

$$\log\left(\mathbb{E}\left[\frac{\sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))}}{w(\boldsymbol{\delta})}\right]\right) = \log\left(\mathbb{E}\left[\sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))}\right]\right) - \log \Gamma(D), \quad (5)$$

327 implying that bounding the volume can be equivalently achieved by bounding the behaviour of
 328 $\log\left(\mathbb{E}\left[\sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))}\right]\right)$, and then appending the $\log \Gamma(D)$ term. The upper and lower bounds on this
 329 volume integral, are shown in Appendix A.2.

330 **A.2 Proof of Theorem 2**

331 In this subsection we proceed to find lower and upper bounds for the $\log V$ in the case of the
 332 prescribed lattice geometry, by exploiting the re-parameterization of the η co-ordinate system.

333 *Proof. Volume Upper Bound:*

334 From Hadamarad's inequality: $|\det(A)| \leq \prod_{i=1}^D A_{ii}$, for $A \in \mathbb{R}^{D \times D}$. Thus:

$$\begin{aligned}
 |\det(\mathcal{G}(\boldsymbol{\delta}))| &= \det(\mathcal{G}(\boldsymbol{\delta})) \quad (\mathcal{G}(\boldsymbol{\delta}) \succ 0) \\
 &\leq \prod_{i=1}^D \mathcal{G}_{ii}(\boldsymbol{\delta}), \\
 \Rightarrow \quad \sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))} &\leq \sqrt{\prod_{i=1}^D \mathcal{G}_{ii}(\boldsymbol{\delta})} \\
 \Leftrightarrow \quad \mathbb{E} \left[\sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))} \right] &\leq \mathbb{E} \left[\sqrt{\prod_{i=1}^D \mathcal{G}_{ii}(\boldsymbol{\delta})} \right] \\
 \Leftrightarrow \quad \log \left(\mathbb{E} \left[\sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))} \right] \right) &\leq \log \left(\mathbb{E} \left[\sqrt{\prod_{i=1}^D \mathcal{G}_{ii}(\boldsymbol{\delta})} \right] \right) \\
 \Leftrightarrow \quad \log \left(\mathbb{E} \left[\sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))} \right] \right) - \log \Gamma(D) &\leq \log \left(\mathbb{E} \left[\sqrt{\prod_{i=1}^D \mathcal{G}_{ii}(\boldsymbol{\delta})} \right] \right) - \log \Gamma(D) \\
 \Leftrightarrow \quad \log V &\leq \log \left(\frac{\mathbb{E} \left[\sqrt{\prod_{i=1}^D \mathcal{G}_{ii}(\boldsymbol{\delta})} \right]}{\Gamma(D)} \right),
 \end{aligned}$$

335

336

337 **Volume Lower Bound:**

$$\log \left(\mathbb{E} \left[\sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))} \right] \right) \geq \mathbb{E} \left[\log \left(\sqrt{\det(\mathcal{G}(\boldsymbol{\delta}))} \right) \right] \quad (6)$$

$$= \frac{1}{2} \mathbb{E} \left[\log \circ \det(\mathcal{G}(\boldsymbol{\delta})) \right], \quad (7)$$

338 where Inequality (6) is Jensen's inequality. Moreover, $\mathcal{G} \succ 0 \Rightarrow \exists \mathcal{M}$ s.t. $\mathcal{G} = \mathcal{M}\mathcal{M}^\top$, where \mathcal{M} is a
 339 triangular matrix (Cholesky decomposition). Thus,

$$\begin{aligned}
 \det(\mathcal{G}) &= \det(\mathcal{M}\mathcal{M}^\top) \\
 &= \det(\mathcal{M}) \cdot \det(\mathcal{M}^\top) \\
 &= \det(\mathcal{M}) \cdot \det(\mathcal{M}) \\
 &= \det(\mathcal{M})^2 \\
 &= \left(\prod_{i=1}^D \mathcal{M}_{ii} \right)^2, \quad (8)
 \end{aligned}$$

340

$$\begin{aligned}
&\Rightarrow \frac{1}{2} \mathbb{E} [\log \circ \det (\mathcal{G}(\boldsymbol{\delta}))] = \frac{1}{2} \mathbb{E} \left[\log \left(\prod_{i=1}^D \mathcal{M}_{ii}(\boldsymbol{\delta}) \right)^2 \right] \\
&= \mathbb{E} \left[\sum_{i=1}^D \log (\mathcal{M}_{ii}(\boldsymbol{\delta})) \right], \\
&\Rightarrow \log \left(\mathbb{E} \left[\sqrt{\det (\mathcal{G}(\boldsymbol{\delta}))} \right] \right) \geq \mathbb{E} \left[\sum_{i=1}^D \log (\mathcal{M}_{ii}(\boldsymbol{\delta})) \right] \\
&\Leftrightarrow \log \left(\mathbb{E} \left[\sqrt{\det (\mathcal{G}(\boldsymbol{\delta}))} \right] \right) - \log \Gamma(D) \geq \mathbb{E} \left[\sum_{i=1}^D \log (\mathcal{M}_{ii}(\boldsymbol{\delta})) \right] - \log \Gamma(D) \\
&\Leftrightarrow \log V \geq \mathbb{E} \left[\sum_{i=1}^D \log (\mathcal{M}_{ii}(\boldsymbol{\delta})) \right] + \log \left(\frac{1}{\Gamma(D)} \right)
\end{aligned}$$

341

□

342 A.3 Proof of Remark 1

343 Here we show that as $D \rightarrow \infty$, $V \rightarrow 0$ which implies that $\log V \rightarrow -\infty$. This is an important
344 indicator in the increase of generalization performance, as sufficiently large D can therefore overpower
345 the $\mathcal{O}(D)$ model complexity term, present in traditional AIC and BIC.

346 *Proof.* As this represents a volume integral a trivial lower bound is zero. It follows that

$$\begin{aligned}
0 \leq V &\leq \mathbb{E} \left[\frac{\sqrt{\prod_{i=1}^D \mathcal{G}_{ii}(\boldsymbol{\delta})}}{\Gamma(D)} \right] \\
&\Rightarrow \lim_{D \rightarrow \infty} 0 \leq \lim_{D \rightarrow \infty} V \leq \lim_{D \rightarrow \infty} \mathbb{E} \left[\frac{\sqrt{\prod_{i=1}^D \mathcal{G}_{ii}(\boldsymbol{\delta})}}{(D-1)!} \right].
\end{aligned}$$

347 Since each $\mathcal{G}_{ii}(\boldsymbol{\delta}) \in [0, 1]$, the factorial in the denominator strongly dominates, so that:

$$\lim_{D \rightarrow \infty} \mathbb{E} \left[\frac{\sqrt{\prod_{i=1}^D \mathcal{G}_{ii}(\boldsymbol{\delta})}}{(D-1)!} \right] = 0.$$

348 Thus from via an application of squeeze theorem we see that,

$$\begin{aligned}
0 &\leq \lim_{D \rightarrow \infty} V \leq 0, \\
&\Rightarrow \lim_{D \rightarrow \infty} V = 0.
\end{aligned}$$

349

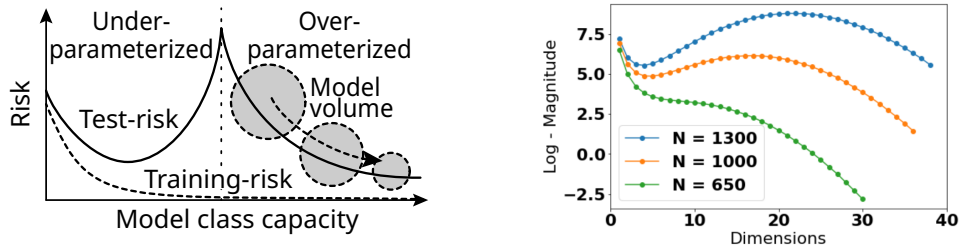
□

350 A.4 Visual Summary of the Effect of $\log V$

351 A visual summary of the effect of the $\log V$ term in the modeling behaviour for test is made clear
352 in Figure A2b. This plot has been produced with respect to toy data, where the behaviour of the
353 log likelihood has been modeled wlog as: $\log(10^3/D^{10})$ to model rapidly increasing log likelihood
354 as D increases. The $\mathcal{O}(D)$ term has been modeled in accordance to BIC as: $D/2 \log N$, and the
355 value of N in Figure A2b has been changed accordingly (for three different cases). Lastly, the $\log V$
356 value has been determined by constructing the metric tensor from first principles and evaluating the
357 volume integral from first Monte Carlo sampling. Firstly, it is clear that the presence of $\log V$ term
358 seems to drive a generalization behaviour that bears striking resemblance to the double-descent risk

359 curves. Similar intuitions were shown by Sun & Nielsen [35] for deep neural networks, which is
 360 shown in Figure A3. However, in their analysis it was necessary to consider an additional term in
 361 the MDL expansion, and instead of decreasing volume, they report having an *increasing* volume
 362 term. Thus it would appear that these higher order terms in the MDL-like expansion are highly
 363 model specific, and are required to be examined on a case-by-case basis. Indeed a blanket statement
 364 such as: *the log-volume term drives a double descent behavior* is not correct, but it can be said that
 365 for fully-observed statistical lattice models this appears to be the case. Regardless, it would appear
 366 that for a more complete understanding of the peculiarities of the *modern ML regime*, a geometric
 367 perspective seems to be certainly invaluable.

368 Secondly, in Figure A2b the double descent peak is observed to shift to the right with increasing N .
 369 Information theoretically, increasing N proliferates the total number of possible encodings which
 370 may be able to explain the observed data. This is an interpretation which is consistent with Rissanen’s
 371 original derivation of MDL, in which he states: “the number of distinguishable models grows with the
 372 length of the data, which seems reasonable. In view of this we define the model complexity (as seen
 373 through the data)” [27]. Interestingly, this can imply that when a model’s total distinguishability is
 374 insufficient (weak regularization, and or insufficient noise), it is possible for the model to generalize
 375 well on one quantity of data, N_1 , and then upon re-training on some new data such that $N_2 > N_1$, to
 376 then generalize poorly, due to the ability of the double-descent cusp to shift towards the right. Similar
 377 ideas have been uttered recently by Nakkiran et al., in that: “for a fixed architecture and training
 378 procedure, more data (can) actually hurt” [25]. Ultimately, such analysis seem to suggest that for a
 379 more holistic understanding on the double-descent phenomenon a geometric approach provides a
 380 more complete understanding, with invaluable intuition and assistance.



(a) Double descent risk proposed behavior for lattice-type statistical models.

(b) Three different data cases on MDL evaluations (y-values).

Figure A2: Comparing double descent risk shapes to MDL calculations for statistical lattice machines (of which Boltzmann machines form a subset).

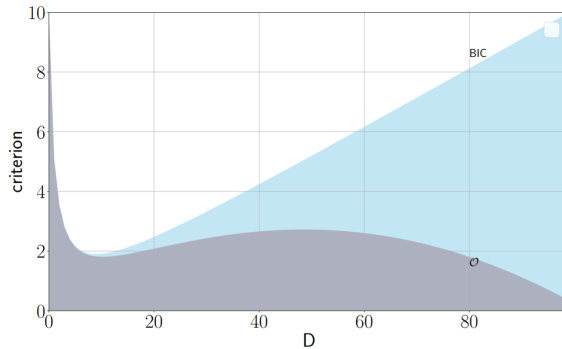


Figure A3: Deep learning razor based on inspecting the geometric structure of MDL-like expressions in Sun & Nielsen [35].

381 A.5 Proof of Remark 2

382 In Rissanen’s MDL derivation for decision trees [28] (and finite alphabet processes), he derives
 383 a log-volume term as follows, $\log V = \log \frac{\pi^{D/2}}{\Gamma(D/2)}$. Since there is a topological similarity in the

384 structure between lattice models and tree structures (not necessarily one-to-one), it is tempting
 385 to determine under what conditions Risannen's volume expression can be derived from the poset
 386 formulations shown here. Note that in Risannen's tree structure there is no inequality symbol needed
 387 for his log-volume expression. This is because the FIM for his model is diagonal, and so there is no
 388 need to invoke Hadamarad's inequality to make the volume integral tractable (which is done here).

389 *Proof.* Consider an upper-bound on $\log V$, established through Hadamarad's inequality of determi-
 390 nants: $\sqrt{\det(\mathcal{G}(\eta))} \leq \sqrt{\prod_{i=1}^D \mathcal{G}(\eta)_{ii}} = \sqrt{\prod_{i=1}^D \eta_i(1 - \eta_i)}$, which results because each diagonal
 391 term in our formulation, $\mathcal{G}(\eta)_{ii} = \eta_i(1 - \eta_i)$. Thus,

$$\begin{aligned} V &= \int_{\Delta_{D-1}} \sqrt{\det(\mathcal{G}(\eta))} d\eta \\ &\leq \int_{\Delta_{D-1}} \sqrt{\prod_{i=1}^D \mathcal{G}(\eta)_{ii}} d\eta \\ &= \int_{\Delta_{D-1}} \sqrt{\prod_{i=1}^{D-1} \eta_i(1 - \eta_i)} d\eta \end{aligned}$$

392 where we denote the D-1 dimensional simplex via Δ_{D-1} . Our aim here is to transform the expression
 393 $\sqrt{\prod_{i=1}^{D-1} \eta_i(1 - \eta_i)}$ into a more amenable term, from which the above integral can be evaluated.
 394 This can be established by expanding and collecting polynomial terms. Consider,

$$\begin{aligned} \prod_{i=1}^{D-1} \eta_i(1 - \eta_i) &= \prod_{i=1}^{D-1} \eta_i(1 - \eta_1)(1 - \eta_2)\dots(1 - \eta_{D-1}) \\ &= \prod_{i=1}^{D-1} \eta_i(1 - \eta_1 - \eta_2 + \eta_1\eta_2)(1 - \eta_3)\dots(1 - \eta_{D-1}) \\ &= \prod_{i=1}^{D-1} \eta_i \left(1 - \sum_{i=1}^{D-1} \eta_i + \sum_{i \neq j} \eta_i\eta_j - \sum_{i \neq j \neq k} \eta_i\eta_j\eta_k + \dots + (-1)^D \prod_{i=1}^{D-1} \eta_i \right) \\ &= \left(\prod_{i=1}^{D-1} \eta_i \right) \left(1 - \sum_{i=1}^{D-1} \eta_i \right) + \left(\prod_{i=1}^{D-1} \eta_i \right) \left(\sum_{i \neq j} \eta_i\eta_j - \sum_{i \neq j \neq k} \eta_i\eta_j\eta_k + \dots + (-1)^D \prod_{i=1}^{D-1} \eta_i \right) \end{aligned}$$

395 where the above is a result of expanding the polynomial terms and collecting like expressions into
 396 their respective groups, and expanding the brackets. Therefore the integrand can be considered to
 397 evaluate as:

$$\begin{aligned} \sqrt{\prod_{i=1}^D \mathcal{G}(\eta)_{ii}} &= \sqrt{\prod_{i=1}^{D-1} \eta_i(1 - \eta_i)} \\ &= \sqrt{\left(\prod_{i=1}^{D-1} \eta_i \right) \left(1 - \sum_{i=1}^{D-1} \eta_i \right) + \left(\prod_{i=1}^{D-1} \eta_i \right) \left(\sum_{i \neq j} \eta_i\eta_j - \sum_{i \neq j \neq k} \eta_i\eta_j\eta_k + \dots + (-1)^D \prod_{i=1}^{D-1} \eta_i \right)} \\ &= \sqrt{\left(\prod_{i=1}^{D-1} \eta_i \right) \left(1 - \sum_{i=1}^{D-1} \eta_i \right) + \mathcal{O}(\eta^3)} \\ &\approx \sqrt{\left(\prod_{i=1}^{D-1} \eta_i \right) \left(1 - \sum_{i=1}^{D-1} \eta_i \right)} \end{aligned}$$

398 where the approximation results from observing that each $\eta_i \in [0, 1]$. In particular, high polynomial
 399 orders of η_i must quickly approach zero, with the speed of approach increasing as the order becomes
 400 higher. This would only not be the case if many $\eta_i = 1$. However, such a system would imply a trivial
 401 poset structure, since our partial ordering needs $\eta_i > \eta_j$ for $i > j$ for non-trivial lattice structures. In
 402 other words, since the geometric co-ordinates for the lattice structure are expressed through η , for
 403 non-trivial co-ordinate space it is required for the majority of $\eta \in (0, 1)$.

404 The purpose of re-writing the integrand in this way, was to motivate the usage of a Type I Dirichlet
 405 Integral for its evaluation. In particular, such integrals subscribe to the following evaluation rule:

$$\int_{\Delta_D} \prod_{i=1}^D x_i^{\alpha_i-1} f\left(\sum_{i=1}^D x_i\right) d^D x_i = \frac{\prod_{i=1}^D \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^D \alpha_i)} \int_0^1 f(\tau) \tau^{\sum_{i=1}^D \alpha_i - 1} d\tau.$$

406 In order to evaluation our volume integral we must consider, (i) Defining the function $f : y \mapsto \sqrt{1-y}$,
 407 where $y := \sum_{i=1}^{D-1} \eta_i$ and (ii) Each $\alpha_i = 3/2$. These reasons are made explicit as follows:

$$\begin{aligned} \int_{\Delta_{D-1}} \prod_{i=1}^{D-1} x_i^{\alpha_i-1} f\left(\sum_{i=1}^{D-1} x_i\right) d^{D-1} x_i &= \int_{\Delta_{D-1}} \prod_{i=1}^{D-1} \eta_i^{1.5-1} f\left(\sum_{i=1}^{D-1} \eta_i\right) d\eta \\ &= \int_{\Delta_{D-1}} \prod_{i=1}^{D-1} \eta_i^{0.5} \left(1 - \sum_{i=1}^{D-1} \eta_i\right)^{0.5} d\eta \\ &= \int_{\Delta_{D-1}} \sqrt{\prod_{i=1}^{D-1} \eta_i (1 - \eta_i)} d\eta \end{aligned}$$

408 where evidently it can be seen that $\alpha_i = 3/2$, and $f\left(\sum_{i=1}^{D-1} \eta_i\right) = \sqrt{\left(1 - \sum_{i=1}^{D-1} \eta_i\right)}$. Therefore,
 409 we can evaluate the volume integral upper-bound as follows:

$$\int_{\Delta_{D-1}} \sqrt{\prod_{i=1}^{D-1} \eta_i (1 - \eta_i)} d\eta = \frac{\prod_{i=1}^{D-1} \Gamma\left(\frac{3}{2}\right)}{\Gamma\left(\sum_{i=1}^{D-1} \frac{3}{2}\right)} \int_0^1 (1-\tau)^{1/2} \tau^{\sum_{i=1}^{D-1} (3/2)-1} d\tau \quad (9)$$

410 From this evaluation, it can be noticed that the new integral on the RHS is in fact an Euler-Beta
 411 integral, which evaluates according to the following rule:

$$\int_0^1 (1-x)^{\beta-1} x^{\alpha-1} dx = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)},$$

412 where $\beta = 3/2$ and $\alpha = \sum_{i=1}^{D-1} (3/2)$ in this case. Thus:

$$\begin{aligned} \int_{\Delta_{D-1}} \sqrt{\prod_{i=1}^{D-1} \eta_i (1 - \eta_i)} d\eta &= \frac{\prod_{i=1}^{D-1} \Gamma\left(\frac{3}{2}\right)}{\Gamma\left(\sum_{i=1}^{D-1} \frac{3}{2}\right)} \cdot \frac{\Gamma\left(\frac{3}{2}\right) \Gamma\left(\sum_{i=1}^{D-1} \frac{3}{2}\right)}{\Gamma\left(\sum_{i=1}^{D-1} \frac{3}{2} + \frac{3}{2}\right)} \\ &= \frac{\prod_{i=1}^D \Gamma\left(\frac{3}{2}\right)}{\Gamma\left(\sum_{i=1}^D \frac{3}{2}\right)}. \end{aligned}$$

413 From the Legendre relation for the Γ function, the following recurrence relation holds: $\Gamma(x)\Gamma(x +$
 414 $1/2) = 2^{1-2x} \sqrt{\pi} \Gamma(2x)$. This allows one to evaluation $\Gamma(3/2)$ as follows:

$$\begin{aligned}\Gamma(1)\Gamma(3/2) &= 2^{-1}\sqrt{\pi}\Gamma(2) \\ \Rightarrow \Gamma(3/2) &= 2^{-1}\sqrt{\pi}.\end{aligned}$$

415 Moreover, since by definition of the Gamma function we can write $\Gamma\left(\sum_{i=1}^D \frac{3}{2}\right) = \Gamma\left(D + \frac{D}{2}\right) =$
 416 $\Gamma\left(\frac{D}{2}\right)(D-1)!$, it can therefore be established that:

$$\frac{\prod_{i=1}^D \Gamma\left(\frac{3}{2}\right)}{\Gamma\left(\sum_{i=1}^D \frac{3}{2}\right)} = \frac{\pi^{D/2}}{\Gamma(D/2)} \cdot \frac{1}{2^D(D-1)!},$$

417 Establishing now the $\log V$ term one can arrive at the following conclusions:

$$\begin{aligned}\log V &\leq \log\left(\frac{\pi^{D/2}}{\Gamma(D/2)}\right) + \log\left(\frac{1}{2^D(D-1)!}\right) \\ &= \log\left(\frac{\pi^{D/2}}{\Gamma(D/2)}\right) - \log(2^D\Gamma(D)) \\ &\leq \log\left(\frac{\pi^{D/2}}{\Gamma(D/2)}\right)\end{aligned}$$

418 Evidently, the second expression $\frac{1}{2^D\Gamma(D)}$ goes to zero much faster than $\frac{\pi^{D/2}}{\Gamma(D/2)}$, especially for large
 419 D . Thus, as before we take the dominant (first) term in this approximation and arrive at:

$$\log V \leq \log\left(\frac{\pi^{D/2}}{\Gamma(D/2)}\right)$$

420 as required. □