

Mathematics of AI

Optimisation in deep learning

Goodfellow et al

Connections to previous weeks

- Simon in Week 1 + me in Week 3: “ReLU works good!”
 - Let’s see it in action: <https://udlbook.github.io/udlfigures/>
- Matt in Week 2: “Eigenvalues are super important!”
 - They sure are. Will see in action soon.

Deep Mysteries of Deep Learning

- Why does it converge at all?
- Why does it work well?
 - non-convex optimisation
 - many local minima
 - generalises to unseen data
- How does regularisation happen?
- Why do high dimensions seem to help?
 - segue into project...

What is (Supervised) ML²?

Given data $\{x_i, y_i\}_{i=1}^N$ and a model $g(x_i; w)$,
 $G(x_i)$

Find w s.t. g is a good approximator of G .

Training: $\hat{w} = \arg \min_w \frac{1}{N} \sum_{i=1}^N \ell(G(x_i), g(x_i, w))$

Testing: What we actually want is

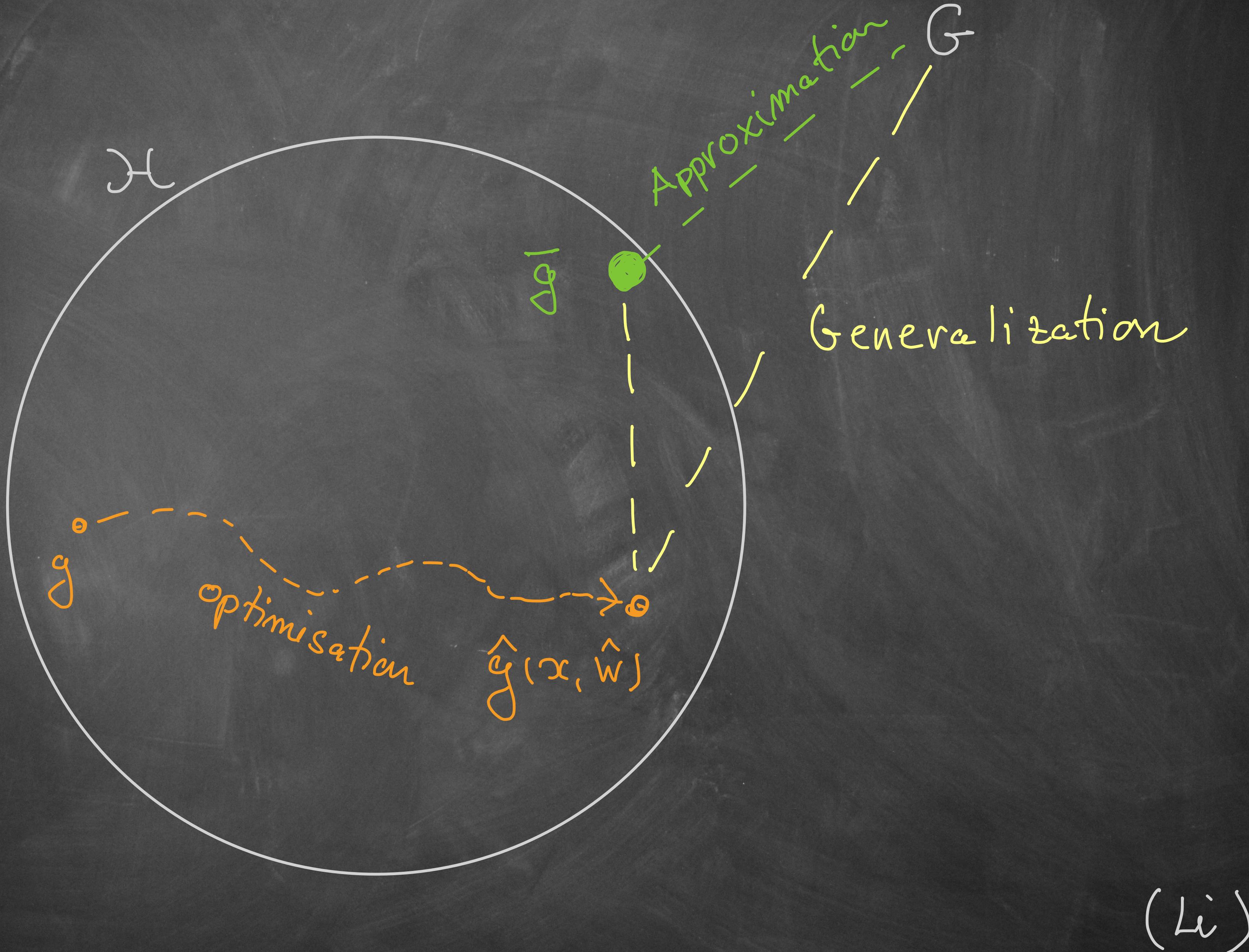
$$\hat{w}^* = \arg \min \mathbb{E}_{x \sim P} \ell(G(x), g(x, w))$$

This poses 3 separate (but interlinked) aspects:

Approximation: choice of model (Hypothesis space)

Optimization: how to find \hat{w}

Generalization: how good is $\hat{g}(x, \hat{w})$ on unseen
data



How to think about ML

Optimisation

- Why does gradient descent work so well?
- Statistical mechanics
- Why do momentum methods work even better?
continuous time perspective

Deep Residual Networks

Architecture

- Why do some architectures work better than others?
continuous time perspective

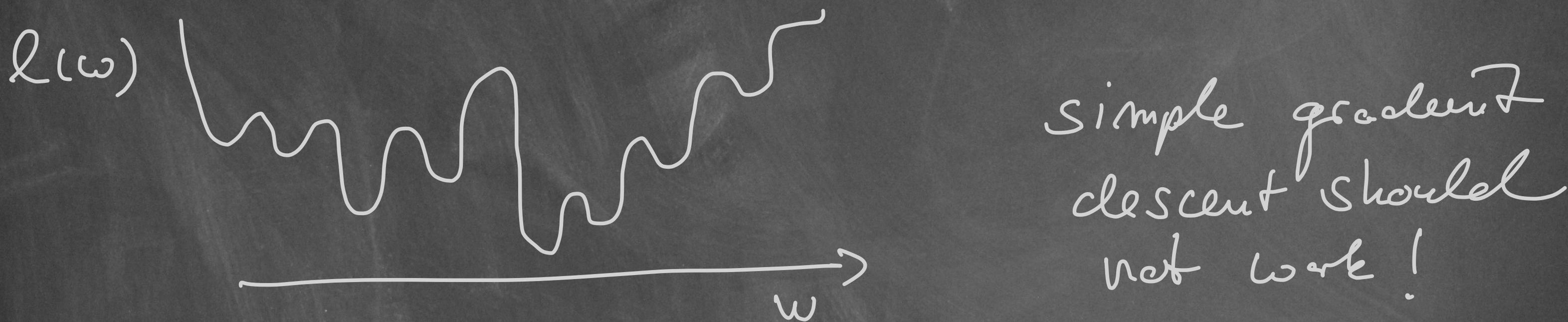
Approximation Error

- What type of functions can a neural network approximate?
Function spaces

The landscape of loss functions

Mean-square error : $\ell(w) = \sum_{j=1}^N \|y_j - \hat{g}(x_j; w)\|^2$

This is a nonconvex function, so in principle we can have multiple minima (alongside maxima and saddle points)



- Observation:
- * simple gradient descent works well to train high-dimensional deep neural nets and works less well in shallow networks
 - * shallow networks can match accuracy of deep networks (similar expressivity)
(Ba & Carreira 2014)

High-dimension non-convex cost functions

What type of critical point is most common?

- Global minima?
- Local minima?
- Saddle points?
- Local maxima?
- Global maxima?

(What does the distribution of eigenvalues of a large, random Hessian look like?)

Gradient Descent and the Structure of Neural Network Cost Functions

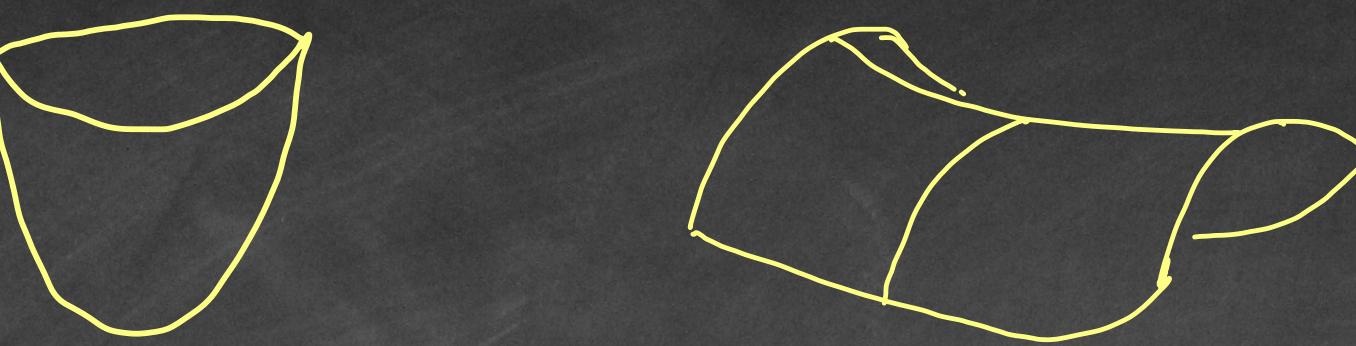
presentation by Ian Goodfellow

adapted for www.deeplearningbook.org
from a presentation to the
CIFAR Deep Learning summer school on August 9, 2015

Are saddle points or local minima more common?

- Imagine for each eigenvalue, you flip a coin
- If heads, the eigenvalue is positive, if tails, negative
- Need to get all heads to have a minimum
- Higher dimensions -> exponentially less likely to get all heads
- Random matrix theory:
- The coin is weighted; the lower J is, the more likely to be heads
 - So most local minima have low J !
 - Most critical points with high J are saddle points!

Basic intuition:



For a minimum in \mathbb{R}^d all d eigenvalues of the Hessian are positive; at a saddle point some are positive, some are negative (ignoring the case of semi-definite Hessians).

Draw the sign of the eigenvalues of the Hessian by coin tossing

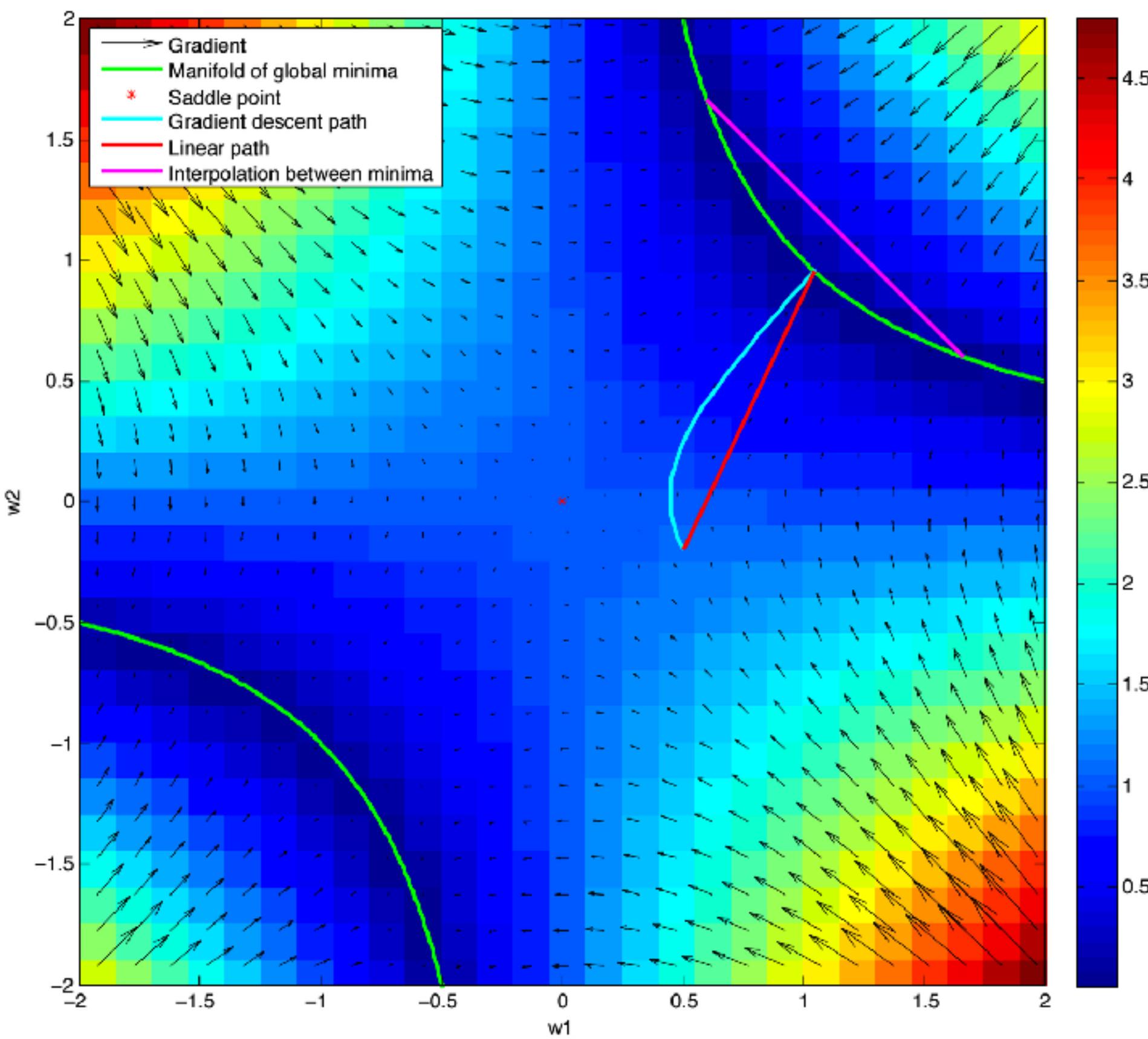
→ Getting a minimum becomes exponentially less likely for increasing dimension d .

But we need more:

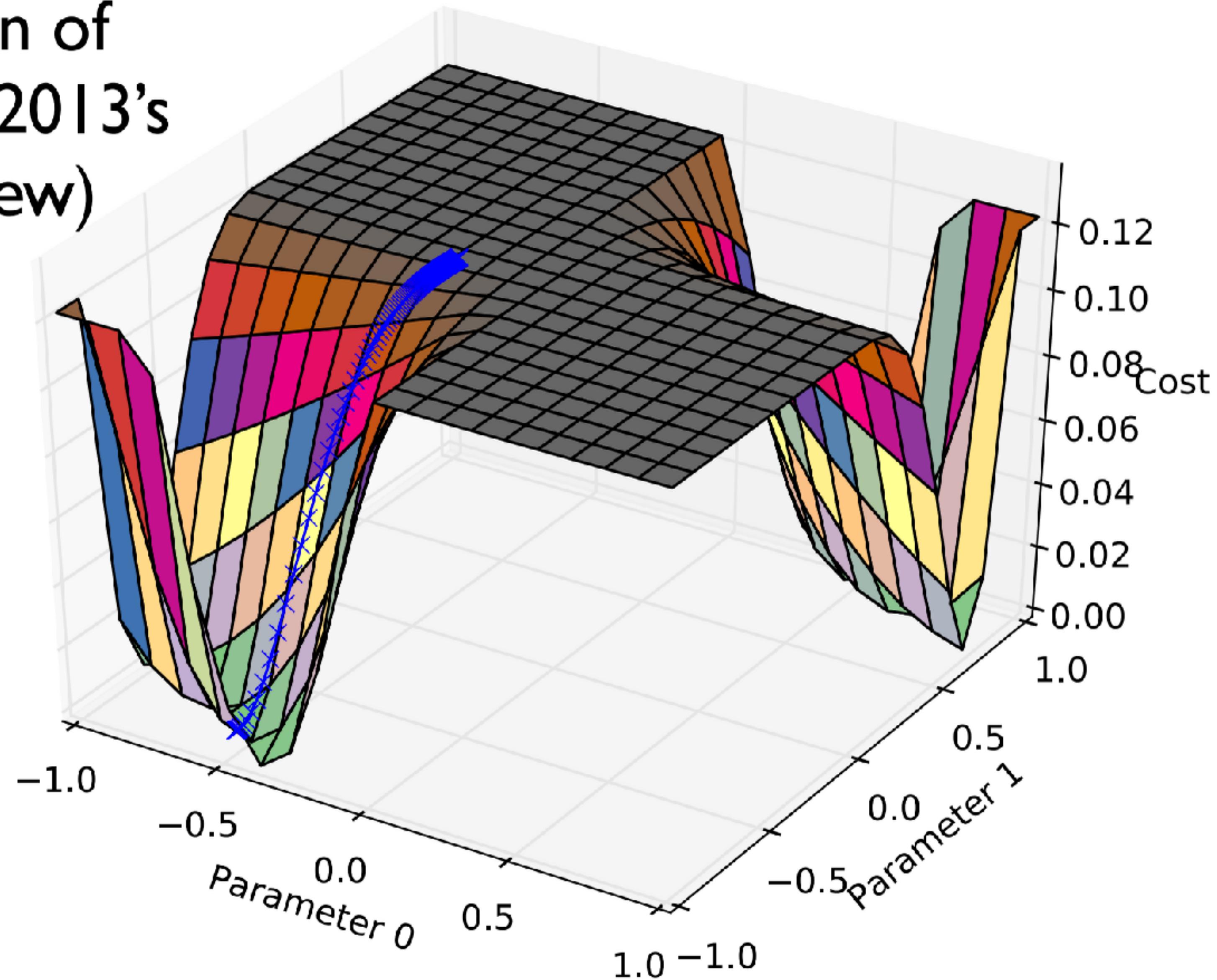


Do neural nets have saddle points?

- Saxe et al, 2013:
- neural nets without non-linearities have many saddle points
- all the minima are global
- all the minima form a connected manifold



(Cartoon of
Saxe et al 2013's
worldview)



Do neural nets have saddle points?

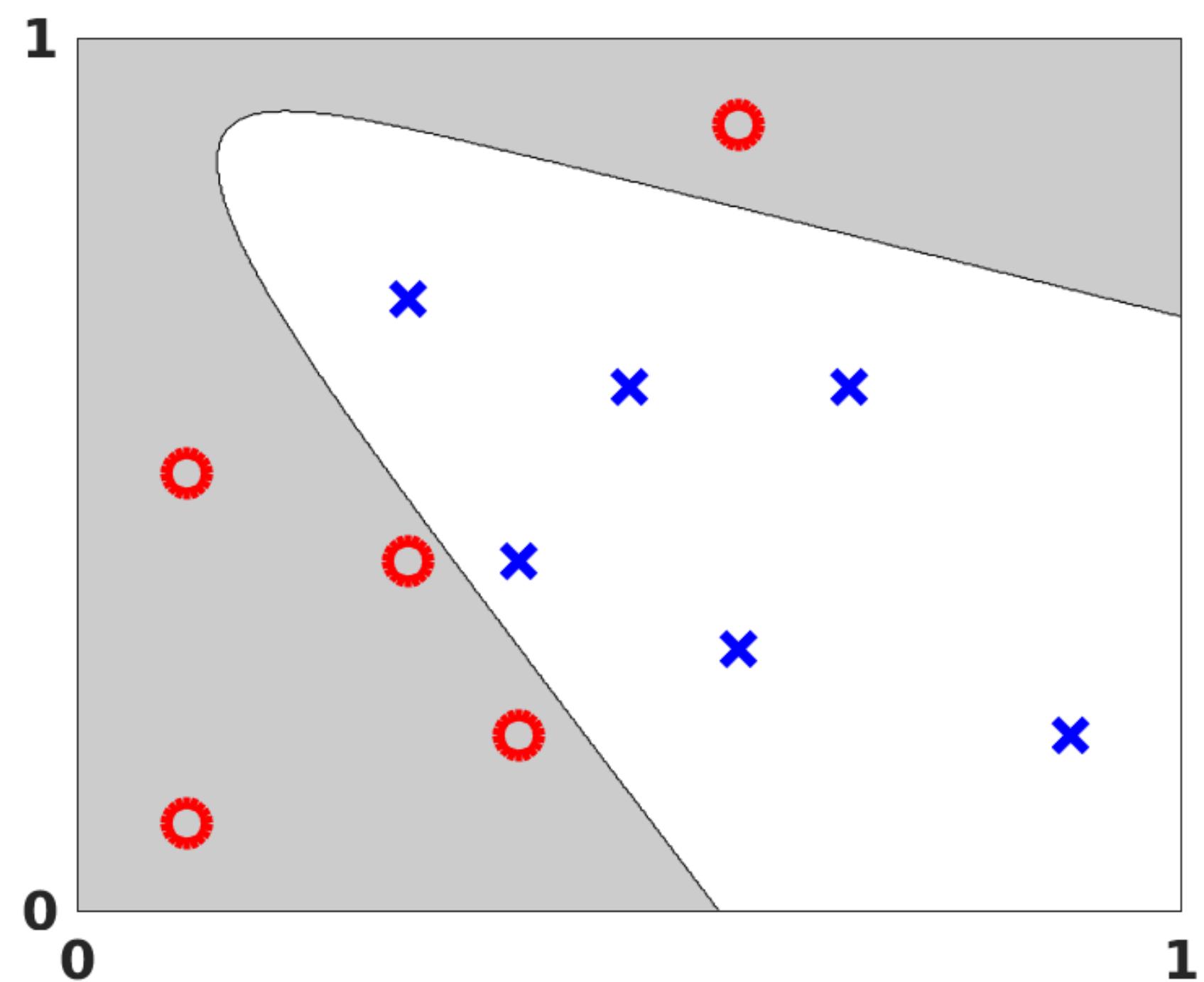
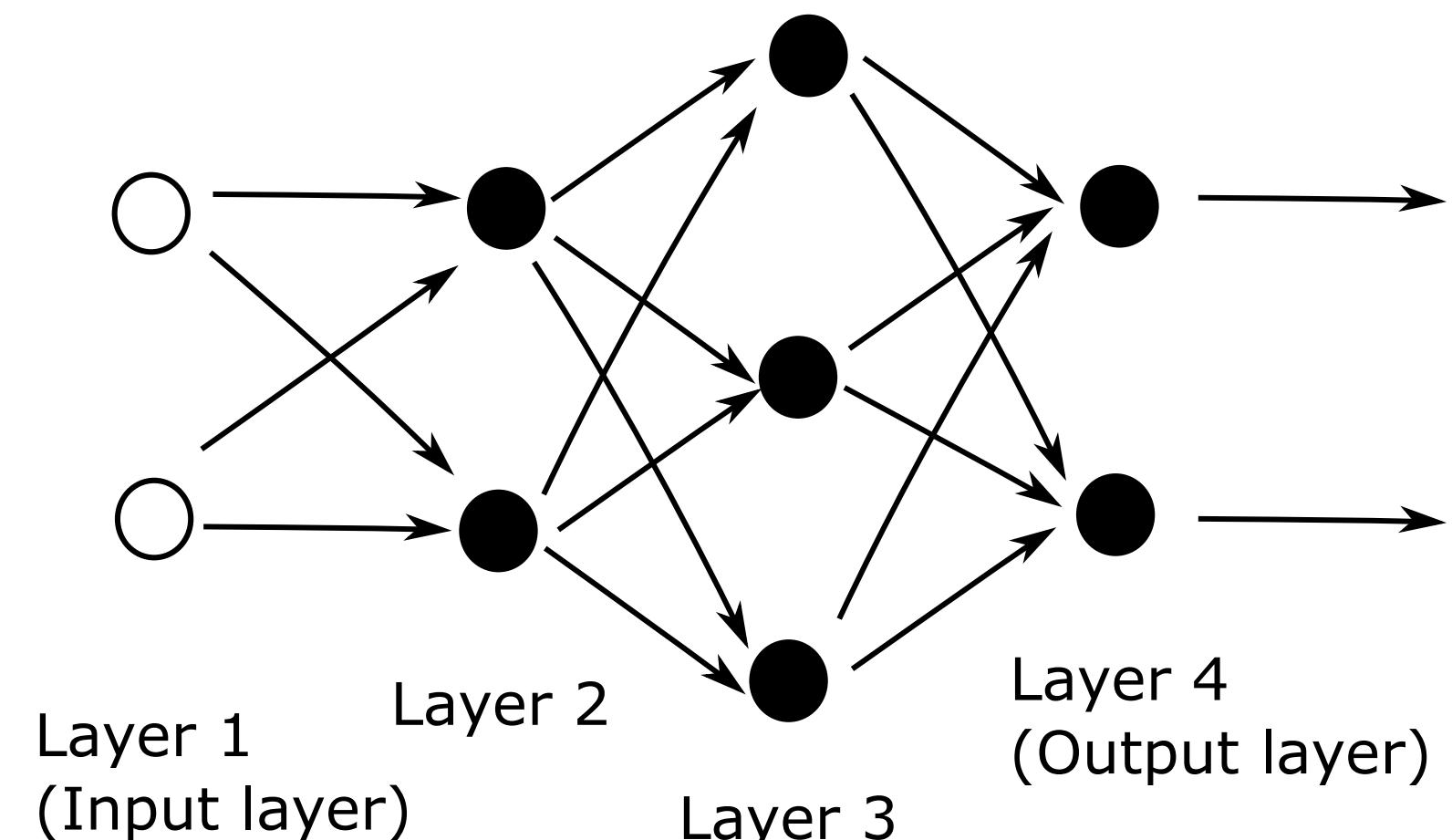
- Dauphin et al 2014: Experiments show neural nets do have as many saddle points as random matrix theory predicts
- Choromanska et al 2015: Theoretical argument for why this should happen
- Major implication: **most minima are good, and this is more true for big models.**
- Minor implication: the reason that *Newton's method* works poorly for neural nets is its attraction to the ubiquitous saddle points.

Mathematics > History and Overview*[Submitted on 17 Jan 2018]*

Deep Learning: An Introduction for Applied Mathematicians

Catherine F. Higham, Desmond J. Higham

Multilayered artificial neural networks are becoming a pervasive tool in a host revolution are familiar concepts from applied and computational mathematics; linear algebra. This article provides a very brief introduction to the basic ideas perspective. Our target audience includes postgraduate and final year undergr the area. The article may also be useful for instructors in mathematics who wis deep learning techniques. We focus on three fundamental questions: what is a stochastic gradient method? We illustrate the ideas with a short MATLAB code state-of-the art software on a large scale image classification problem. We fin



$$\text{Cost} \left(W^{[2]}, W^{[3]}, W^{[4]}, b^{[2]}, b^{[3]}, b^{[4]} \right) = \frac{1}{10} \sum_{i=1}^{10} \frac{1}{2} \|y(x^{(i)}) - F(x^{(i)})\|_2^2.$$

Download:

- PDF
- Other formats
(license)

Current browse context:
math.HO[< prev](#) | [next >](#)[new](#) | [recent](#) | [1801](#)

Change to browse by:

[cs](#)[cs.LG](#)[math](#)[math.NA](#)[stat](#)[stat.ML](#)

Mathematics > History and Overview*[Submitted on 17 Jan 2018]***Deep Learning: An Introduction for Applied Mathematicians**[Catherine F. Higham](#), [Desmond J. Higham](#)**Download:**

- [PDF](#)
- [Other formats
\(license\)](#)

Current browse context:
[math.HO](#)

Stochastic Gradient Descent: 2 approaches

$$p \rightarrow p - \eta \nabla \text{Cost}(p). \quad \text{where} \quad \nabla \text{Cost}(p) = \frac{1}{N} \sum_{i=1}^N \nabla C_{x^{\{i\}}}(p). \quad \text{Too expensive!}$$

- Approach 1:** 1. Shuffle the integers $\{1, 2, 3, \dots, N\}$ into a new order, $\{k_1, k_2, k_3, \dots, k_N\}$.
 2. For $i = 1$ upto N , update

$$(4.7) \qquad \qquad \qquad p \rightarrow p - \eta \nabla C_{x^{\{k_i\}}}(p).$$

- Approach 2:** 1. Choose m integers, k_1, k_2, \dots, k_m , uniformly at random from $\{1, 2, 3, \dots, N\}$.
 2. Update

$$(4.8) \qquad \qquad \qquad p \rightarrow p - \eta \frac{1}{m} \sum_{i=1}^m \nabla C_{x^{\{k_i\}}}(p).$$

Stochastic gradient descent

Algorithm 3.4: pseudo-code of a Stochastic Gradient Algorithm to minimise $f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x})$.

input: Gradient functions $\nabla f_i(\mathbf{x})$, and starting point \mathbf{x}_0 .

output: An estimate of the minimiser \mathbf{x}_k .

- 1 initialise $k = 0$;
 - 2 **while** *not reached compute time deadline* **do**
 - 3 Choose index $i_k \in \{1, \dots, m\}$ uniformly at random;
 - 4 Choose $\lambda_k > 0$ somehow;
 - 5 Set $\mathbf{x}_{k+1} = \mathbf{x}_k - \lambda_k \nabla f_{i_k}(\mathbf{x}_k)$;
 - 6 $k = k + 1$;
 - 7 **end**
-

SGD with momentum

Algorithm 3.5: pseudo-code of a Stochastic Gradient Descent Algorithm with momentum to minimise $f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x})$.

input: Gradient functions $\nabla f_i(\mathbf{x})$, velocity estimate \mathbf{v}_0 and starting point \mathbf{x}_0 .

output: An estimate of the minimiser \mathbf{x}_k .

- 1 initialise $k = 0$;
 - 2 **while** *not reached compute time deadline* **do**
 - 3 Choose index $i_k \in \{1, \dots, m\}$ uniformly at random;
 - 4 Choose $\lambda_k > 0$, $\alpha_k > 0$ somehow;
 - 5 Let $\mathbf{v}_{k+1} = \alpha_k \mathbf{v}_k - \lambda_k \nabla f_{i_k}(\mathbf{x}_k)$;
 - 6 Set $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{v}_{k+1}$;
 - 7 $k = k + 1$;
 - 8 **end**
-

SGD with momentum

Algorithm 3.5: pseudo-code of a Stochastic Gradient Descent Algorithm with momentum to minimise $f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x})$.

input: Gradient functions $\nabla f_i(\mathbf{x})$, velocity estimate \mathbf{v}_0 and starting point \mathbf{x}_0 .

output: An estimate of the minimiser \mathbf{x}_k .

- 1 initialise $k = 0$;
 - 2 **while** *not reached compute time deadline* **do**
 - 3 Choose index $i_k \in \{1, \dots, m\}$ uniformly at random;
 - 4 Choose $\lambda_k > 0, \alpha_k > 0$ somehow;
 - 5 Let $\mathbf{v}_{k+1} = \alpha_k \mathbf{v}_k - \lambda_k \nabla f_{i_k}(\mathbf{x}_k)$; “velocity”, accumulates gradients
 - 6 Set $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{v}_{k+1}$;
 - 7 $k = k + 1$;
 - 8 **end**
-

SGD with Nesterov momentum

Algorithm 3.6: pseudo-code of a Stochastic Gradient Descent Algorithm with Nesterov momentum to minimise $f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x})$.

input: Gradient functions $\nabla f_i(\mathbf{x})$, velocity estimate \mathbf{v}_0 and starting point \mathbf{x}_0 .

output: An estimate of the minimiser \mathbf{x}_k .

- 1 initialise $k = 0$;
 - 2 **while** *not reached compute time deadline* **do**
 - 3 Choose index $i_k \in \{1, \dots, m\}$ uniformly at random;
 - 4 Choose $\lambda_k > 0, \alpha_k > 0$ somehow;
 - 5 Let $\mathbf{v}_{k+1} = \alpha_k \mathbf{v}_k - \lambda_k \nabla f_{i_k}(\mathbf{x}_k + \alpha_k \mathbf{v})$;
 - 6 Set $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{v}_{k+1}$;
 - 7 $k = k + 1$;
 - 8 **end**
-

A birds-eye view of optimization algorithms

FABIAN PEDREGOSA
Google AI

November
2018

<http://fa.bianp.net/teaching/2018/eecs227at/>

also <https://distill.pub/2017/momentum/>

Do first, then understand

1. Run the compareOptimisationAlgos.ipynb notebook
 - From Liquet, Moka, & Nazarathy: deeplearningmath.org
 - This does “deterministic” GD, but same principle
2. Heavy ball analogy!
3. Adapt the code to do SGD, on the world’s most simple ML problem
 - Start with “online” SGD
 - Then do minibatch SGD (probably with more datapoints)

Continuous Time Analysis of Momentum Methods

Nikola B. Kovachki

*Computing and Mathematical Sciences
California Institute of Technology
Pasadena, CA 91125, USA*

NKOVACHKI@CALTECH.EDU

Andrew M. Stuart

*Computing and Mathematical Sciences
California Institute of Technology
Pasadena, CA 91125, USA*

ASTUART@CALTECH.EDU

Editor: Suvrit Sra

Abstract

Gradient descent-based optimization methods underpin the parameter training of neural networks, and hence comprise a significant component in the impressive test results found in a number of applications. Introducing stochasticity is key to their success in practical problems, and there is some understanding of the role of stochastic gradient descent in this context. Momentum modifications of gradient descent such as Polyak's Heavy Ball method (HB) and Nesterov's method of accelerated gradients (NAG), are also widely adopted. In this work our focus is on understanding the role of momentum in the training of neural networks, concentrating on the common situation in which the momentum contribution is fixed at each step of the algorithm. To expose the ideas simply we work in the deterministic setting.

Our approach is to derive continuous time approximations of the discrete algorithms; these continuous time approximations provide insights into the mechanisms at play within the discrete algorithms. We prove three such approximations. Firstly we show that standard implementations of fixed momentum methods approximate a time-rescaled gradient descent flow, asymptotically as the learning rate shrinks to zero; this result does not distinguish momentum methods from pure gradient descent, in the limit of vanishing learning

Adaptive gradient methods

Algorithm 8.4 The AdaGrad algorithm

Require: Global learning rate ϵ

Require: Initial parameter $\boldsymbol{\theta}$

Require: Small constant δ , perhaps 10^{-7} , for numerical stability

Initialize gradient accumulation variable $\mathbf{r} = \mathbf{0}$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)})$.

 Accumulate squared gradient: $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$.

 Compute update: $\Delta \boldsymbol{\theta} \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$. (Division and square root applied element-wise)

 Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$.

end while

Adaptive gradient methods

Algorithm 8.4 The AdaGrad algorithm

Require: Global learning rate ϵ

Require: Initial parameter θ

Require: Small constant δ , perhaps 10^{-7} , for numerical stability

Initialize gradient accumulation variable $r = \mathbf{0}$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$.

 Accumulate squared gradient: $r \leftarrow r + \mathbf{g} \odot \mathbf{g}$.

 Compute update: $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{r}} \odot \mathbf{g}$. (Division and square root applied element-wise)
 rescales gradient over time

 Apply update: $\theta \leftarrow \theta + \Delta \theta$.

end while

Adaptive gradient methods

Algorithm 8.7 The Adam algorithm

Require: Step size ϵ (Suggested default: 0.001)

Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1]$.
(Suggested defaults: 0.9 and 0.999 respectively)

Require: Small constant δ used for numerical stabilization (Suggested default:
 10^{-8})

Require: Initial parameters $\boldsymbol{\theta}$

Initialize 1st and 2nd moment variables $\mathbf{s} = \mathbf{0}$, $\mathbf{r} = \mathbf{0}$

Initialize time step $t = 0$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with
 corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

 Update biased first moment estimate: $\hat{\mathbf{s}} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

 Update biased second moment estimate: $\hat{\mathbf{r}} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

 Correct bias in first moment: $\hat{\mathbf{s}} \leftarrow \frac{\hat{\mathbf{s}}}{1 - \rho_1^t}$

 Correct bias in second moment: $\hat{\mathbf{r}} \leftarrow \frac{\hat{\mathbf{r}}}{1 - \rho_2^t}$

 Compute update: $\Delta\boldsymbol{\theta} = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}}} + \delta}$ (operations applied element-wise)

 Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta\boldsymbol{\theta}$

end while

Batch Normalization

OpenAI

Ian Goodfellow

Deep Learning Study Group
San Francisco
September 12, 2016

Batch Normalization

$$Z = XW$$

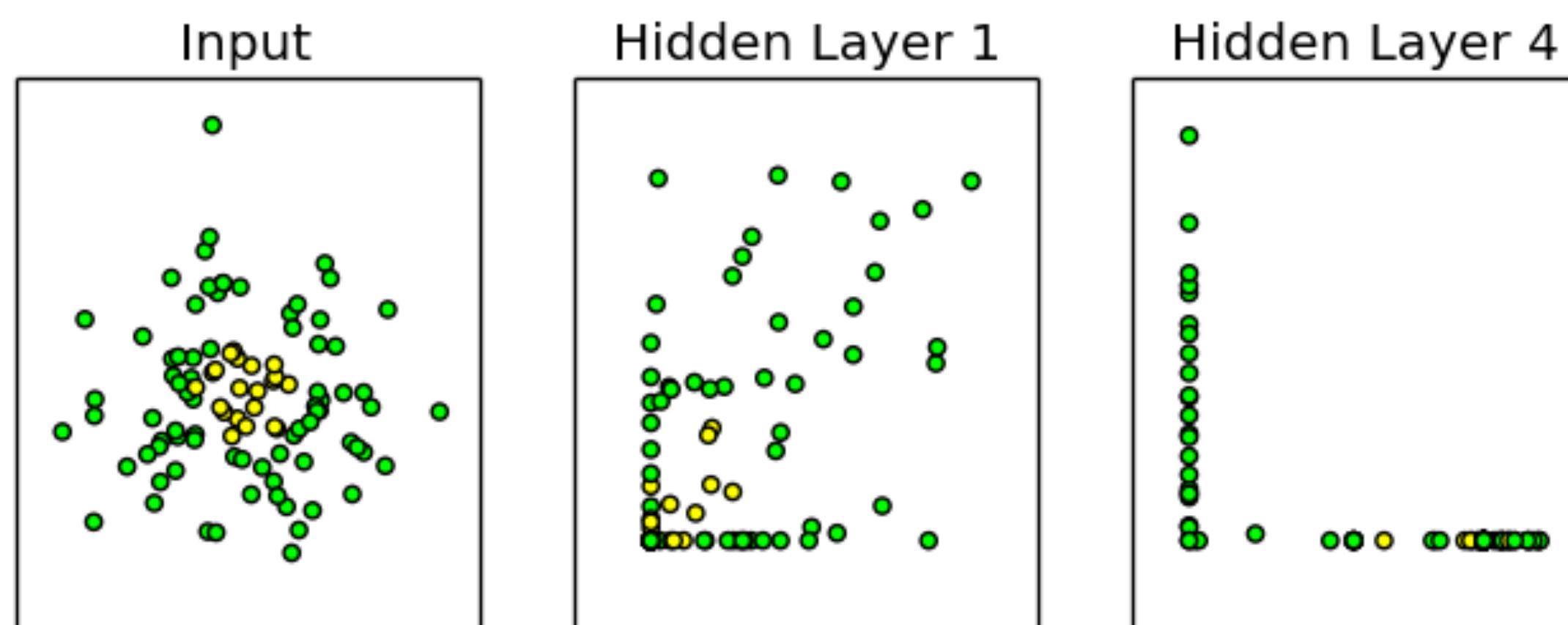
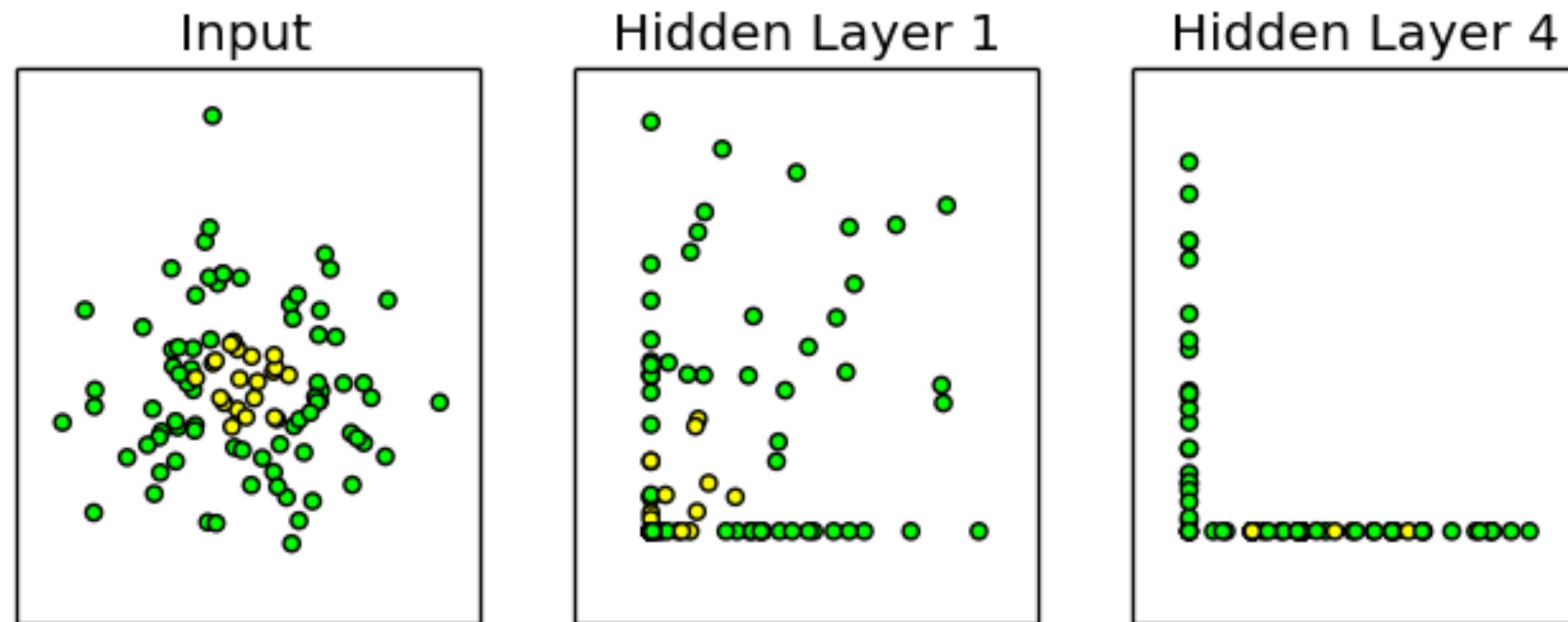
$$\tilde{Z} = Z - \frac{1}{m} \sum_{i=1}^m z_{i,:}$$

$$\hat{Z} = \frac{\tilde{Z}}{\sqrt{\epsilon + \frac{1}{m} \sum_{i=1}^m \tilde{Z}_{i,:}^2}}$$

$$H = \max\{0, \gamma \hat{Z} + \beta\}$$

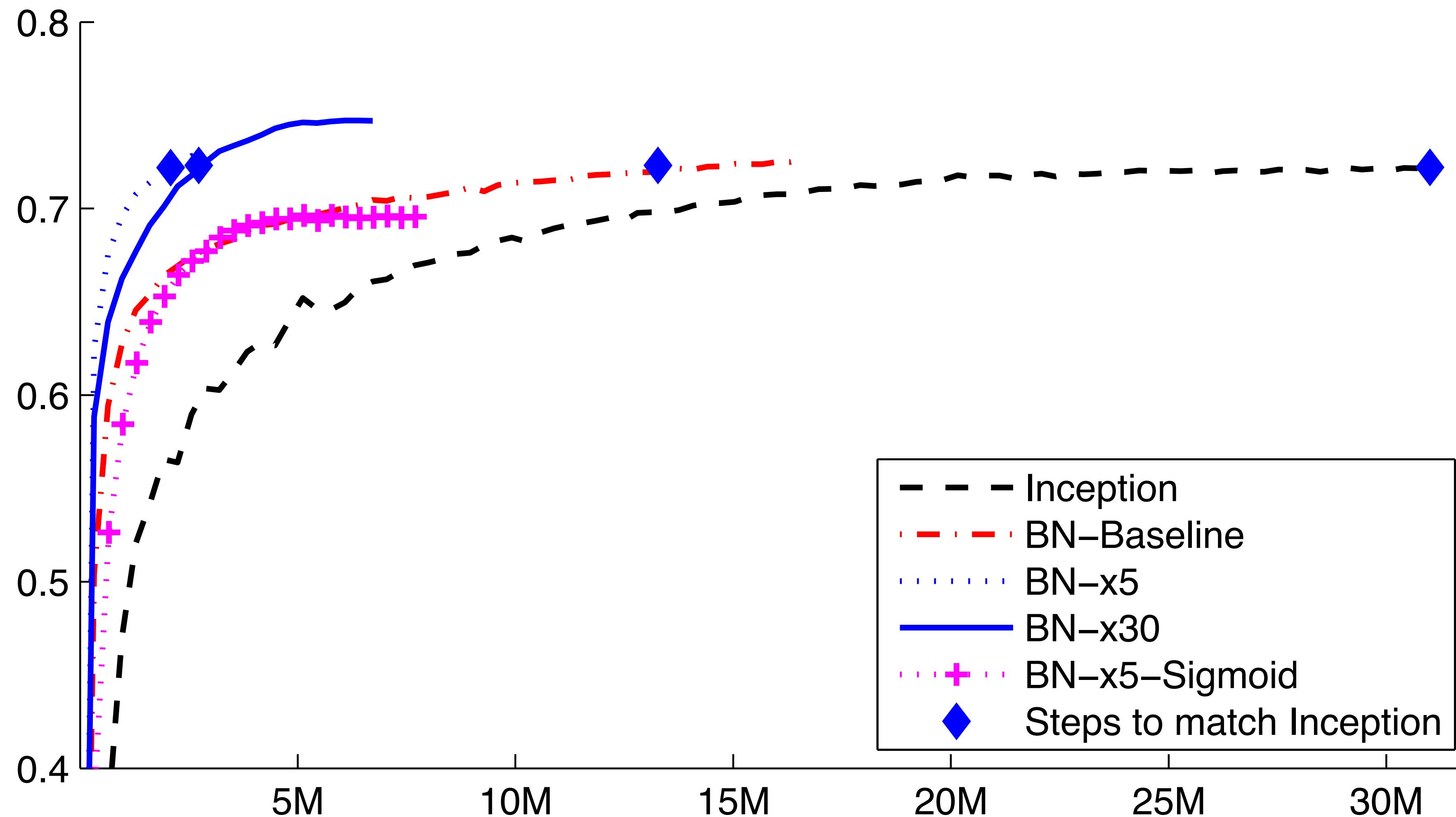
“Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift,” Ioffe and Szegedy 2015

Before SGD step



After SGD step

“Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift,” Ioffe and Szegedy 2015



“Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift,” Ioffe and Szegedy 2015