#### The role of over-parametrisation in NNs

Levent Sagun, EPFL

#### **Classical bias-variance dilemma**



#### Classical bias-variance dilemma, or?



# Observation 1 GD vs SGD

1. Take an iid dataset and split into two parts  $\mathcal{D}_{train} \& \mathcal{D}_{test}$ 2. Form the loss using only  $\mathcal{D}_{train}$   $\mathcal{L}_{train}(\theta) = \frac{1}{|\mathcal{D}_{train}|} \sum_{(x,y)\in\mathcal{D}_{train}} \ell(y, f(\theta; x))$ 3. Find:  $\theta^* = \arg \min \mathcal{L}_{train}(\theta)$ 4. ...and hope that it will work on  $\mathcal{D}_{test}$ 

- N : number of parameters  $heta \in \mathbb{R}^N$
- P : number of examples in the *training* set  $|\mathcal{D}_{train}|$

## Moving on the fixed landscape

1. Take an iid dataset and split into two parts  $\mathcal{D}_{train} \& \mathcal{D}_{test}$ 2. Form the loss using only  $\mathcal{D}_{train}$   $\mathcal{L}_{train}(\theta) = \frac{1}{|\mathcal{D}_{train}|} \sum_{(x,y)\in\mathcal{D}_{train}} \ell(y, f(\theta; x))$ 3. Find:  $\theta^* = \arg \min \mathcal{L}_{train}(\theta)$ 

4. ...and hope that it will work on  $\mathcal{D}_{test}$ 

- N : number of parameters  $heta \in \mathbb{R}^N$
- P : number of examples in the *training* set  $|\mathcal{D}_{train}|$

#### "Stochastic gradient learning in neural networks" Léon Bottou, 1991

• The total gradient (3) converges to a *local minimum* of the cost function. The algorithm then <u>cannot escape this local minimum</u>, which is sometimes a poor solution of the problem.

In practical situations, the gradient algorithm may get stuck in an area where the cost is extremely ill conditionned, like a deep ravine of the cost function. This situation actually is a local minimum in a subspace defined by the largest eigenvalues of the Hessian matrix of the cost.

The stochastic gradient algorithm (4) usually is able to escape from such bothersome situations, thanks to its random behavior (Bourrely, 1989).

#### GD is bad use SGD

#### Bourrely, 1988

#### 5 CONCLUSION

It has been shown that the difficulty in parallel learning is due to the fact that the parallel algorithm does not really use the stochastic algorithm. Two solutions are presently proposed to prevent the system from falling into a local minimum.

1) Add momentum to the algorithm such that it can "roll past" a local minimum. Thus the algorithm then becomes:

 $W_{t+1} = (1-\alpha) W_t - \varepsilon \alpha f(W_t, X_t)$ 

where f is the error gradient Q relative to W

2) One can add a random "noise" to the gradient calculations. One method of performing this task is to calculate the gradients in an approximate manner. This variation could be modelled as a type of 'Brownian motion', using a temperature function (similar to simulated annealing). This temperature could be lowered relative to the remaining system error. For example, the variation in gradients could follow a Gaussian distribution. Thus, for example:

 $W_{t+1} = W_t - \varepsilon N(f(W_t, X_t), k\sqrt{Temp})$ 

where f is the error gradient Q relative to W and N is a function giving a Gaussian random variable.

Both of these approaches are presently under research.

#### Fully connected network on MNIST: $N\sim 450 {\rm K}$



Sagun, Guney, LeCun, Ben Arous 2014

### Different regimes depending on ${\cal N}$

#### Bourrely, 1988

Evaluation of computational time and learning time is achieved by training the network for the handwritten numbers recognition task. The network is designed as follows : 400 input units (a 20x20 grid), <u>5 hidden</u> units and 10 output units. It must perform a classification task: for each input number, the network must activate the correct output unit (0 to 9). In addition, it must overcome distortions such as vertical and horizontal translations, scaling, rotation and random white noise. Numbers are coded on matrices of 20x20 real grey-levels. Table 1 to

#### Fully connected network on MNIST: $N\sim 450 {\rm K}$



Average number of mistakes: SGD 174, GD 194

Further empirical confirmations

- Teacher-Student networks
- landscape of the *spin glass* model
- GD vs SGD on fully-connected MNIST
- GD vs SGD on noisy inputs, scrambled labels...

#### Regime where SGD is really special?

Where common wisdom may be true (Keskar et. al. 2016):  $\rightarrow$  Similar training error, but gap in the test error.



- Jastrzębski et. al. 2017
- Goyal et. al. 2017
- Shallue and Lee et. al. 2018
- McCandlish et. al. 2018
- Smith et. al. 2018

- Jastrzębski et. al. 2017
- Goyal et. al. 2017
- Shallue and Lee et. al. 2018
- McCandlish et. al. 2018
- Smith et. al. 2018



- Jastrzębski et. al. 2017
- Goyal et. al. 2017
- Shallue and Lee et. al. 2018
- McCandlish et. al. 2018
- Smith et. al. 2018



- Jastrzębski et. al. 2017
- Goyal et. al. 2017
- Shallue and Lee et. al. 2018
- McCandlish et. al. 2018
- Smith et. al. 2018





Why is it important?

#### Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour

Priya Goyal Piotr Dollár Lukasz Wesolowski Aapo Kyrola

ollár Ross Girshick rola Andrew Tulloch Pieter Noordhuis Yangqing Jia Kaiming He

Facebook

#### Abstract

Deep learning thrives with large neural networks and large datasets. However, larger networks and larger datasets result in longer training times that impede research and development progress. Distributed synchronous SGD offers a potential solution to this problem by dividing



Now anyone can train Imagenet in 18 minutes Written: 10 Aug 2018 by Jeremy Howard Note from Jeremy: I'll be teaching Deep Learning for Coders at the University of Noordhuis San Francisco starting in October; if you've got at least a year of coding 2 Jia Kaiming He experience, you can <u>apply here</u>. A team of fast.ai alum Andrew Shaw, DIU researcher Yaroslav Bulatov, and I have managed to train Imagenet to 93% accuracy in just 18 minutes, using 16 public AWS cloud instances, each with 8 <u>NVIDIA V100</u> GPUs, running the <u>fastai</u> and <u>PyTorch</u> libraries. This is a new speed record for training Imagenet to this accuracy on publicly available infrastructure, and is 40% faster than Google's DAWNBench record on their proprietary TPU Pod cluster. Our approach uses the same number of processing units as Google's benchmark (128) and costs around \$40 to run.

Now arr minu

Written: 10

*Note from Jer* San Francisco experience, you

A team of fast.ai alum managed to train <u>Imag</u> cloud instances, each wit libraries. This is a new spe publicly available infrastruc on their proprietary <u>TPU Pod</u> processing units as Google's be





- Optimization of the training function is easy *... as long as there are enough parameters*
- Effects of SGD is a little bit more subtle ... but exact reasons are somewhat unclear

# Observation 2 A look at the bottom of the loss

## Different kinds of minima

- Continuing with Keskar et al (2016): LB  $\rightarrow$  sharp, SB  $\rightarrow$  wide...
- Also see Jastrzębski et. al. (2017), Chaudhari et. al. (2016)...
- Older considerations Pardalos et. al. (1993)
- Sharpness depends on parametrization: Dinh et. al. (2017)

#### Different kinds of minima

- Continuing with Keskar et al (2016): LB  $\rightarrow$  *sharp*, SB  $\rightarrow$  *wide...*
- Also see Jastrzębski et. al. (2018), Chaudhari et. al. (2016)...
- Older considerations Pardalos et. al. (1993)
- Sharpness depends on parametrization: Dinh et. al. (2017)

$$\tilde{H}_{\Lambda,f}(R) \equiv f^{-1} \left\{ \int S_{\Lambda}(R-R') f[H(R')] dR' \right\}$$
(2)

where R is a multidimensional vector representing all the coordinates in the molecule.

One of the simplest and most useful forms for  $S_{\Lambda}$  is a Gaussian

$$S_{\Lambda}(R) \equiv C(\Lambda)e^{-R\Lambda^{-2}R}$$

$$C(\Lambda) \equiv \pi^{-d/2}Det^{-1}(\Lambda) \qquad (3)$$

where d is the total dimensionality of R. The function f included in (2) allows for nonlinear averaging. Two choices motivated by physical considerations are f(x) = x and  $f(x) = e^{-x/k_BT}$ . These choices correspond respectively to the "diffusion equation" and "effective energy" methods which are described below. Wu [77] has presented a general discussion of transformations of the form of (2).

A highly smoothed  $H_{\Lambda,f}$  (from which all high spatial-frequency components have been removed) will in most cases have fewer local minima than the unsmoothed ("bare") function, so it will be much easier to identify its global minimum. If the strong spatial-scaling hypothesis is correct, the position of this minimum can then be iteratively tracked by localminimization as  $\Lambda$  decreases. As  $\Lambda \to 0$ , the position will approach the global minimizer of the bare objective function.

### A look through the local curvature

Eigenvalues of the Hessian at the beginning and at the end



### A look through the local curvature

Eigenvalues of the Hessian at the beginning and at the end



### A look through the local curvature

Increasing the batch-size leads to larger outlier eigenvalues:



#### A look at the structure of the loss

Recall the loss per sample:  $\ell(y, f(\theta; x))$ 

- $\ell$  is convex (MSE, NLL, hinge...)
- *f* is non-linear (CNN, FC with ReLU...)

We can see the Hessian of the loss as:

$$abla^2\ell(f)=\ell''(f)
abla f
abla f^T+\ell'(f)
abla^2 f$$

a detailed study on this can be found in Papyan 2019

### More on the lack of barriers

- 1. Freeman and Bruna 2017: barriers of order 1/N
- 2. Baity-Jesi & Sagun et. al. 2018: no barriers in SGD dynamics
- 3. Xing et. al. 2018: no barrier crossing in SGD dynamics
- 4. Garipov et. al. 2018: no barriers between solutions
- 5. Draxler et. al. 2018: no barriers between solutions

## More on the lack of barriers

1. Freeman and Bruna 2017: barriers of order 1/N

2. Baity-Jesi et. al. 2018: no barrier crossing in SGD dynamics

- 3. Xing et. al. 2018: no barrier crossing in SGD dynamics
- 4. Garipov et. al. 2018: no barriers between solutions
- 5. Draxler et. al. 2018: no barriers between solutions



- A large and connected set of solutions *... possibly only for large N*
- Visible effects of SGD is on a tiny subspace ... again, exact reasons are somewhat unclear

# A simple example

#### Lessons from observations

Observation 1: easy to optimize Observation 2: flat bottom



# **Defining over-parametrization**

Several works joint *with: Mario Geiger, Stefano Spigler, Marco* Baity-Jesi, Stephane d'Ascoli, Arthur Jacot, Franck Gabriel, Clement Hongler, Giulio Biroli, & Matthieu Wyart

- 1. For large N the dynamics don't get stuck  $\rightarrow$  When is the training landscape *nice*?
- 2. Often N >> P, yet it doesn't it overfit  $\rightarrow$  Relationship of the landscape with generalization?

- N : number of parameters  $heta \in \mathbb{R}^N$
- P : number of examples in the *training* set  $|\mathcal{D}_{train}|$

Switch to squared-hinge from cross-entropy  $\ell(y, f( heta; x)) = rac{1}{2}max(0, 1 - yf( heta; x))^2$ 

- precise stopping condition
- clear stability condition

$$abla^2\ell(f) = 
abla f 
abla f^T + (1-yf) 
abla^2 f$$

- sum over unsatisfied constraints
- a local minimum is only possible if: N/2 < P (very loose)

#### Sharp transition to OP in NNs



- P : number of examples in the *training* set  $|\mathcal{D}_{train}|$
- $N^*$  : critical number of parameters that fits  $\mathcal{D}_{train}$

#### Sharp transition to OP in NNs



- N : number of parameters  $heta \in \mathbb{R}^N$
- P : number of examples in the *training* set  $|\mathcal{D}_{train}|$
- $N^*$  : critical number of parameters that fits  $\mathcal{D}_{train}$

#### Sharp transition to OP in NNs



- P : number of examples in the *training* set  $|\mathcal{D}_{train}|$
- $N^*$  : critical number of parameters that fits  $\mathcal{D}_{train}$

#### Jamming is linked to Generalization



#### Jamming is linked to Generalization



The peak itself is also observed in Advani and Saxe 2017

#### Jamming is linked to Generalization



#### **Recent independent work**

Belkin et. al. December 31, 2018



See also Neal et al. 18, Neyshabur et al. 15 & 17 for related work

## **Ensembling improves generalization**

Key: reducing fluctuations or increased regularization with N



### **Ensembling improves generalization**



Sagun, Geiger, d'Ascoli, Spigler, Biroli, Wyart 2019 (unpublished)

### **Ensembling improves generalization**



Sagun, Geiger, d'Ascoli, Spigler, Biroli, Wyart 2019 (unpublished)

#### *Potential impact:*

- Clear definition of OP can help guide design of models
- At finite *P* we have a proposal for the best generalization
- New directions for theoretical understanding

   → Belkin et. al. 18 March 2019
   → Hastie et. al. 19 March 2019

#### Future work

*On the model-data-algorithm interactions:* 

- Can we disentangle the algorithm?
- Can we entangle the model-data interactions to unite
  - model complexity measure
  - data complexity measure
    - the role of priors on performance!

# **Thank You!**