Why do deep neural networks generalise in the overparameterised regime?

Ard Louis



Will answer questions in the chat

Guillermo Valle Perez



Learning machines?

Instead of trying to produce a programme to simulate the adult mind, why not rather try to produce one which simulates the child's ? If this were then subjected to an appropriate course of education one would obtain the adult brain

We have thus divided our problem into two parts. The child-programme and the education process.

Alan Turing, *Computing Machinery and Intelligence*, Mind **59**, 433 (1950)



Hyper-astronomical numbers and Boolean choices

Doctor's decision table for COVID-19

Send to hospital?	Fever?	Cough?	Lost sense of smell?	Body ache?	Recent travel to hotspot?	Over 50?	Heart problem?	Obese?	Diabetes?
	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	0	1	1	1
	1	1	1	1	1	0	0	0	0
	1	1	1	0	0	1	0	1	1
	1	1	0	0	0	1	0	1	1
			Can we	learn the r	est of the f	unction?			

n questions; 2^n possible answers; 2^{2^n} possible Boolean functions

For n=9 $2^9 = 512$ answers; $2^{512} = 1.34 \times 10^{154}$ possible functions

(10⁸⁶ elementary particles in observable universe).

Hyper-astronomical numbers in history



Mari 08613 tablet (Old Babylonian. 1900-1600 BC) A barley-corn: to a single barley-corn I increased, 2 barley-corns in the 1st day; 4 barley-corns in the 2nd day; 8 barley-corns in the 3rd day;

" "

"

30 2 'thousand' 7 'hundred' 37 talents 1/2 mina 2 1/3 shekels 4 barley-corns in the 30th day.

(47 tons of Barley)

Jöran Friberg (2005), Unexpected Links Between Egyptian and Babylonian Mathematics,

Hyper-astronomical numbers in history



Invetion of Chess: Shah-nama "The Book of Kings" by Ferdowsi, (940 - 1020)

Barley corn and chess board story recounted by Ibn Khallikan c.a. 1254.

 2^{64} -I = 18,446,744,073,709,551,615 (18 quintillion) > 1000 times current annual wheat production

Hyper-astronomical numbers in biology



Proteins: 20^L

All

- L=37 proteins weigh more than earth
- L=58 proteins weigh more than the visible universe
- L=476 proteins: 10⁵⁰⁰ times the mass of the visible universe



RNA: 4^L

All

- L=55 RNA strands weigh 10¹⁰ kg
- L=79 RNA strands weigh more than the Earth
- L=126 RNA weigh more than the visible universe

AAL, <u>Contingency, convergence and hyper-astronomical numbers in biological evolution</u> Studies in History and Philosophy of Biological and Biomedical Sciences 58, 107 (2016)

Evolving the same structure again and again

Question: Convergence: Natural selection or the arrival of variation?



Hammerhead ribozyme L \approx 55, so all 4^L combinations weigh 10¹⁰ kg Experiment is on ~ 0.01g of material, about ~ 1/10¹⁵ fraction of the space

Salehi-Ashtiani K, Szostak JW. In vitro evolution suggests multiple origins for the hammerhead ribozyme. Nature. 414, 82 (2001).

RNA Genotype \rightarrow Phenotype map

GAAAGUCUGGGCUAAGCCACUGAUGGUGUCUGAAAUGAGAGGAAAACUUUUG



Tertiary structure (3D)

Secondary structure (who bonds to whom)

Model G \rightarrow P map: RNA secondary structures

 $N_G = 4^L$



L=15: $4^{15} \sim 10^9$ sequences; these map to 431 phenotypes (note big bias!)



Hammerhead ribosome is anomalously many sequences mapping to it (large neutral set)



0.001 = fraction of all possible structures take up > 50% of genotypes

 10^{33} sequences $\rightarrow 10^{13}$ structures.

Green: Set of 504 sequences from functional RNA database

Kamaludin Dingle, Steffen Schaper, and Ard A. Louis, Interface Focus 5 (6), 20150053 (2015)



Why the strong bias?

INTUITION:

What is the probability that a monkey types out X digits of $\mathbf{\pi}$ on an N key typewriter ?



But what if the monkey types into C ?

 $P(X) = (1/N)^{(X+1)}$

3.14159265358979323846264338327950288419716939 937510582097494459230781640628620899862803482 534211706798214808651328230664709384460955058 223172535940812848111745028410270193852110555 964462294895493038196442

 $\mathsf{P}(\mathsf{X}) \lesssim (1/\mathsf{N})^{133}$

133 character (obfuscated) C code to calculate first 15,000 digits of π

```
a[52514],b,c=52514,d,e,f=1e4,g,h;
main(){for(;b=c-=14;h=printf("%04d", e+d/f))
for(e=d%=f;g=--b*2;d/=g)d=d*b+f*(h?a[b]:f/5),a[b]=d%--g;}
```

 $\pi = \sum_{i=0}^{\infty} \frac{(i!)^2 2^{i+1}}{(2i+1)!}$

C program due to Dik Winter and Achim Flammenkamp (See Unbounded Spigot Algorithms for the Digits of Pi, by Jeremy Gibbons (Oxford CS), Math. Monthly, April 2006, pages 318-328.)

Formalizing the monkey intuition with AIT: Kolmogorov complexity





A.N. Kolgomorov 1903-1987

G.J. Chaitin 1947--

Kolmogorov/Chaitin complexity K(X) is the length in bits of the shortest program on a UTM that generates X

K is universal, (not UTM dependent) because you can always write a compiler => O(1) terms.

$$K_U(X) = K_W(X) + O(1) \approx K(X)$$
 asymptotically

K is not computable due to Halting problem.

Warning: you don't know for sure that it is complex, t could be encoding $\pi = 3.141592653589793238462 \dots$

new intutions

- -- A random number is one for which $\mathrm{K}(\mathrm{X}) \gtrsim |\mathrm{X}|$
- -- The complexity of a set can be << than complexity of elements of the set

Formalizing the monkey intuition with AIT: Algorithmic Probability



R. Solomonoff 1926-2009



Intuitively: simpler (small K(X)) outputs are much more likely to appear

It seems to me that the most important discovery since Gödel was the discovery by Chaitin, Solomonoff and Kolmogorov of the concept called Algorithmic Probability,. Everybody should learn all about that and spend the rest of their lives working on it. Marvin Minsky (2014) https://www.youtube.com/watch?v=DfY-DRsE86s&feature=youtu.be&t=1h30m02s

Solomonoff, R., "A Preliminary Report on a General Theory of Inductive Inference", Report V-131, Zator Co., Cambridge, Ma. Feb 4, 1960, revision, Nov., 1960.

Formalizing the monkey intuition with AIT: The Coding Theorem



 $2^{-K(x)} \le P(x) \le 2^{-K(x) + O(1)}$

We should teach this much more widely!

L. Levin, 1948 --

Intuitively: simpler (small K(X)) outputs are much more likely to appear

Serious problems for applying coding theorem

- 1) Many systems are not Universal Turing Machines
- 2) Kolmogorov complexity K(x) is formally incomputable
- 3) Only holds in the asymptotic limit of large x...

L.A. Levin. Laws of information conservation (non-growth) and aspects of the foundation of probability theory. Problems of Information Transmission, 10:206–210, 1974.

AIT coding (like) theorem for non UTM maps

Proof sketch:

- 1) For simple maps **f**, with input size **n** we can calculate the whole set of input \rightarrow output at O(1) cost (complexity of a set << elements of set)
- 2) Encode this with a Shannon-Fano-Elias (SFE) code for which $P(x) \sim \frac{1}{2}$ length -O(1)
- 3) This procedure gives a bound on the Kolmogorov complexity given f and n

K. Dingle, C. Camargo and A.AL, Nature Communications 9,761 (2018)

AIT coding (like) theorem for input-output maps



(2 Dphils of work)

$$P(x) \lesssim 2^{-a\tilde{K}(x)-b}$$



Kamal Dingle

NOTE: upper bound only!

Chico Camargo

- 1) Computable input-output map f: $I \rightarrow O$
- 2) Map **f** must be simple e.g. K(f) grows slowly with system size -- then $K(x|f,n) \approx K(x) + O(1)$
- K(x) is approximated, for example by Lempel Ziv compression or some other suitable measure.
- 2) Constants a and b depend on mapping only and can be approximated fairly easily.
- 3) Bound is tight for most inputs, but not most outputs.
- 4) Maps must be a) simple, b) redundant, c) non-linear, d) well-behaved (e.g. not a pseudorandom number generator)
- 5) There is also a statistical lower bound.

RNA map shows simplicity bias

Frequent



More rare

Simplicity Bias: $P(x) \lesssim 2^{-a\tilde{K}(x)-b}$

Is black line (red dashed with b=0)



K. Dingle, C. Camargo and A.AL, Nature Communications 9,761 (2018)



Can simplicity bias be applied to supervised learning with deep neural networks?

Neural networks are highly expressive

We have thus divided our problem into two parts. The child-programme and the education process.

- A Turing (1950)





Universal approximation theorem for NN

Neural networks are fundamentally function approximators. The following theorem holds:

For any Lebesgue-integrable function $f: \mathbb{R}^n \to \mathbb{R}$ and any $\epsilon > 0$, there exists a fullyconnected ReLU network \mathcal{A} with width $d_m \leq n+4$, such that the function $F_{\mathcal{A}}$ represented by this network satisfies

$$\int_{\mathbb{R}^n} |f(x) - F_\mathcal{A}(x)| \,\mathrm{d} x < \epsilon$$

Neural networks are highly expressive -

B. Hanin Approximating Continuous Functions by ReLU Nets of Minimal Width. arXiv preprint arXiv:1710.11278.

Neural networks are typically highly over-parameterized: number of parameters >> number of data points



5 parameters

With four parameters I can fit an elephant, and with five I can make him wiggle his trunk -- John von Neuman (according to Fermi)

> F. Dyson, *A meeting with Enrico Fermi*. Nature. 427, 287 (2004)



1923-2020

Drawing an elephant with four complex parameters Jürgen Mayer; Khaled Khairy; Jonathon Howard; *American Journal of Physics* **78,** 648-649 (2010) Neural networks are typically highly over-parameterized: number of parameters >> number of data points



Why do the DNNs not over-fit?

Comparison of a polynomial fit to a DNN fit (with thousands of parameters)

$$y(x) = a_0 + a_1 x + a_2 x^2 + a_3 x^3 + \dots a_n x^n$$

Bias-Variance tradeoff in Machine Learning

Why does this concept not work well for DNNs?



FIG. 5 Bias-Variance tradeoff and model complexity.

From A high-bias, low-variance introduction to Machine Learning for physicists

Pankaj Mehta, Ching-Hao Wang, Alexandre G. R. Day, and Clint Richardson

Arxiv/1803.08823

Conundrum: if DNNs are highly expressive, why do they pick functions that generalize so well?

airplane	1	-X-	-	X	*	+	2	-1	-	-
automobile					-	Tet	-		1-0	*
bird	S	ſ	2		(Constanting	-	1		2	4
cat	2		-	50		12	E.	Å,	the second	2
deer	6	48	X	RA		Y	Y	1	1	
dog	19%.	C	-	% .	1			V?	A	1.
frog	-7	a.	-		23		No.	57		5-
horse	-ph	T.	1	2	P	TAB	-	2	G	1
ship	-		dinte	-	<u>M</u> A		Ż	18	1	
truck			1							in the

C. Zhang et al., **Understanding deep learning requires rethinking generalization**. arXiv:1611.03530 (2016) Showed that you could randomise the labels, and still easily

If a DNN can "memorize" a dataset, why does it pick functions that generalise so well?

train to zero training error.

CIFAR-10 dataset

Supervised learning of a Boolean function with DNNs

Doctor's decision table for COVID-19

	Send to hospital?	Fever?	Cough?	Lost sense of smell?	Body ache?	Recent travel to hotspot?	Over 50?	Heart problem?	Obese?	Diabetes?
tion	1	1	1	1	1	1	1	1	1	1
Inct	1	1	1	1	1	1	0	1	1	1
Boolean tu	0	1	1	1	1	1	0	0	0	0
	1	1	1	1	0	0	1	0	1	1
	0	1	1	0	0	0	1	0	1	1

Given some examples, can we learn the rest of the function?

n questions; 2^n possible answers; 2^{2^n} possible Boolean functions

For n=9 $2^9 = 512$ answers; $2^{512} = 1.34 \times 10^{154}$ possible functions

 $(10^{86} \text{ elementary particles in observable universe})$.

Parameter-function map for deep learning



Let the space of functions that the model can express be \mathcal{F} . If the model has p real valued parameters, taking values within a set $\Theta \subseteq \mathbb{R}^p$,

the parameter-function map, \mathcal{M} , is defined as:

$$\mathcal{M}: \Theta o \mathcal{F} \ heta \mapsto f_ heta$$

where f_{θ} is the function implemented by the model with choice of parameter vector θ .

G.Valle-Perez, C. Camargo and A.A. Louis, arxiv: 1805.08522 – ICLR 2019

A-Priori probability: If we randomly sample parameters θ , how likely are we to produce a particular function f?





Chris Mingard

Theorem 4.1. For a perceptron f_{θ} with b = 0 and weights w sampled from a distribution which is symmetric under reflections along the coordinate axes, the probability measure $P(\theta : \mathcal{T}(f_{\theta}) = t)$ is given by

 $P(\theta: \mathcal{T}(f_{\theta}) = t) = \begin{cases} 2^{-n} & \text{if } 0 \le t < 2^n \\ 0 & \text{otherwise} \end{cases}.$

We also prove that this bias towards simple function gets stronger with more layers

Neural networks are a priori biased towards Boolean functions with low entropy, Chris Mingard, Joar Skalse, Guillermo Valle-Pérez, David Martínez-Rubio, Vladimir Mikulik, Ard A. Louis arxiv: 1909.11522

P(f): If we randomly sample parameters θ , how likely are we to produce a particular function f?

Model problem for a 7 bit string, study all Boolean functions f. There are $2^7 = 128$ different strings, and $2^{128} \simeq 10^{38}$ different functions. You might expect a 10^{-38} chance of finding any function. Instead, we find strong simplicity bias.





10⁸ samples of parameters for (7,40,40,1) vanilla fully connected DNN system (FCN) with ReLU.

G.Valle Perez, C. Camargo and A.A. Louis, arxiv: 1805.08522 – ICLR 2019

Does simplicity bias help generalisation?



but less well on complex functions

Trained on 64 of 128 possible functions.

Problem; DNNs are not trained by randomly sampling parameters



DNNs are trained using Stochastic gradient descent (SGD) on a loss function.

Problem; DNNs are not trained by randomly sampling parameters

Intuition: for very strong bias: Basin of attraction \sim Basin size (P(f))



Chris

Mingard

Bayesian picture of bias for training set S

Prior over functions P(f)

If we wish to infer (i.e. no noise) at some points, then we need a 0-1 likelihood on training data $S = \{(x_i, y_i)\}_{i=1}^{m}$

$$P(S|f) = \begin{cases} 1 \text{ if } \forall i, \ f(x_i) = y_i \\ 0 \text{ otherwise }. \end{cases}$$

Posterior follows from Bayes rule

$$P(f|S) = \frac{P(S|f)P(f)}{P(S)},$$

Where the marginal likelihood or evidence is

Functions that fit S

$$P(S) = \sum_{f} P(S|f) P(f) = \sum_{f \in C(S)} P(f)$$

For fixed S, bias in P(f|S) translates over from bias in prior P(f)

SGD acts (almost) like a Bayesian sampler



FCN on binarized MNIST – training set 10,000, test set $10,000 - 2^{100} = 10^{30}$ possible functions fit the test set.

SGD acts (almost) like a Bayesian sampler



(b) vanilla CNN

Other data sets and architectures

SGD acts (almost) like a Bayesian sampler



Figure 2: Effects of hyperparameter changes: All examples are for a 2-layer 1024 node FCN trained on MNIST with CE loss. (a) shows RMSprop which for batch size = 128 achieves an error on

Training set size dependence









10-20

 10^{-40}

0.6 0.8 1.0

 $\langle \epsilon_G \rangle = 3.33\%$













(h) Experiment 3. 20000 training examples

4

6

 10^{0}

5

Test set dependence



Upshot so far:

- The probabilities that functions are found by SGD is remarkably close to a Bayesian sampler.
- Exact reason not understood, but heuristically this must stem from large bias. (for unbiased systems the correlations are very weak)
- Nice, because we can use Bayesian picture to understand generalisation.
- Also, we observe very strong bias (consistent with simplicity bias)

Two questions in generalisation



- Why do DNNs generalise at all in the overparameterised regime?
- 2) Given DNNs that generalise, can we further fine-tune the hyperparameters to improve generalisation?

Two questions in generalisation



Why do DNNs generalise at all in the overparameterised regime?

Can we derive generalisation bounds ?

PAC-Bayes bounds

Corollary 1. (*Realizable PAC-Bayes theorem (for Bayesian classifier)*) Under the same setting as in Theorem 1, with the extra assumption that D is realizable, we have:

$$-\ln\left(1-\epsilon(Q^*)\right) \le \frac{\ln\frac{1}{P(U)} + \ln\left(\frac{2m}{\delta}\right)}{m-1}$$

where $Q^*(c) = \frac{P(c)}{\sum_{c \in U} P(c)}$, U is the set of concepts in \mathcal{H} consistent with the sample S, and where $P(U) = \sum_{c \in U} P(c)$

Adapted from David McAllester

P(U) – Marginal likelihood we saw before.

G. Valle Perez, C. Camargo and A.A. Louis, arxiv: 1805.08522 - ICLR 2019

Guillermo Valle Perez

Tight PAC-Bayes bounds:



(a) for a 4 hidden layers convolutional network

(b) for a 1 hidden layer fully connected network

Corollary 1. (*Realizable PAC-Bayes theorem (for Bayesian classifier)*) Under the same setting as in Theorem 1, with the extra assumption that D is realizable, we have:

$$-\ln\left(1-\epsilon(Q^*)\right) \le \frac{\ln\frac{1}{P(U)} + \ln\left(\frac{2m}{\delta}\right)}{m-1}$$

where $Q^*(c) = \frac{P(c)}{\sum_{c \in U} P(c)}$, U is the set of concepts in \mathcal{H} consistent with the sample S, and where $P(U) = \sum_{c \in U} P(c)$ G.Valle Perez, C. Camargo and A.A. Louis, arxiv: 1805.08522 – ICLR 2019



Guillermo Valle Perez

Tight PAC-Bayes bounds for scaling of error with training set size m (learning curves)



Guillermo Valle Perez



Observed: error ~ m^{-α}
1) α decreases with data complexity (bad news for machine learning)
2) α appears independent of algorithm
3) We can reproduce this scaling with PAC-Bayes theory approach we have derived.

But, WHY this scaling?

Corollary 1. (*Realizable PAC-Bayes theorem (for Bayesian classifier)*) Under the same setting as in Theorem 1, with the extra assumption that D is realizable, we have:

$$-\ln\left(1-\epsilon(Q^*)\right) \le \frac{\ln\frac{1}{P(U)} + \ln\left(\frac{2m}{\delta}\right)}{m-1}$$

where $Q^*(c) = \frac{P(c)}{\sum_{c \in U} P(c)}$, U is the set of concepts in \mathcal{H} consistent with the sample S, and where $P(U) = \sum_{c \in U} P(c)$

Bayesian predictions

Averaged over datasets

$$\langle P(f|\mathcal{D}^m)\rangle = P(f) \left\langle \frac{P(\mathcal{D}_i^m|f)}{P(\mathcal{D}_i^m)} \right\rangle_{\mathcal{D}_i^m} \approx \frac{P(f) \left(1 - \epsilon(f)\right)^m}{\langle P(\mathcal{D}^m) \rangle} = t_F(f,m) P(f) \approx \frac{P(f) e^{-m\epsilon(f)}}{\langle P(\mathcal{D}^m) \rangle}$$

Training factor tells you how much the function probability is affected by training

$$t_F(f,m) = \frac{(1-\epsilon(f))^m}{\langle P(\mathcal{D}^m) \rangle} \approx \frac{e^{-m\epsilon(f)}}{\langle P(\mathcal{D}^m) \rangle} = \frac{e^{-m\epsilon_G}}{\langle P(\mathcal{D}^m) \rangle} e^{-m(\epsilon(f)-\epsilon_G)}$$

Error distribution on data $(1-\epsilon(f))^m$ Log(P(f) 1 $(1-\varepsilon(f))$ 0.5 $K_{17}(f)$ $\left\langle P(f|\mathcal{D}^m)\right\rangle = P(f) \left\langle \frac{P(\mathcal{D}_i^m|f)}{P(\mathcal{D}_i^m)} \right\rangle_{\mathcal{D}_i^m} \approx \frac{P(f) \left(1 - \epsilon(f)\right)^m}{\langle P(\mathcal{D}^m) \rangle} = t_F(f,m) P(f) \approx \frac{P(f) e^{-m\epsilon(f)}}{\langle P(\mathcal{D}^m) \rangle}$

Chaotic regime (tanh nonlinearity)



(a) Tanh

(**b**) *ReLU*

Figure 3: Mean field phase diagrams for tanh and ReLU activation functions showing various phase regimes as a function of σ_w and σ_b .



(a) Increasing σ_w has no significant effect on the bias in a simple ReLU DNN. $\sigma_b = 0.2$. The DNN has 2 hidden layers each containing 40 ReLU neurons.



(a) Increasing σ_w decreases the bias in a simple DNN. $\sigma_b = 0.2$. The DNN has 2 hidden layers of 40 tanh neurons.



(b) In the chaotic regime, increasing the number of layers decreases the bias in a simple DNN. $\sigma_b = 0.2$ and $\sigma_w = 2$. Each hidden layer contains 40 tanh neurons.

Greg Yang and Hadi Salman. A fine-grained spectral perspective on neural networks. arXiv preprint arXiv:1907.10599, 2019.

Increasing chaos decreases bias

Henry Rees



J. Empirical probability versus LZ complexity plots



K. Generalisation versus LZ complexity plots

Figure 17: Empirical probability of individual functions versus their LZ complexity for networks initialised with various σ_w and numbers of layers. Despite suffering from finite-size effects, points with a probability of 10^{-8} are not removed since in plots ($\sigma_w = 4$, d = 10) and ($\sigma_w = 8$, d = 10) only points of this type are found. Details are the same as Figure 6.

Figure 19: Effect of increasing σ_w and network depth on the generalising ability of Boolean networks trained on functions $\{0,1\}^7 \rightarrow \{0,1\}$. Details are the same as Figure 7.

Bayesian Approach

Henry Rees

$$\langle P(f|\mathcal{D}_m)\rangle = P(f)\left\langle \frac{P(\mathcal{D}_m|f)}{P(\mathcal{D}_m)}\right\rangle \approx \frac{P(f)\left(1-\epsilon(f)\right)^m}{\langle P(\mathcal{D}_m)\rangle}$$

 $(1 - \varepsilon(f)) = 1 - 2/128 = 0.984...$

Bayesian Approach -

Henry Rees







Bayesian Approach

$$\langle P(f|\mathcal{D}_m)\rangle = P(f)\left\langle \frac{P(\mathcal{D}_m|f)}{P(\mathcal{D}_m)}\right\rangle \approx \frac{P(f)\left(1-\epsilon(f)\right)^m}{\langle P(\mathcal{D}_m)\rangle}$$





Bayesian Approach

Henry Rees

 $P(K|D^{64})$



Experiments with SGD

- Train on a Specific Target Function with training set of size 64
- Calculate the complexity and Generalisation Error of function implemented by model
- Repeat many times randomising the examples in the training set each time.



Target LZ = 66.5Target LZ = 66.5 $\sigma_w \mid N^o$ Layers 0.14 0.14 $\sigma_w \mid N^o$ Layers Normalised Probability 1 | 10 | 10 0.12 0.12 8 | 10 8 | 10 0.10 0.10 P(K|D)0.08 0.08 0.06 0.06 0.04 0.04 0.02 0.02 0.00 0.00 50 100 150 50 75 125 150 0 25 100 0 Lempel-Ziv Complexity Output Lempel-Ziv Complexity

Comparison of SGD to Bayesian

Comparison of SGD to Bayesian



Bayesian Approach

90 80

 $\phi \mid N^o$ Layers

1.2

tanh | 2

tanh | 5

tanh | 10

 $\phi \mid N^o$ Layers

1.2

1.0

tanh | 2

tanh | 5

tanh | 10

1.4

1.6

1.8

Weight Standard Deviation σ_w Training Set Size = 10000

2.0

2.2

2.4

1.4

1.6

80

70

Henry Rees

1.8 2.0 2.2 2.4 Weight Standard Deviation σ_w Training Set Size = 1000



MNIST

THANKYOU

Conclusions

- Deep learning may work because they have a natural bias towards simple functions (Occam's razor)
- We show evidence that SGD behaves like a Bayesian optimiser
- PAC-Bayes provides tight bounds
- A Bayesian approach can help explain how generalisation works
- Work done by Guillermo Valle Perez, Chris Mingard, Joar Skalse, Henry Rees and more