

# Lecture 14 Deep Learning Theory

IEMS 402 Statistical Learning

Northwestern

# References

[https://www.di.ens.fr/~fbach/ltfp\\_book.pdf](https://www.di.ens.fr/~fbach/ltfp_book.pdf)

- Section 12

objective function

$$\min_{\theta} \mathbb{E}_{x,y} \|NN_{\theta}(x) - y\|^2$$

is non-convex, because  $NN_{\theta}$  is highly non-linear

- why we can train the NN use GD

# Neural Tangent Kernel

- GD trained NN = kernel method.  
when NN is very wide,

# Neural Tangent Theory

Minimizing  $F(w) := R(h(w))$

↑ weight/parameter  
 ↑ risk NN

$h(w)$  is highly non-linear

feature

Consider a linearized model  $\bar{F}(w) := R(h(w_0) + \underbrace{\nabla_w h(w_0)}_{\text{gradient}}(w - w_0))$  → a linear regression.

initialization

Taylor expansion

↳ use a linear model to approximate the non-linear  $h(w)$

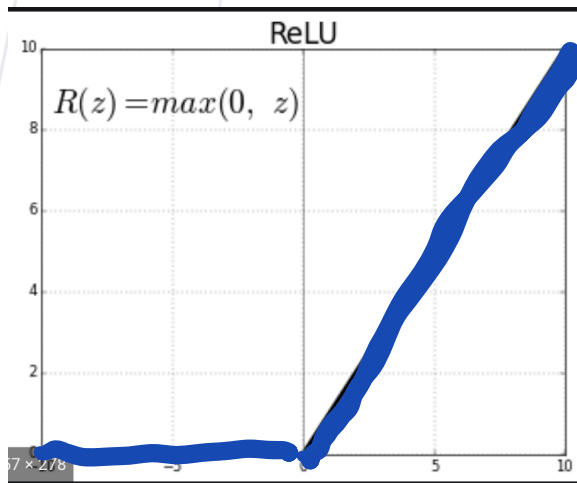
lazy training the less expected situation where these two paths remain close until the algorithm is stopped.

Neural tangent kernel:  $k(x, y) = \langle \nabla_w h(w_0)(x), \nabla_w h(w_0)(y) \rangle$

a vector same size as  $w$  only the data  $x$

Thm. When the NN is wide enough, then the approximation is good.

# Homogenous activation



$$\text{relu}(5x) = 5 \text{relu}(x)$$

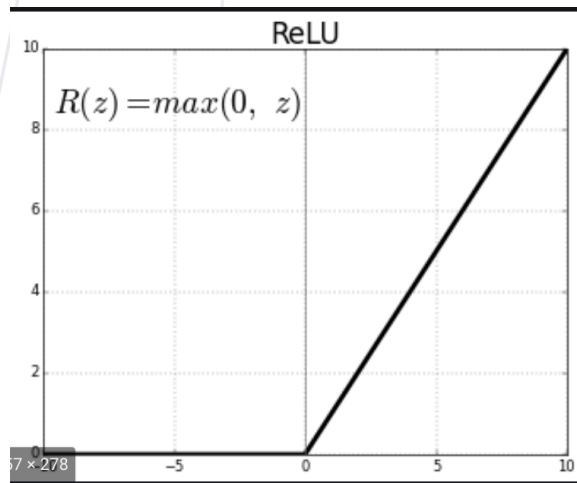
the grad of first layer is small

$$\underbrace{w_2}_{\tilde{w}_2} \text{relu}(\underbrace{5w_1}_{\tilde{w}_1} x) = \underbrace{5w_2}_{\tilde{w}_2} \text{relu}(w_1 x)$$

What's the thing different?  $\nabla_{\tilde{w}_1} \neq \nabla_{w_1}$ ,  $\nabla_{w_2} \neq \nabla_{\tilde{w}_2}$

gradient descent, depend on how you initialize the network!

# Homogenous activation



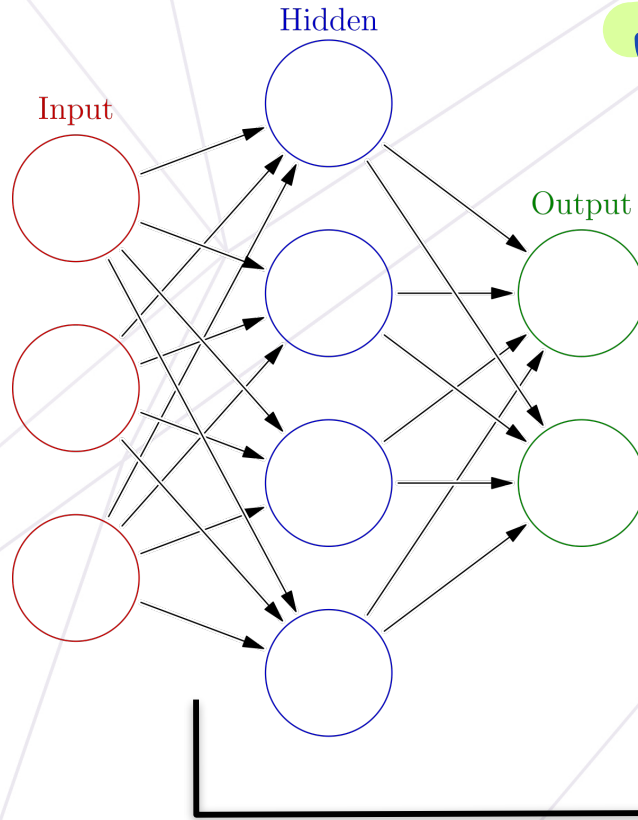
$$\text{relu}(5x) = 5\text{relu}(x)$$

$$\begin{array}{c} \text{Network 1} \\ w_2 \text{relu}(\underbrace{5w_1}_{\tilde{w}_1} x) \end{array} = \begin{array}{c} \text{Network 2} \\ \underbrace{(5w_2)}_{\tilde{w}_2} \text{relu}(w_1 x) \end{array}$$

What's the thing different?  $\nabla_{\tilde{w}_1} \neq \nabla_{w_1}, \nabla_{w_2} \neq \nabla_{\tilde{w}_2}$

*Is Adam/muon dynamics the same for two network?*

# Take Home Message

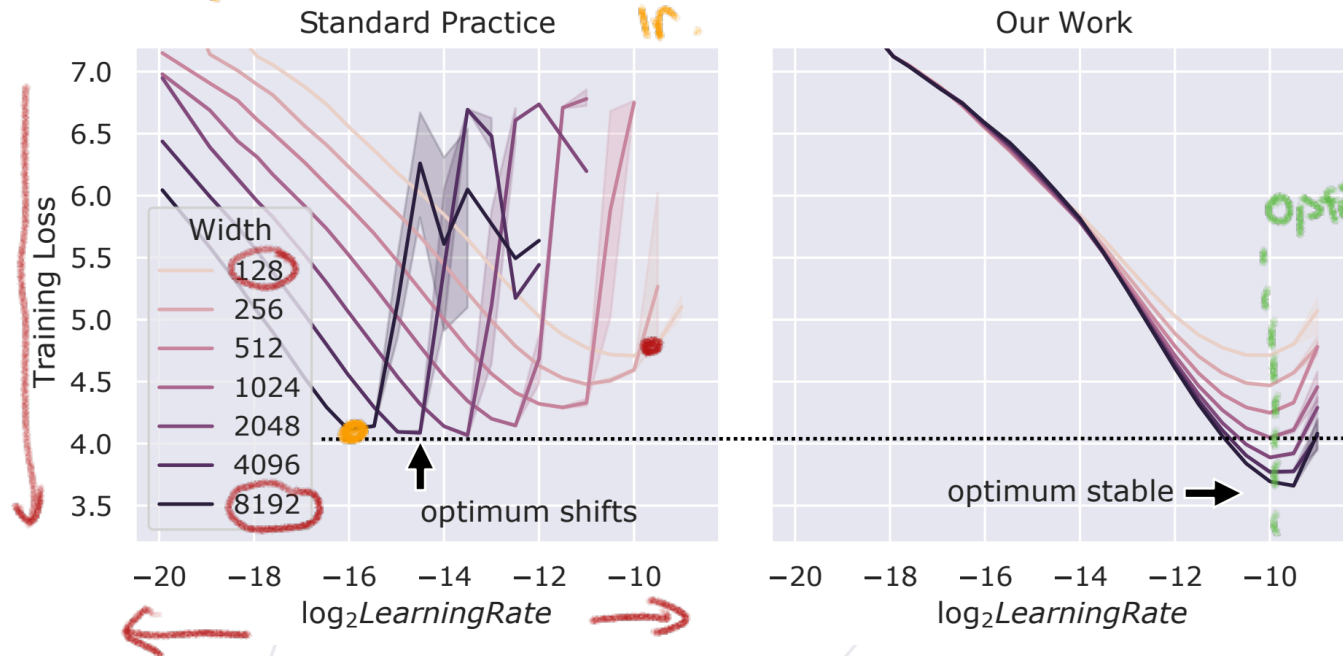


when NN is wider, the gradient of the first layer will become smaller (if you use standard initialization).

Wider network has smaller gradient on first layer

# Learning rate transfer

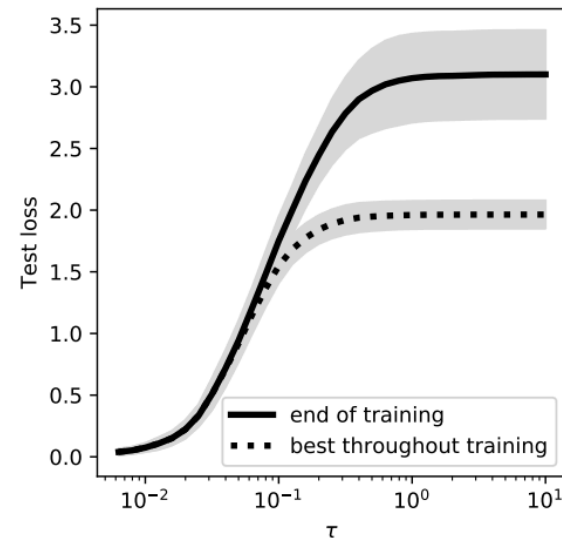
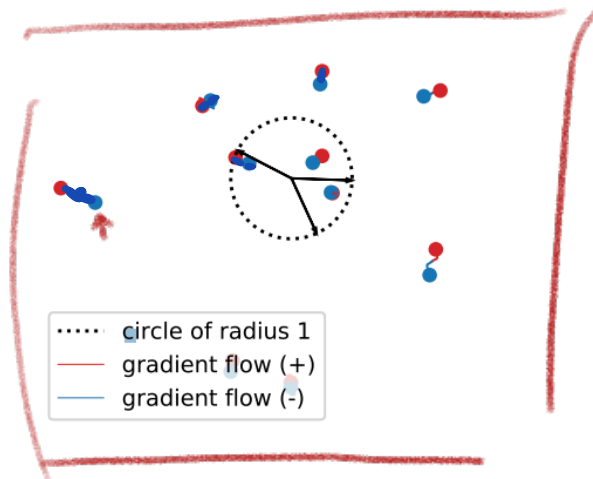
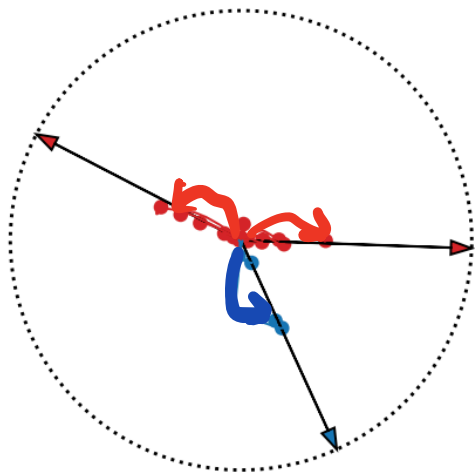
larger network have smaller lr.



<https://arxiv.org/pdf/2203.03466>



# Feature Learning and Lazy Learning



(a) Non-lazy training ( $\tau = 0.1$ )

(b) Lazy training ( $\tau = 2$ )

(c) Generalization properties

Chizat, Lenaic, Edouard Oyallon, and Francis Bach. "On lazy training in differentiable programming." Advances in neural information processing systems 32 (2019).

# When Lazy Training occurs?

$$F(w_1) = F(w_0 - \eta \nabla F) \approx F(w_0) - \nabla F(w_0) \cdot ((w_0 - \eta \nabla F) - w_0) = \eta \|\nabla F\|^2$$

Gradient descent  $w_1 := w_0 - \eta \nabla F(w_0)$ ,

$$\approx F(w_0) - \nabla F(w_0) \cdot ((w_0 - \eta \nabla F) - w_0)$$

Relative change of objective function

$$\Delta(F) := \frac{|F(w_1) - F(w_0)|}{F(w_0)} \approx \eta \frac{\|\nabla F(w_0)\|^2}{F(w_0)}$$

Relative change of linearization

$$\Delta(Dh) := \frac{\|Dh(w_1) - Dh(w_0)\|}{\|Dh(w_0)\|} \leq \eta \frac{\|\nabla F(w_0)\| \cdot \|D^2 h(w_0)\|}{\|Dh(w_0)\|}$$



$$\kappa_h(w_0) := \|h(w_0) - y^*\| \frac{\|D^2 h(w_0)\|}{\|Dh(w_0)\|^2} \ll 1,$$

condition of lazy training

# Example: Lazy Training For Homogeneous Model

**Homogeneous models.** If  $h$  is  $q$ -positively homogeneous<sup>4</sup> then multiplying the initialization by  $\lambda$  is equivalent to multiplying the scale factor  $\alpha$  by  $\lambda^q$ . In equation,

$$\kappa_h(\lambda w_0) = \frac{1}{\lambda^q} \|\lambda^q h(w_0) - y^*\| \frac{\|D^2 h(w_0)\|}{\|Dh(w_0)\|^2}.$$

$$\lambda \rightarrow \infty, \quad \kappa_h(w_0) \rightarrow 0,$$

finally lazy training appears.

# Mean Field Theory

feature learning for two-layer -

# Mean Field Theory

$$h(x) = \frac{1}{m} \sum_{j=1}^m \eta_j \sigma(w_j^\top x + b_j),$$

*individual neurons*

$$h = \frac{1}{m} \sum_{j=1}^m \Psi(v_j).$$

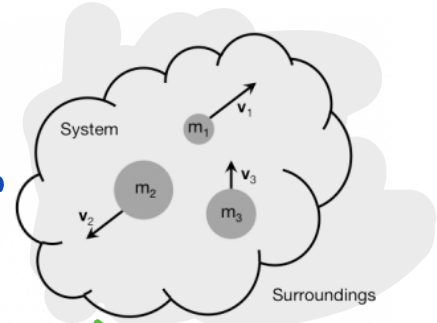
*average*  
*neuron*

Reformulate as probability distribution:

$$h = h(\cdot, v_1, \dots, v_m) = \int_v \Psi(v) d\mu(v),$$

*replace, average with expectation.*

*distribution of the weights.*



*-  $\mu(v) \rightarrow h$  is actually a linear mapping.*

*- distribution is in high-dimensional, hard to analyze*

# Gradient Flow in Wasserstein Space

Gradient descent:  $x_{t+1} = x_t - \alpha \nabla f(x_t)$

$$\Leftrightarrow x_{t+1} = \operatorname{argmin}_x \langle \nabla f, x - x_t \rangle + \|x - x_t\|^2$$

$x_{t+1} = \operatorname{argmin} \langle \nabla f, x - x_t \rangle + \|x - x_t\|_{\infty} \rightarrow \text{adam}$

$x_{t+1} \dots \dots \dots \| \dots \|_{\text{op}} \rightarrow \text{Muon}$

Change to different norms.

$\Rightarrow$  the state-of-the-art of for trans LLM.

# Gradient descent in weight = Gradient flow in Wasserstein space

We run gradient descent on  $\| \frac{1}{m} \sum_{i=1}^m \psi(V_i) - \kappa^* \|$  (1)

What is the trajectory of  $\| \int \psi(v) \mu_t - \kappa^* \|$   
the  $v$

$\mu_{t+1} =$  linear term + Optimal transport  $(\mu_{t+1}, \mu_t)$

HW2 .

$wx = y$  # data is smaller than # feature .

↳ infinite solution for  $w$  .

If you run gradient descent, it will converge to  $\min \|w\|$  .

s.t.  $wx = y$

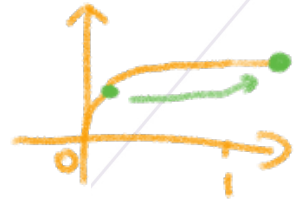
# Implicit Bias

Gradient Descent  
will select the data  
with minimum norm .



# Convergence in direction

$$F(\theta) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T \theta)),$$



$y_i x_i^T \theta \geq 0$ , I classify everything correct

$\theta \rightarrow 100\theta$  the loss will decrease.

↳ increase the confidence on the training data,

even if classification boundary is not changing!

Convergence In direction!

$$\nabla F(\theta) = \lambda \theta$$

# KKT condition for Largest Margin

$$\min \|\theta\|_2^2 \text{ subject to } y_i(\theta_i \cdot x) \geq 1$$

$\lambda_i = 0$ , when  $y_i(\theta_i \cdot x) > 1$   
 $\lambda_i \neq 0$  when  $y_i(\theta_i \cdot x) = 1 \Rightarrow$  support vector

$$\Rightarrow L = \|\theta\|_2^2 + \sum \lambda_i [y_i(\theta_i \cdot x) - 1]$$

classification confidence.

$$\Rightarrow \nabla_{\theta} L = 2\theta - \sum \lambda_i \nabla_{\theta} [y_i(\theta_i \cdot x)]$$
$$\Rightarrow \underline{2\theta} = \underline{\sum \lambda_i} \nabla_{\theta} [\text{classification confidence}]$$

# SVM=Logistic Regression

grad of ↑ confidence ·      gradient of logistic loss ·

$$G'(\theta) = \frac{-\sum_{i=1}^n y_i x_i \exp(-y_i x_i^\top \theta)}{\sum_{i=1}^n \exp(-y_i x_i^\top \theta)} = -\sum_{i=1}^n \alpha_i y_i x_i \rightarrow \text{gradient of the confidence}$$

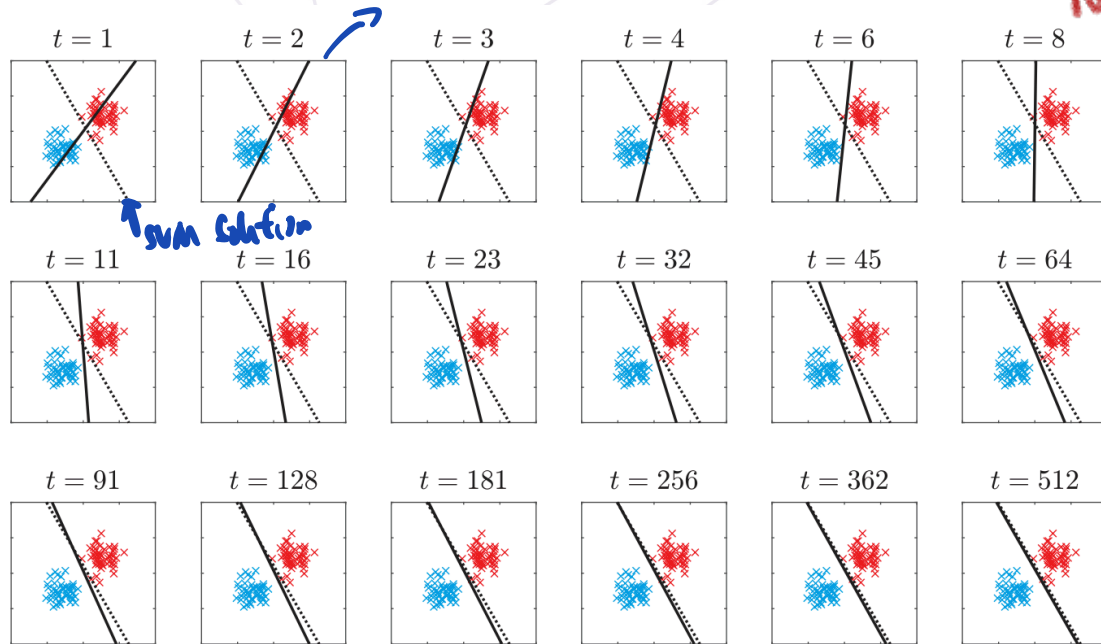
$$\nabla_{\theta} l(\text{confidence}(\theta)) = \nabla l \cdot \nabla_{\theta} \text{confidence}$$

Convergence in direction,  $\nabla F(\theta) = \lambda \theta$ . will change to.

$$\lambda \theta = \sum_{i=1}^n \alpha_i \text{grad}(\text{confidence})$$

# Converge to SVM solution

GD on logistic regression



SVM solution

holds  $\Downarrow$   
for all  
homogenous  
models

# Where is the support vectors

$$F'(\theta_t) \sim -\frac{1}{n} \sum_{i \in I} \exp(-\|\theta_t\|_2 y_i x_i^\top \eta) x_i \cdot y_i$$

$\eta = \frac{\theta_t}{\|\theta_t\|} \cdot \|\eta\| = 1$

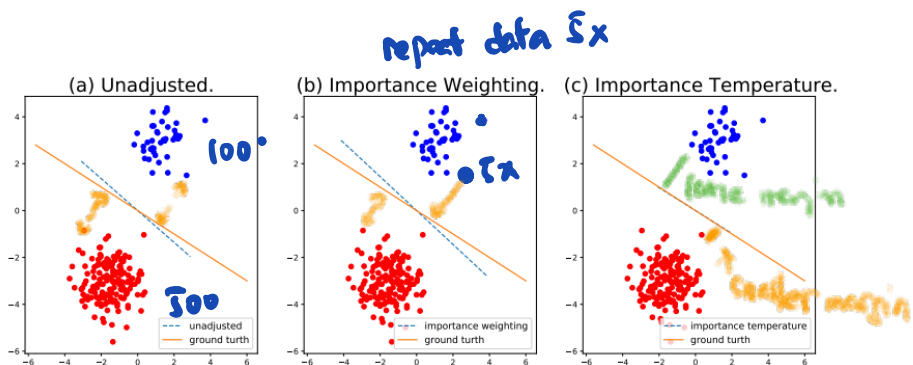
$x_i$       gradient of confidence.

$$\|\theta_t\| \rightarrow \infty$$

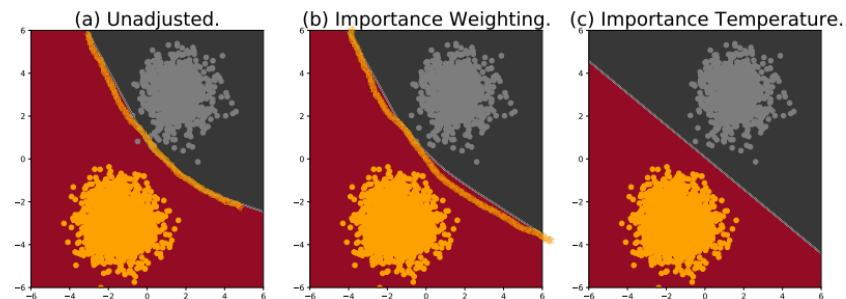
$$i = \operatorname{argmin}_i y_i x_i^\top \eta$$

$\exp(-\|\theta_t\|_2 y_i x_i^\top \eta)$  decays to 0 slowest  
 $\hookrightarrow$  dominates  $\rightarrow$  makes  $x_i$  support vector

# Failure of Importance Weighting



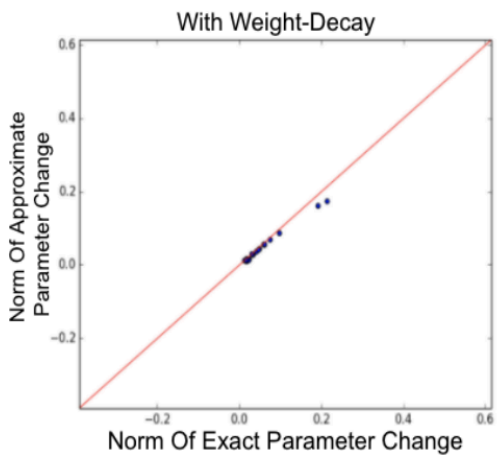
(a) Linear Model for Separable Data



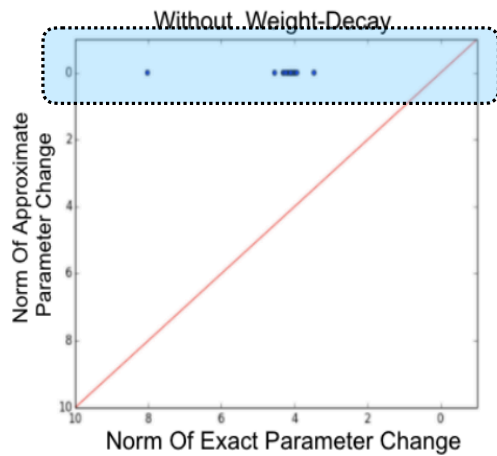
(b) Multilayer Perceptron with two hidden layers of size 200

Byrd J, Lipton Z. What is the effect of importance weighting in deep learning? International conference on machine learning. PMLR, 2019: 872-881.

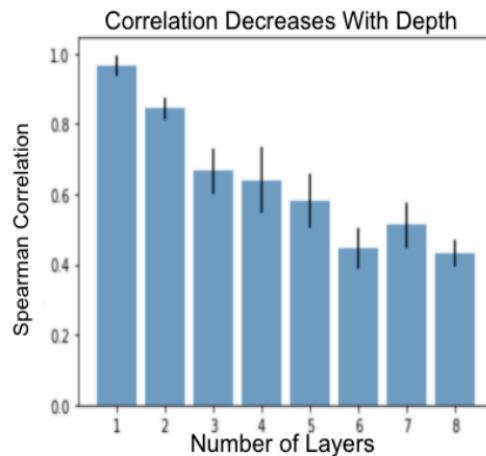
# Failure of Influence Function



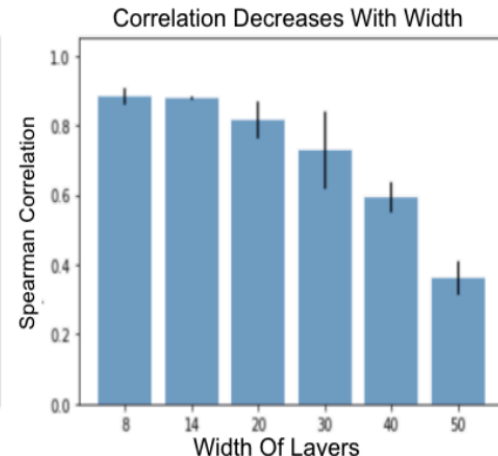
(a)



(b)



(c)



(d)

<https://arxiv.org/pdf/2006.14651>