

# Caution remark

① Ongoing work



② Possible <sup>likely!</sup> errors !

# Solutions for DNNs

## dynamics and norms

1. Dynamical system approach : square loss

Linear case one layer

Multilayer nonlinear

( Conflict control, minimum norm )

2. Dynamical system approach : sep losses

- 1 layer linear, 1 layer nonlinear, Srebro, us,  
nonlinearization, regularization, norm control  
min norm

Dynamical systems approach :

square loss

We are interested in characterizing which solutions  $G D$  converges to. One reason eventually is to understand which kind of complexity control takes place in deep nets

## Notation

$$f(x) = W_k^{i,m} \dots \left( \dots \sigma(W_2^{qi} (\sigma(W_i^{ij} x_j))) \right)$$

$$\text{Jacobian of } f \mapsto \frac{\partial f}{\partial W_h^{ij}}$$

$$\text{Hessian of } f \mapsto \frac{\partial^2 f}{\partial W_h^{ij} \partial W_{h'}^{ab}}$$

## Why learning rate?

An intuition for a non-discrete math guy is given by

$$\dot{x} + \gamma(t)x = 0$$

Solution is  $x(t) = x_0 e^{-\int \gamma(t) dt}$ .

For convergence to minimum ( $\phi$ ) we need  $\int \gamma(t) dt \rightarrow \infty$

For robustness against noise such as discretization noise in

$$\dot{x} + \gamma(t)(x + u(t)) = 0$$

we also need  $\gamma(t) \rightarrow 0$



## Gradient dynamical systems and GD

The loss function is  $\equiv$  potential (Lagrangian) function

$$L(w) = \sum_i l_i(f_w(x_i), y_i)$$

The gradient dynamical system induced by GD is

$$\dot{w} = -\eta \nabla_w L(w) = -\eta F(w)$$

Square loss  $\Rightarrow l_i = (f_w(x_i) - y_i)^2$

We are interested in equilibrium points (convergence of GD) that is  $w^*$  s. t.  $\dot{w} = 0$

We are interested in stability of  $w^*$ . For this

we look at Hessian of  $L$  that is

$$H = \frac{\partial^2 L}{\partial w_k^{i^*} \partial w_{k'}^{e^b}}$$

$H$  of convex function is pos. def.



## Linear one-layer networks

In this case  $K = 1$  and

$$f(x) = W^{1T} x_T = W^T x$$

Thus

$$L_u(f) = \sum_n^N (2y_n - W^{1T} x_J^n)^2 \quad \begin{array}{l} y \in \mathbb{R}(-1,1) \\ x \in \mathbb{R}^d \end{array}$$

$$\dot{W}^{1T} = 2 \sum_n^N E_n x_J^n$$

$\Downarrow$

$$W_{\min}^{1T} \text{ s.t. } Y = WX, \quad W = YX^{-1}$$

Suppose overparametrization that is  $N < d$

Then  $f$  fits data that is  $E_n = 0 \quad \forall n = 1, \dots, N$   
at which point  $\dot{W}^{1T} = 0$

Notice that during dynamics only the weights that are in the span of the  $x^n$  change. The others in the null space of  $x^n$  do not change, the gradient is zero.

The Hessian of  $L$  is

$$H_{i,j} = \sum_n \begin{pmatrix} x_J^n & x_i^n \end{pmatrix}$$

$$= \text{at min} = \sum x^n x^{nT}$$

that is the sum of autocovariance matrices, each one being semidefinite (positive). Thus there are zero eigenvalues.

Show slide 

In this situation - degenerate  $H$  - there is no unique minimum and no explicit control of norm.

Regularization. Adding an arbitrarily small regularization term solves the problem here

Instead of

$$L = \sum_n (y_n - W^{1T} x_n^u)^2$$

I use  $L = \sum_n (y_n - W^{1T} x_n^u)^2 + \frac{\lambda}{2} \sum_T \|W^{1T}\|^2$

$\|W\|^2$

Thus

$$W^{1T} = \left( \sum_n E_n x_n^u + \lambda W^{1T} \right)$$

which is called "weight decay".

The equilibrium point shifts to

$$\lambda W^{1T} = \sum_n E_n x_n^u - \sum (y_n - W^{1T} x_n^u) x_n^u \quad \text{that is}$$

$$\lambda W^{i,j} = \sum_n \left[ y_n x_n^i x_n^j - W^{i,l} x_l^n x_j^n \right]$$

corresponding to

$$\begin{aligned} & \nabla (-Y + W X) X^T + \lambda W = 0 \\ & -Y X^T + W (\lambda I + X X^T) = 0 \\ & Y X^T = W (X X^T + \lambda I) \\ & W = Y X^T (X X^T + \lambda I)^{-1} \end{aligned}$$

The Hessian is now

$$H_{ij} = \sum_n x_i^n x_j^n + \lambda I$$

always positive definite for  $\forall \lambda > 0$

---

## Implicit regularization

For square loss and linear networks GD converges to the same solution of regularization with  $\lambda \rightarrow 0$  which is the pseudoinverse solution

$$W = Y X^T (X X^T)^{-1} = Y X^+ \text{ which is}$$

the minimum norm solution.

The condition for this to happen is that the initial condition for GD is  $W^0 \sim 0$ .

Then the degenerate components of the gradient  $\frac{\partial L}{\partial W^{ij}} = 0$  do not change the weights  $W^{ijT}$

which remain  $\sim 0$  thus min norm

stick

## Norm control, regularization, implicit regularization

① The regularization min is hyperbolic and independent of initial conditions. This also means one can perturb and get back to same  $W$ .

② The implicit regularization min is degenerate and depends on initial conditions. This means that perturbations during GD will change final  $W$ .

A similar situation holds for  
multilayer linear networks.

slides

# Deep nets, square loss

For general networks  $f(x) = W_{x^l}^{i,l} \sigma(\dots \sigma(W_{x^1}^{i,1} x_1))$

the gradient equations are

$$W_{k'}^{i,j} = - \frac{\partial}{\partial W_{k'}^{i,j}} L = - \sum_n E_n \frac{\partial f}{\partial W_{k'}^{i,j}} \left( - \lambda W_{k'}^{i,j} \right)$$

$$H_{abk'}^{i,jk'} = \frac{\partial^2}{\partial W_{k'}^{i,j} \partial W_{k'}^{ab}} L =$$

$$= \sum_n \frac{\partial f}{\partial W_{k'}^{i,j}} \frac{\partial f}{\partial W_{k'}^{ab}} - E_n \frac{\partial^2 f}{\partial W_{k'}^{i,j} \partial W_{k'}^{ab}}$$

$$\left( - \lambda I \right)$$

$$\delta_{ia} \delta_{jb} \delta_{kk'}$$

As for linear networks there is no finite equilibrium point because  $H$  is in general degenerate. This is clear in the deep polynomial case because only the  $W_{k'}^{i,j}$  in the span of the  $\sum_n \frac{\partial f}{\partial W_{k'}^{i,j}}$  change. If  $N < \# W_{k'}^{i,j}$  there is degeneracy.

More formally Kenji Kawaguchi proved

### Theorem 4 (Theory III)

Under standard assumption the Hessian matrix of a N.A which is overparametrized has at least one zero eigenvalue.

This is similar to linear case. However, unlike the linear case we cannot guarantee that starting with  $W' \approx 0$  the solution will be min. norm.

First, in general there will be many <sup>degenerate</sup> minima.  
Second, the degenerate directions near a minimum will not be degenerate during G.D.

This issue is related to the conditions under which I can use linearization tools - such as the Hessian - around minima of a dynamical system to characterize its behaviour. It turns out that a general condition for validity of linearization is that the Hessian is non-degenerate. The Hartman Grobman theorem says



## Hartman Grobman Theorem

Dynamical system  $\dot{W} = -F(W) = -\nabla_W L$   
If  $F$  has hyperbolic equilibrium  $W^*$  and Hessian of  $L$  does not have zero eigenvalues there, then  
 $\exists N$  of  $W^*$  and homeomorphism  $h: N \rightarrow \mathbb{R}^d$  s.t.  
 $h(W^*) = 0$  and in  $N$  the flow of  $\dot{W} = -F(W)$   
is topologically conjugate by the continuous map  
 $U = h(W)$  to the flow of linearized system  
 $\dot{U} = -H U$

Homeomorphism cont. map with inverse between  
topological spaces

Extension to deep nets via H.G for square loss

Regularization<sup>can</sup> transform degenerate Hessian at a min into hyperbolic min

Consider

$$\dot{W}_k^{i,j} = - \frac{\partial}{\partial w_k^{i,j}} L = - \sum_n E_n \frac{\partial f}{\partial w_k^{i,j}} - \lambda W_k^{i,j}$$

$$W^* \text{ s.t. } \dot{W} = 0 \iff E_n = \epsilon \forall n, f(x_n; w^*) - y_n = \epsilon$$

Then

$$H = \sum_n \frac{\partial^2 f}{\partial w_k^{i,j} \partial w_{k'}^{a,b}} - E_n \frac{\partial^2 f}{\partial w_k^{i,j} \partial w_{k'}^{a,b}} - \lambda \mathbf{I}$$

$\delta_{ia} \delta_{jb} \delta_{kk'}$

can be pos. def.

at a minimum  
If  $H \approx$  non degenerate  $\wedge$  then Hartman Grobman  
applies and minimum is hyperbolic.

This may happen when  $F_u \approx \varepsilon$ ;  
then Hessian may be pos def. for  
small.

min. norm and complexity control. Otherwise not, as shown  
in slide (fig 4 Theory III).

But ... it depends on  $H$  if  $F_u \neq 0$

# Dynamical systems approach :

## exponential type loss

$$\text{Exponential loss} \Rightarrow L = \sum_n^N l_i$$
$$l_i = e^{-y_i f(x_i)} \quad y_i = \pm 1$$

$$\text{Linear net} \Rightarrow f(x) = W^{jT} x_j$$

Assuming separable data  $y_n f(x_n) > 0 \rightarrow f(x_n)$   
GD generate dynamical system

$$\begin{aligned} \dot{W}^{jT} &= - \frac{\partial}{\partial W^{jT}} \sum_n^N e^{-W^{jT} x_j^n} \\ &= + \sum_n^N e^{-W^{jT} x_j^n} x_j^n \end{aligned}$$

with vector notation  $W^{jT} \rightarrow \underline{w}$

$$\dot{\underline{w}} = \sum_n^N \underline{x}^n e^{-\underline{w}^T \underline{x}^n}$$

Thus the components of  $w$  grow positive or negative to infinity slower and slower

## Frederick result (2017)

Lemma 1 GD with  $b_i = e^{-y_i f(x_i)}$  and  $f$  linear

has  $w(t) = \tilde{w} \log t + g(t)$  s. t.

$$\lim_{t \rightarrow \infty} \frac{w(t)}{\|w(t)\|} = \frac{\tilde{w}}{\|\tilde{w}\|} \quad \text{where } \tilde{w} \text{ is solution of}$$

hardmargin SVM, that is

$$\tilde{w} = \operatorname{argmin}_{w \in \mathbb{R}^d} \|w\|^2 \text{ s.t. } \forall u \quad w^\top x_u \geq 1$$

The convergence is independent of initial conditions

Remarks Notice that as for square loss the degenerate  $w$  components are not changed during GD. However normalization sets them to zero asymptotically...

We will show that

- 1) regularization restores hyperbolicity for linear nets
- 2) a normalization version of GD (similar to weight normalization) also restores pos def Hessian but only under very special data

Because of ① + ② it seems we have

similar situation as for square loss:

regularization (and perhaps normalization)

can guarantee hyperbolicity, thus

validity of Hartman-Grobman and thus

extension to DNNs, including control

of norm.

Next we study the effect of regularization  
and normalization

## Regularization for linear nets with exp loss

$$\dot{W} = - \frac{\partial}{\partial W} (L + \lambda W^2) = \sum_n x_n e^{-W^T x_n} - \lambda W$$

The minimum  $W^*$  is given by

$$\frac{1}{\lambda} \sum_n x_n e^{-W^{*T} x_n} = W^*$$

$$X = \frac{1}{\lambda} \sum e^{-W^{*T} x_n} x_n$$

$$X = \underline{W^*} \rightarrow$$

more in general  $w^* = w^\infty$  is a linear combination of support vectors thus  $w^{*T} x_n = 1$

$$H = - \sum_n^N x_n x_n^T e^{-W^{*T} x_n} - \lambda I$$

which is a stable min. ( $x_n x_n^T$  has zero eigenvalues).

This recovers the gist of Srebro's result for  $\lambda \rightarrow 0$ .

This is not linear but we linearize ad use HG

# ReLU and homogeneity properties of DNNs

Definition of ReLU implies

$$\sigma(z) = \frac{\partial \sigma(z)}{\partial z} z$$

and

$$f(W; x) = \prod_{k=1}^K f_k \tilde{f}(\tilde{W}_k^{ij}; x)$$

where  $f_k \tilde{W}_k^{ij} = W_k^{ij}$

and  $f_k = \|W_k\|$

Furthermore (Sasha et al.)

$$W_n^{ij} \frac{\partial f^{(k)}}{\partial W_n^{ij}} = f^{(k)} \quad \forall k$$



## Normalization

$$\tilde{w} = \frac{w}{\|w\|} \quad \|\tilde{w}\| = 1$$

$$1) \quad \frac{\partial \|w\|}{\partial w} = \tilde{w}$$

$$2) \quad \frac{\partial \tilde{w}}{\partial w} = \frac{I - \tilde{w} \tilde{w}^T}{\|w\|} = S$$

$$3) \quad S w = S \tilde{w} = 0$$

Let us reparametrize GD in terms of  $\tilde{w}, \rho$

(1 layer) starting from

$$\dot{w}^T = \sum_n^N \frac{\partial f(x_n)}{\partial w^T} e^{-f(x_n)}$$

This gives a new dynamical system

$$\|\dot{w}\| = \dot{\rho} = \frac{\partial \|w\|}{\partial w} \dot{w} = \tilde{w}^T \dot{w}$$

$$0 \quad \downarrow \text{tr}$$

$$\dot{\tilde{w}} = \frac{\partial \tilde{w}}{\partial w} \dot{w} = S \dot{w}$$

Thus

$$\dot{J} = \frac{\partial J}{\partial w} \sum_n e^{-\beta f} f = \sum_n e^{-\beta f} \dot{f}(x_n)$$

always giving

$$\dot{\tilde{w}} = \frac{\sum_n e^{-\beta f}}{Z} \left( \frac{\partial f}{\partial \tilde{w}} - \tilde{w} f(x_n) \right)$$

For linear networks

$$\dot{J} = \sum_n e^{-\beta \tilde{w}^T x_n} \tilde{w}^T x_n$$

$\rho = N$

$$\dot{\tilde{w}} = \frac{\sum_n e^{-\beta \tilde{w}^T x_n}}{\rho} \left( x_n - \tilde{w} \tilde{w}^T x_n \right)$$

↓

Define 
$$x = \sum_n \frac{e^{-\beta \tilde{w}^{*T} x_n}}{\mathcal{Z}} x_n$$

Then stationary point in  $w$  satisfies

$$x - \tilde{W}^* \tilde{W}^{*T} x = 0$$

$$\Rightarrow \tilde{W}^* = \frac{x}{\|x\|}$$

~~If~~ 
$$\sum_n e^{-\beta \tilde{w}^{*T} x_n} \approx \mathcal{Z}(x - x^*) \quad \text{then}$$

at minimum

$$H = (\mathbf{I} - \tilde{W} \tilde{W}^T) x x^T + \frac{1}{\beta} (\tilde{W}^T x + \tilde{W} x^T)$$

which is pos. def. unlike

However the argument depends on max effect in  $\sum_n e^{-\beta(\cdot)}$  which is not true in general!

# Summary

## o Square loss

Hyperbolic behaviour is guaranteed for 1 layer;

? can happen for multiple layers by

small  $\lambda$  regularization. In this case

KLG guarantees extension of linear analysis to nonlinear deep nets. Then solutions are

min norm solutions for  $\lambda \rightarrow 0$ .

## o Exponential loss

GD converges to hyperb solution - independent of initial conditions - in terms of normalized weights for linear networks (Leve)

We prove the same is true for regularized GD, not true for normalized GD and true for Early Stopped GD.