

Using Learning Dynamics to Understand Neural Network Generalisation

Arnu Pretorius*

Supervisor: Dr. Steve Kroon (Computer Science*)

Co-supervisor: Dr. Herman Kamper (EE Engineering⁺)

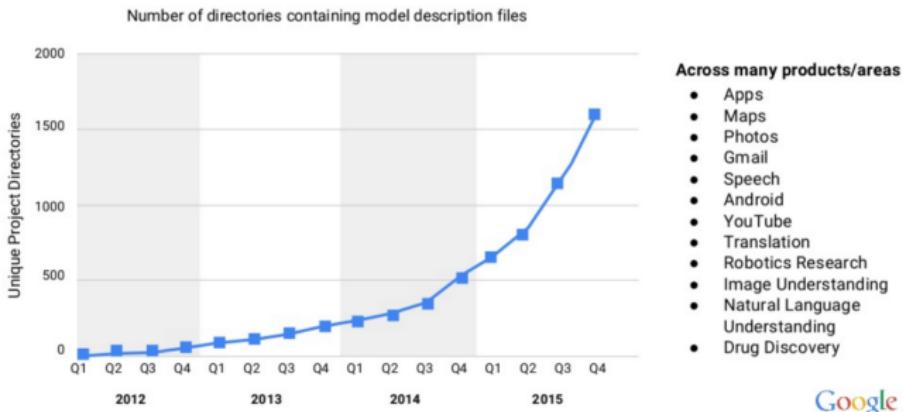
20 October 2017

Maties Machine Learning, Stellenbosch University



The success of Deep Learning

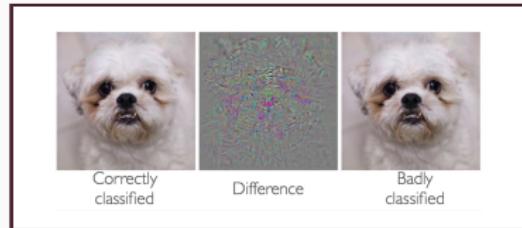
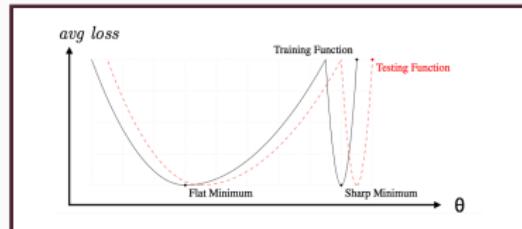
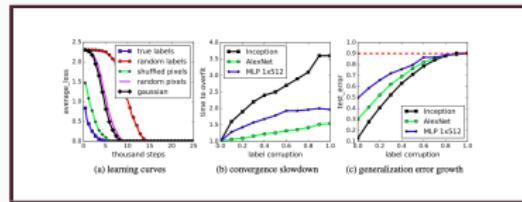
Growing Use of Deep Learning at Google



Scale of deep learning deployment at Google over time.

The mysteries of Deep Learning

- Generalisation
- The role of explicit versus implicit regularisation
- Non-convex optimisation
- Adversarial examples
- and more ...



Understanding deep learning requires rethinking generalization, Zhang, Bengio, Hardt, Recht, Vinyals. ICLR, 2017.

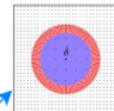
On large-batch training for deep learning: generalization gap and sharp minima, Keskar, Mudigere, Smelyanskiy, Tang. ICLR, 2017.

Intriguing Properties of Neural Networks, Szegedy, Zaremba, Sutskever, Bruna, Erhan, Goodfellow, Fergus. ICLR, 2014.

Approaches Towards Understanding Neural Networks



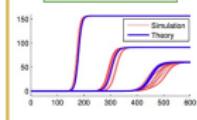
Initialisation [1]



Loss Surface [2]

Understanding Neural Networks

Learning [3]



Complexity [4]



Approaches towards understanding neural networks.

[1] *Understanding the difficulty of training deep feedforward neural networks*, Glorot, Bengio. AISTATS, 2010.

[2] *Sharp minima can generalize for deep nets*, Dinh, Pascanu, Bengio, Bengio. arXiv 2017.

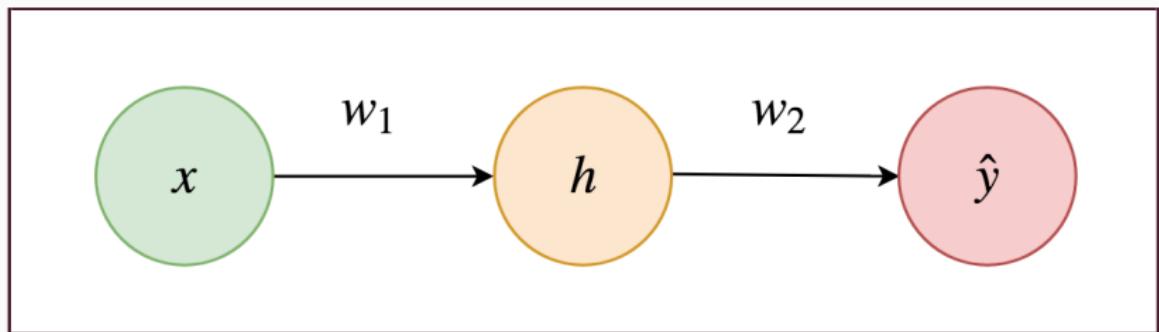
[3] *Exact solutions to the nonlinear dynamics of learning in deep linear neural networks*, Saxe, McClelland, Ganguli. ICLR, 2014.

[4] *On the Expressive Power of Deep Neural Networks*, Raghu, Poole, Kleinberg, Ganguli, Dickstein. ICML, 2017.

Scalar linear neural network

Let $w_2, w_1, x \in \mathbb{R}$,

$$\hat{y} = w_2 w_1 x \quad (1)$$



Scalar linear neural network.

Scalar learning dynamics

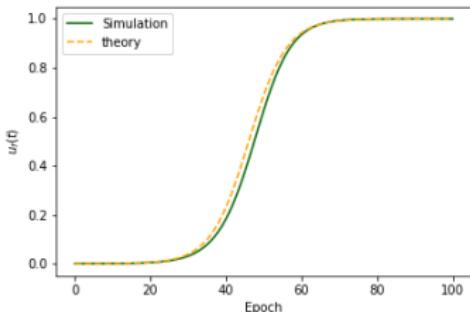
Loss surface $(1 - w_2 w_1)^2$, GD path (orange), learning dynamics (green).

Exact solutions for scalar learning dynamics

The learning dynamics for a scalar linear neural network starting from small initial values $u_0 \equiv w_{2(0)}w_{1(0)}$ is given by

$$u_f(t) = \frac{yE}{x(E-1) + y/u_0}, \quad (2)$$

where $E = e^{2yxt\alpha}$, α is the learning rate and t is measured in epochs.



Simulated versus theoretical learning dynamics.

Learning dynamics for deep linear networks

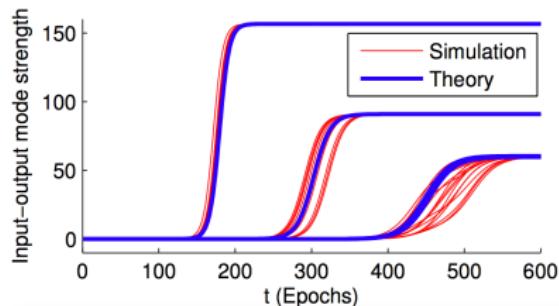
$$\sum^{31} = U S V^T$$

Properties Modes Items

P B S F M 1 2 3 C S O R

Input-output correlation matrix Feature synthesizer vectors Singular values Object analyzer vectors

The diagram shows the decomposition of a 5x5 input-output correlation matrix into three components: feature synthesizer vectors (U), singular values (S), and object analyzer vectors (V^T). The matrix has columns labeled P, B, S, F, M and rows labeled C, S, O, R. The singular values S are represented by three red squares of decreasing size. A color scale bar indicates values from -1 (blue) to 1 (red). The feature synthesizer vectors U and object analyzer vectors V^T are shown as 5x3 matrices where non-zero entries are colored according to the singular value magnitudes.



Linear neural network learning dynamics.

Our current focus

- The role of regularisation

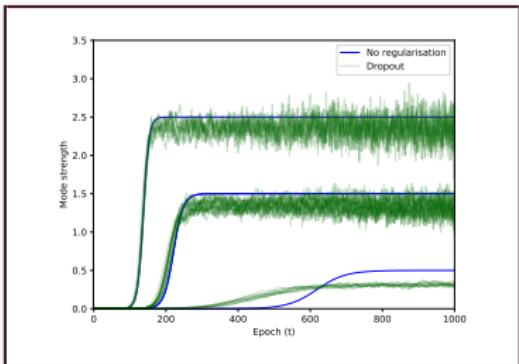
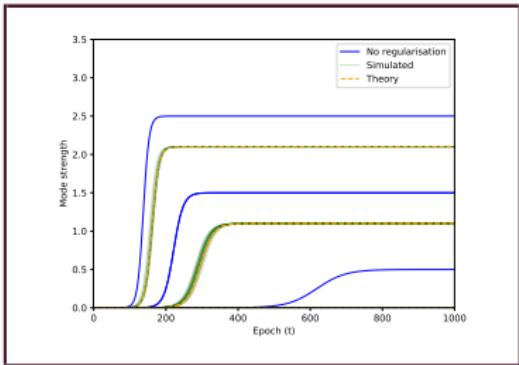
- Weight decay

$$u(t, s, u_0, \lambda^*) = \frac{(s - \lambda^*)e^{2(s - \lambda^*)t/\tau}}{e^{2(s - \lambda^*)t/\tau} - 1 + (s - \lambda^*)/u_0}$$

- Dropout

$$E_{d \sim \text{Bern}(\rho)}[t] \leq \frac{\tau}{sp} \ln \left(\frac{s}{\epsilon} \right) = \mathcal{O} \left(\frac{\tau}{sp} \right)$$

- Autoencoder networks
- Generalisation



Summary

- Deep learning has been hugely successful in solving large and complex machine learning task, however many mysteries remain.
- Better understanding deep neural networks might be achieved via several different routes.
- Studying the learning dynamics of neural networks may help us understand how neural networks learn.
- We hope to use this learning dynamics approach to study generalisation in deep neural networks.

Thank you for listening!