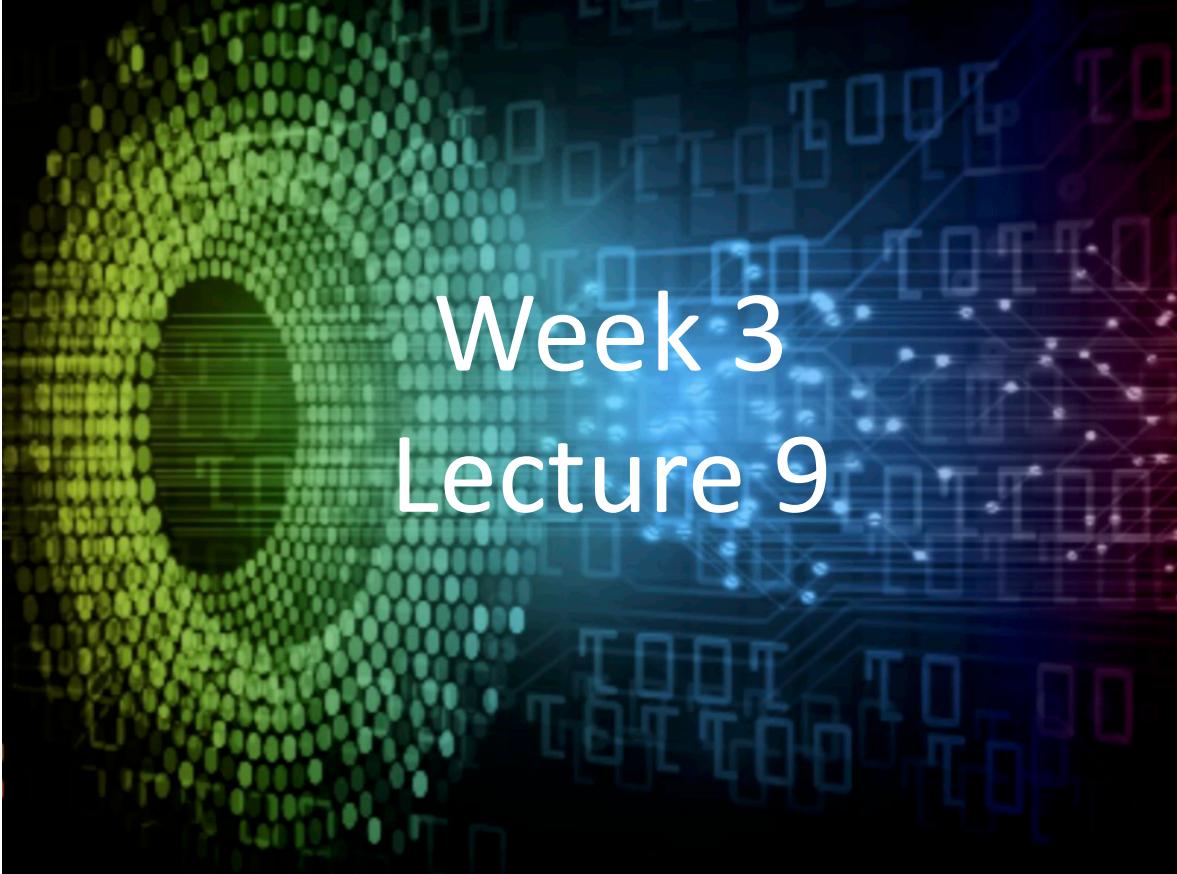


# Introduction to Deep Learning Applications and Theory



Week 3  
Lecture 9

ECE 596 / AMATH 563

# Deep Learning Practices

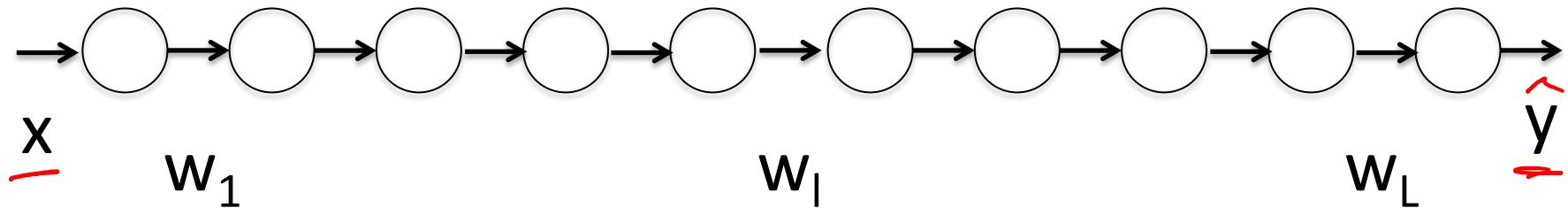
1. Previous Lecture: Regularization
2. Normalization
3. Batch Normalization
4. Initialization
5. Hyperparameters tuning
6. eScience



# Vanishing/Exploding Gradients

- Very deep neural network

$L \# \text{layers}$



$$\hat{y} = \underbrace{w_L \cdot \dots \cdot w_l \cdot \dots \cdot w_2}_{-} \cdot \underbrace{w_1}_{-} \cdot \underbrace{x}_{\uparrow}$$

$$w_l = \underline{w} > 1;$$

$$\hat{y} = w^L x \rightarrow \underline{\infty}$$

$$w_l = w < 1;$$

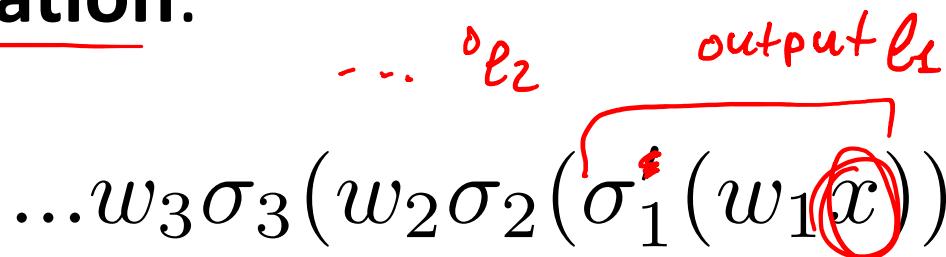
$$\hat{y} = \underline{w^L x} \rightarrow \underline{0}$$

# Vanishing/Exploding Gradients

- With activation:

$$\dots w_3 \sigma_3(w_2 \sigma_2(\sigma_1^*(w_1 x)))$$

$\dots \delta_{l_2}$       output  $l_1$



- For gradients:

$\sigma$

$wx$

$\tanh$

$wx$

$$\frac{\partial J}{\partial w_1} = \sigma'_3(z_3) \underbrace{w_3}_- \sigma'_2(z_2) \underbrace{w_2}_- \underbrace{\sigma'_1(z_1)}_x$$

# Normalization of the datasets

- Zero mean:

$$\underline{\mu} = \frac{1}{m} \sum_{i=1}^m \underline{x^{(i)}}$$

$$\underline{x^{(i)\mu}} = \underline{x^{(i)}} - \underline{\mu}$$

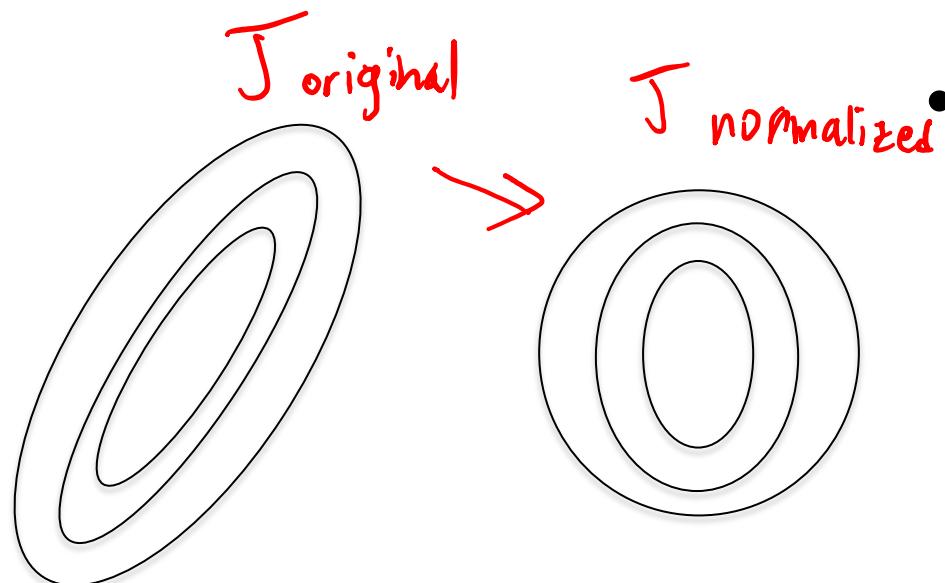
- Normalized Variance:

$$\underline{\sigma^2} = \frac{1}{m} \sum_{i=1}^m x^{(i)2}$$

$$x^{(i)\mu, \sigma^2} = x^{(i)\mu} ./ \sigma^2$$

# Intuition

- If inputs have different scales, the **cost function** will also have to include different scales

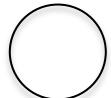


- Remember to normalize all sets: training, validation, testing

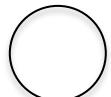
# Batch Normalization

- Normalize the outputs of each layer

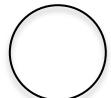
$$\underline{z}^{[l]} = \underline{w}^{[l]} \underline{a}^{[l-1]} + \underline{b}^{[l]}$$



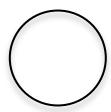
$$\underline{\mu} = \frac{1}{m} \sum_{i=1}^m z^{(i)}$$



$$\underline{\sigma^2} = \frac{1}{m} \sum_{i=1}^m (z^{(i)} - \mu)^2$$



$$z_{norm}^{(i)} = \frac{z^{(i)} - \underline{\mu}}{\sqrt{\underline{\sigma^2} + \epsilon}}$$



# Learning batch normalization

- To make sure that the layers outputs will not be forced to be zero mean and variance 1

$$\tilde{z}^{[l]}(i) = \gamma z_{norm}^{[l]} + \beta$$

- Forward propagation

$$a^{[l-1]}(i) \rightarrow z^{[l]}(i) \xrightarrow{BN} \tilde{z}^{[l]}(i) \rightarrow a^{[l]}(i)$$

- Can do per minibatch

# Pseudo-code

For each mini-batch

compute F-prop:

in each layer replace                                    by

Do B-prop to compute

update

Testing:

$$z[l](i)$$



$$\tilde{z}[l](i)$$

$$\delta, \beta$$

EMA

$$\nabla_{w[l]}, \nabla_{\beta[l]}, \nabla_{\gamma[l]}$$

# Initialization

- Zero → Problematic
- Random Normal (0,1) -> Problematic
- Xavier (tanh):

$$Var(w^{[l]}) : 1/n^{[l-1]}$$
$$w^{[l]} = N(0, 1) \cdot \sqrt{\frac{1}{n^{[l-1]}}}$$

# Initialization

- He (ReLU):

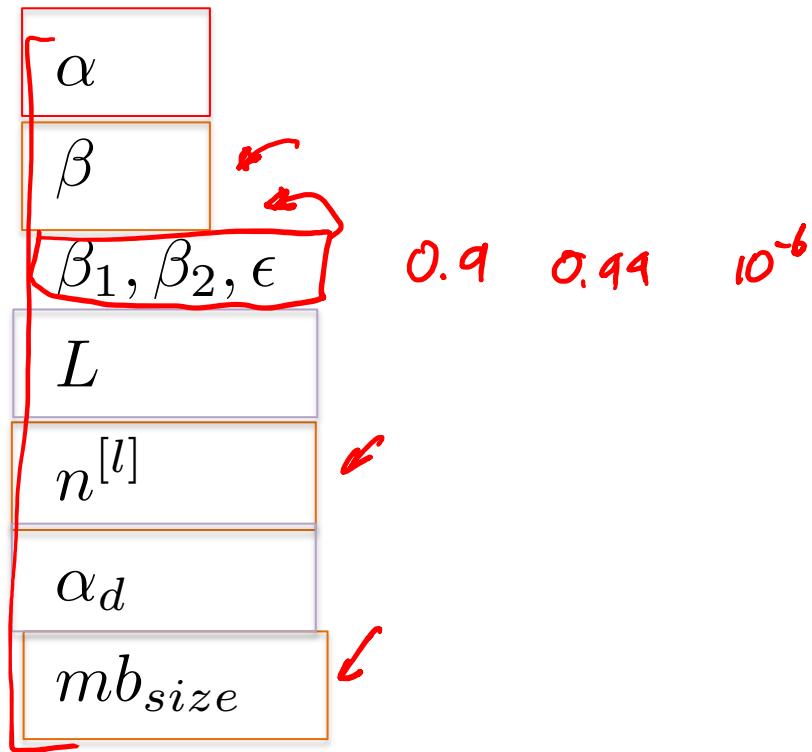
$$Var(w^{[l]}) : \frac{2}{n^{[l-1]}}$$
$$w^{[l]} = N(0, 1) \cdot \sqrt{\frac{2}{n^{[l-1]}}}$$

- Other:

$$Var(w^{[l]}) : \frac{2}{\underline{n^{[l-1]}} + \underline{n^{[l]}}}$$

# Hyper Parameters

Params:

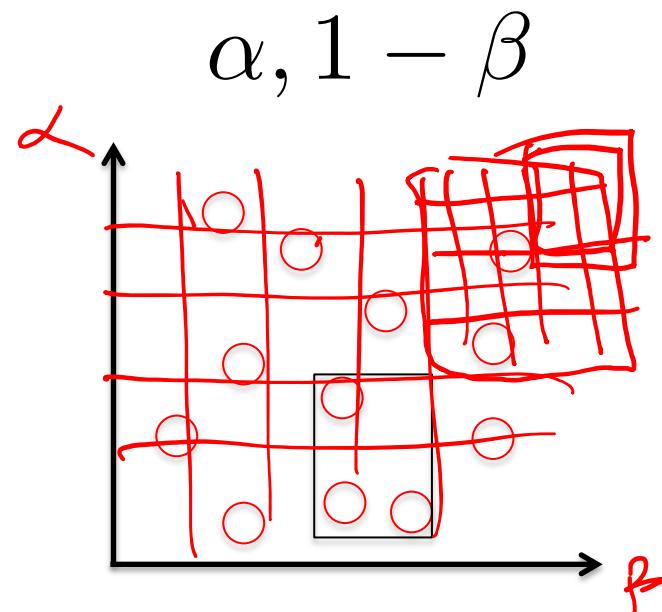


$$\lambda$$

$$P_d$$

# Search

- Grid Search, Random Search
- Coarse to fine
- Linear for : L,  $n^{[l]}$
- Log<sub>10</sub> scale for



# Data Science

- UW eScience lecture



**Sarah Stone**  
eScience Executive Director

Weds  
11:30

