

## Lecture 10

# Feature Normalization and Weight Initialization

STAT 453: Deep Learning, Spring 2020

Sebastian Raschka

<http://stat.wisc.edu/~sraschka/teaching/stat453-ss2020/>

Slides:

[https://github.com/rasbt/stat453-deep-learning-ss20/10\\_norm-and-init/](https://github.com/rasbt/stat453-deep-learning-ss20/10_norm-and-init/)

# "Tricks" for Improving Deep Neural Network Training

## Today:

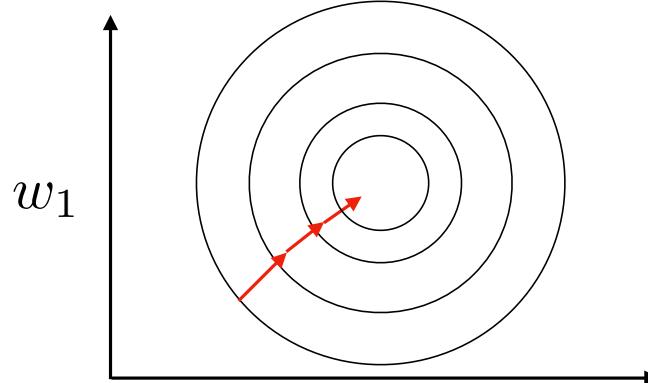
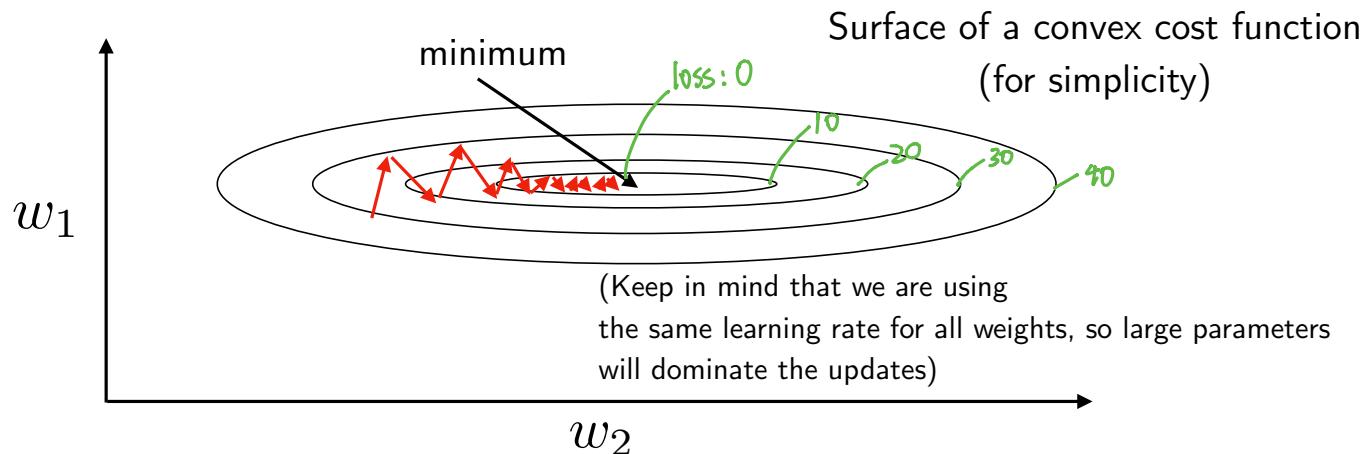
1. Feature/Input Normalization  
(BatchNorm, InstanceNorm, GroupNorm, LayerNorm)
2. Weight Initialization (Xavier Glorot, Kaiming He)

## Next Lecture:

3. Optimization Algorithms (RMSProp, Adagrad, ADAM)

# Part 1: Input Normalization

# Recap: Why We Normalize Inputs for Gradient Descent



"Standardization" of input features

$$x_j' [i] = \frac{x_j^{[i]} - \mu_j}{\sigma_j}$$

*j<sup>th</sup> feature*      *mean*  
                          *s.d.*

(scaled feature will have zero mean, unit variance)

We do normalization when  $w_2$  feature have different ranges

If we always standardize, we will always lose information

However, normalizing  
the inputs to the network  
only affects the first hidden layer ...  
What about the other hidden layers?

# Batch Normalization ("BatchNorm")

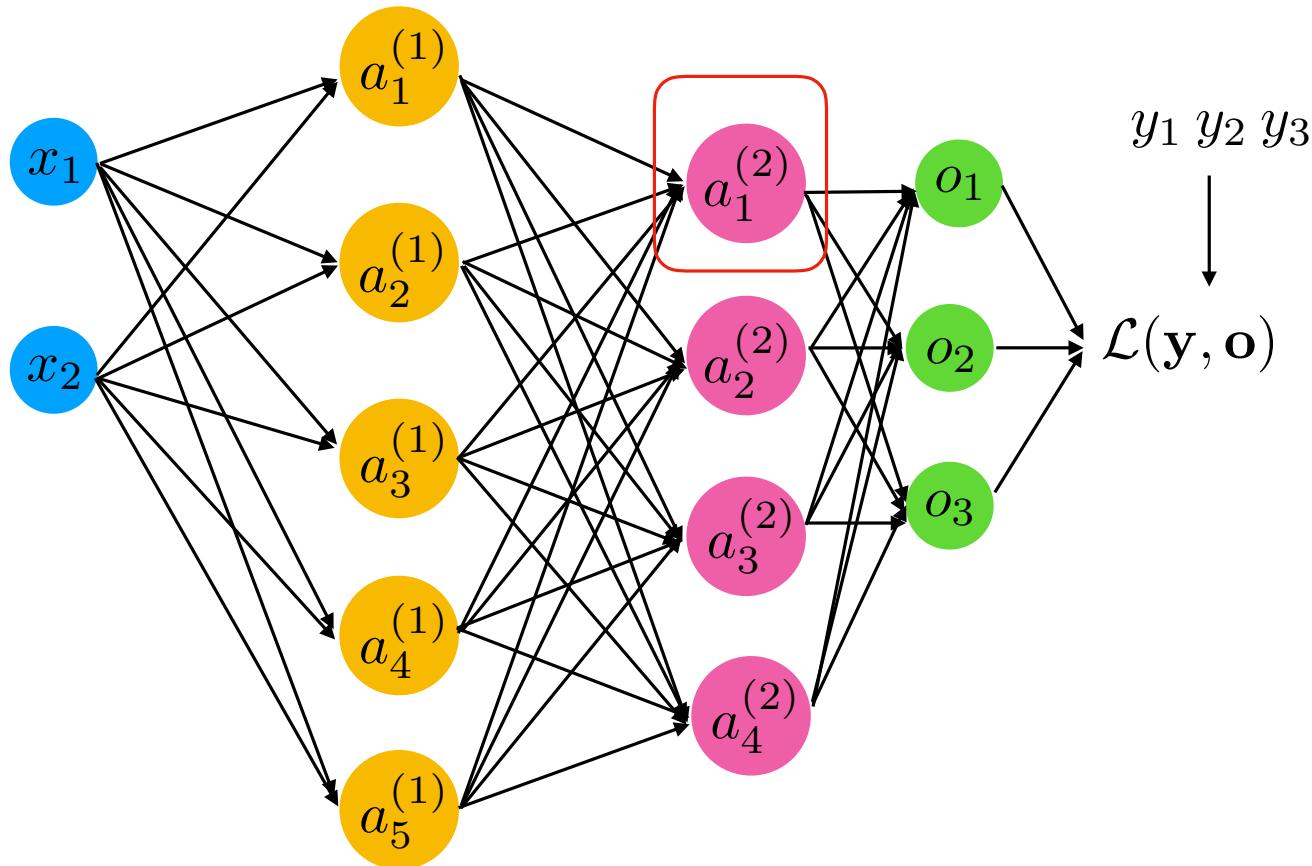
Ioffe, S., & Szegedy, C. (2015, June). Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In *International Conference on Machine Learning* (pp. 448-456).

<http://proceedings.mlr.press/v37/ioffe15.html>

# Batch Normalization ("BatchNorm")

- Normalizes hidden layer inputs
- Helps with exploding/vanishing gradient problems
- Can increase training stability and convergence rate
- Can be understood as additional (normalization) layers  
(with additional parameters)

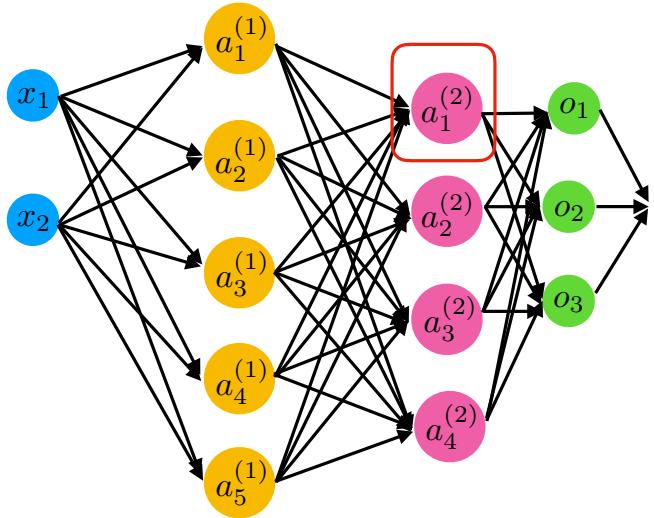
Suppose, we have net input  $z_1^{(2)}$   
associated with an activation in the 2nd hidden layer



Now, consider all examples in a minibatch such that the net input of a given training example at layer 2 is written as  $z_1^{(2)[i]}$

layer index  
minibatch size

where  $i \in \{1, \dots, n\}$



In the next slides, let's omit the layer index, as it may be distracting...

# BatchNorm Step 1: Normalize Net Inputs

sample mean

$$\underline{\mu_j} = \frac{1}{n} \sum_i z_j^{[i]}$$

z-score standardization

do it every layer separately  
for

variance

$$\sigma_j^2 = \frac{1}{n} \sum_i (z_j^{[i]} - \mu_j)^2$$

$$z'_j^{[i]} = \frac{z_j^{[i]} - \mu_j}{\sigma_j}$$

standardize

$$\mu = 0$$

$$\sigma^2 = 1$$

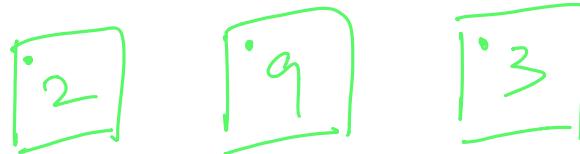
# BatchNorm Step 1: Normalize Net Inputs

$$\mu_j = \frac{1}{n} \sum_i z_j^{[i]}$$

$$\sigma_j^2 = \frac{1}{n} \sum_i (z_j^{[i]} - \mu_j)^2$$

$$z'_j^{[i]} = \frac{z_j^{[i]} - \mu_j}{\sigma_j + \epsilon}$$

Division by  
Zero error



In practice:

$$z'_j^{[i]} = \frac{z_j^{[i]} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

For numerical stability, where epsilon  
is a small number like 1E-5

# BatchNorm Step 2: Pre-Activation Scaling

$$z'_j^{[i]} = \frac{z_j^{[i]} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

$$a'_j^{[i]} = \gamma_j \cdot z'_j^{[i]} + \beta_j \quad \text{Step 2}$$

These are **learnable** parameters

Scale our  $z_j$

Learn by  $\frac{\partial L}{\partial \gamma_j}$        $\frac{\partial L}{\partial \beta_j}$

# BatchNorm Step 2: Pre-Activation Scaling

$$z_j'^{[i]} = \frac{z_j^{[i]} - \mu_j}{\sigma_j}$$

$$a_j'^{[i]} = \gamma_j \cdot z_j'^{[i]} + \beta_j$$

Controls the spread or scale

Controls the mean

It gives the opportunity to undo the  
Standardization

# BatchNorm Step 2: Pre-Activation Scaling

$$z_j'^{[i]} = \frac{z_j^{[i]} - \mu_j}{\sigma_j}$$

$$a_j'^{[i]} = \gamma_j \cdot z_j'^{[i]} + \beta_j$$

Controls the spread or scale

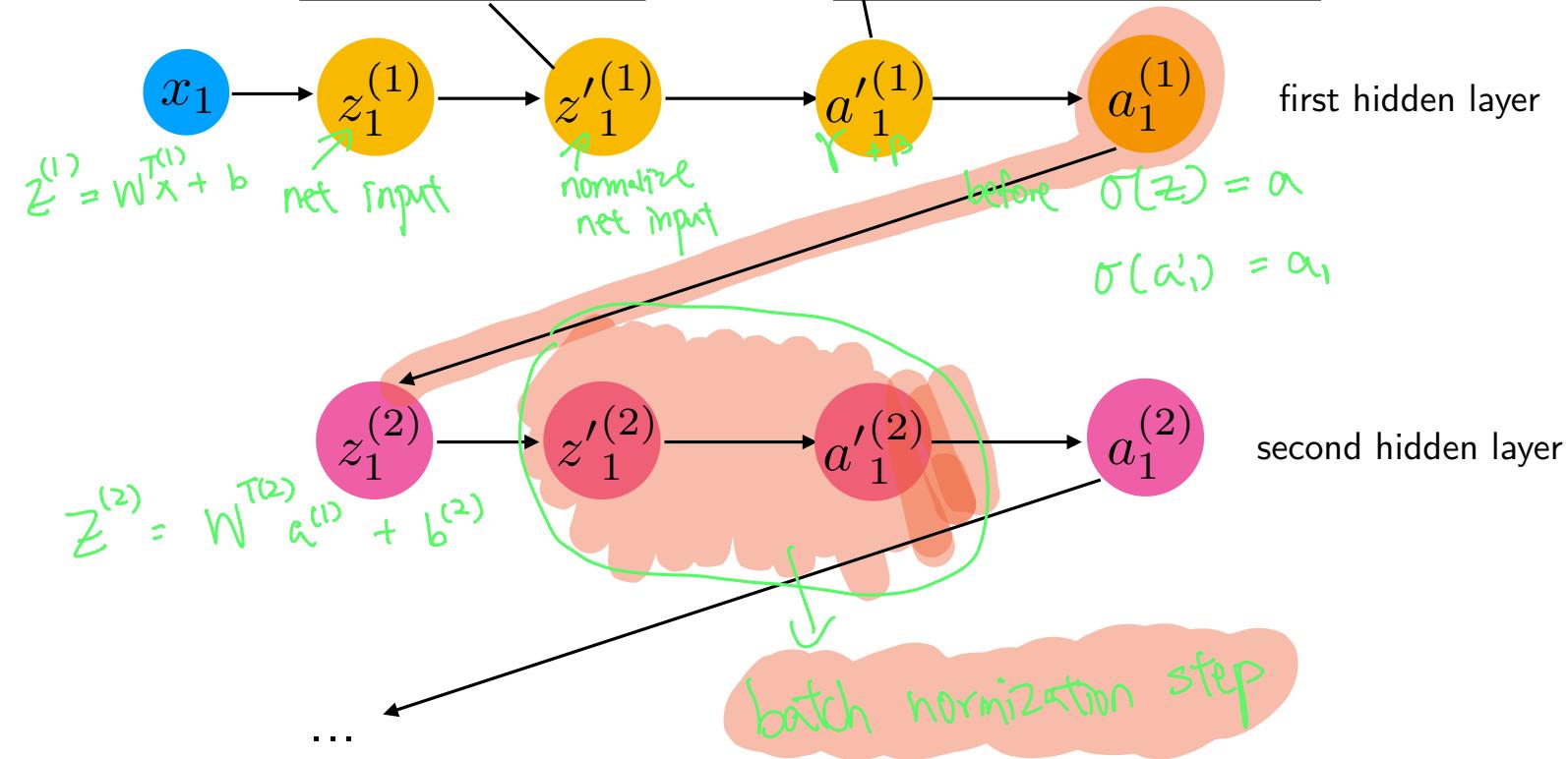
Controls the mean

Technically, a BatchNorm layer could learn to perform "standardization" with zero mean and unit variance

# BatchNorm Step 1 & 2 Summarized

$$z'_j^{[i]} = \frac{z_j^{[i]} - \mu_j}{\sigma_j}$$

$$a'_j^{[i]} = \gamma_j \cdot z'_j^{[i]} + \beta_j$$

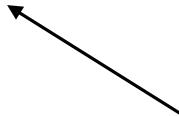


# BatchNorm -- Additional Things to Consider

$$a_j'^{[i]} = \gamma_j \cdot z_j'^{[i]} + \beta_j$$

$$z = W^T x + b + \beta_j$$

-  
無用, useless

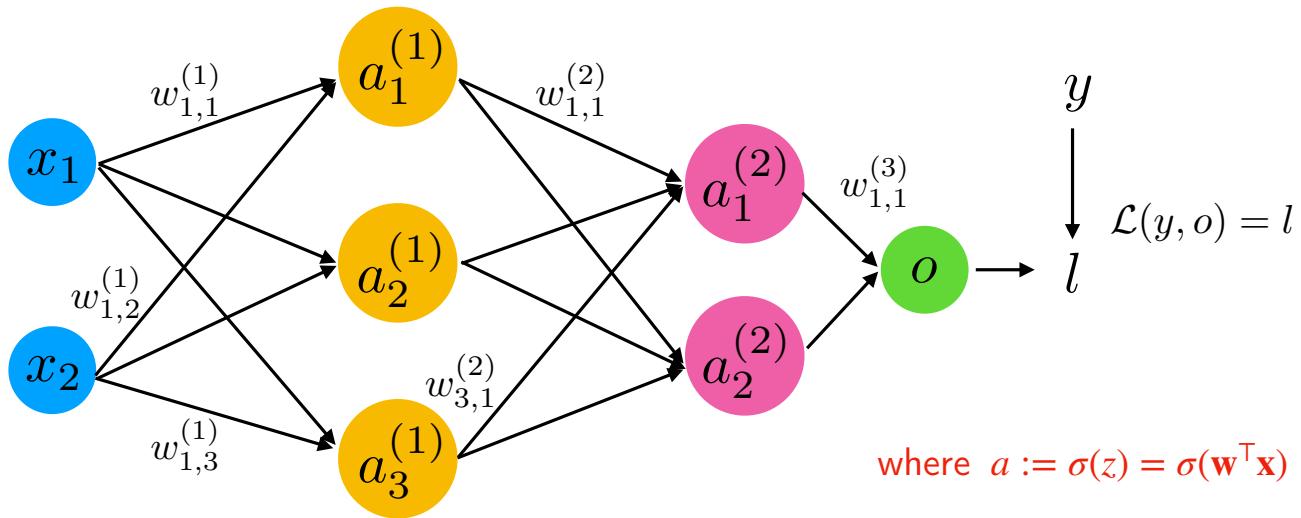


This parameter makes the bias units redundant

Also, note that the batchnorm parameters are vectors with the same number of elements as the bias vector

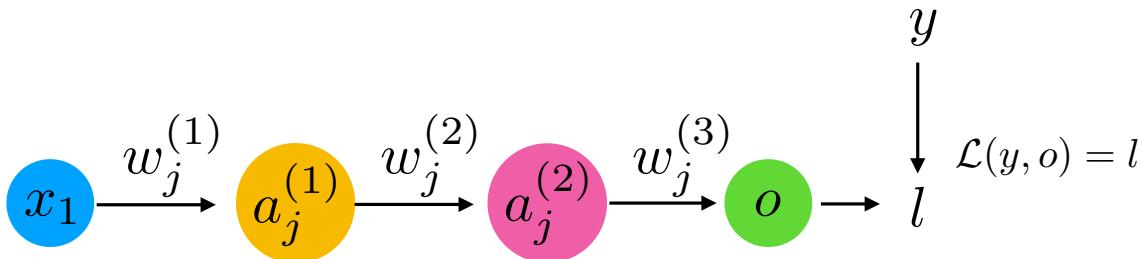
# **Backpropagation for BatchNorm Parameters**

# Reminder: Multilayer Perceptron (from Lecture 9)



$$\begin{aligned}\frac{\partial l}{\partial w_{1,1}^{(1)}} &= \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial a_1^{(1)}} \cdot \frac{\partial a_1^{(1)}}{\partial w_{1,1}^{(1)}} \\ &\quad + \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_2^{(2)}} \cdot \frac{\partial a_2^{(2)}}{\partial a_1^{(1)}} \cdot \frac{\partial a_1^{(1)}}{\partial w_{1,1}^{(1)}}\end{aligned}$$

# Let's consider a simpler case ...

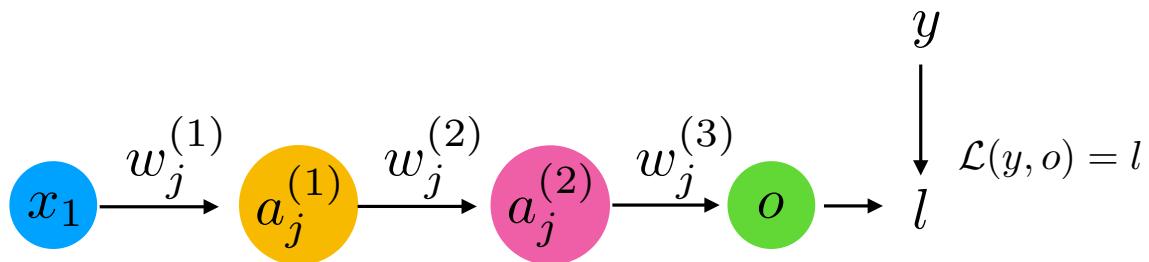


$$\frac{\partial l}{\partial w_j^{(3)}} = \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial w_j^{(3)}}$$

$$\frac{\partial l}{\partial w_j^{(2)}} = \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_j^{(2)}} \cdot \frac{\partial a_j^{(2)}}{\partial w_j^{(2)}}$$

$$\frac{\partial l}{\partial w_j^{(1)}} = \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_j^{(2)}} \cdot \frac{\partial a_j^{(2)}}{\partial a_j^{(1)}} \cdot \frac{\partial a_j^{(1)}}{\partial w_j^{(1)}}$$

# Same as on previous slide, but more verbose ...



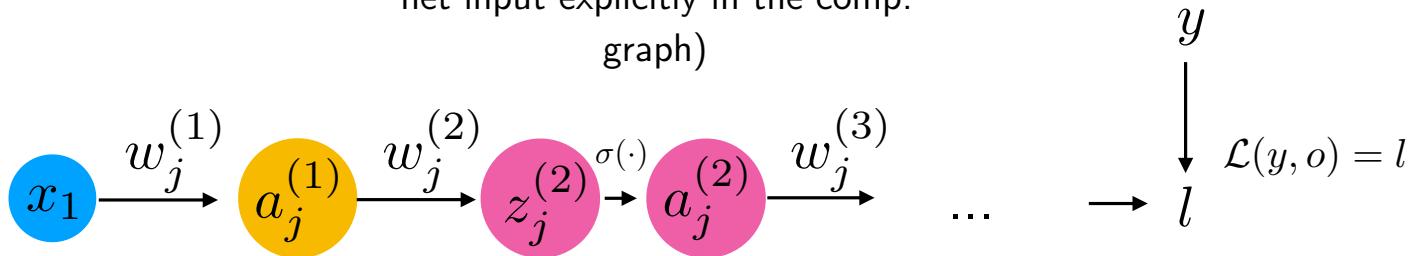
$$\frac{\partial l}{\partial w_j^{(3)}} = \frac{\partial l}{\partial o} \cdot \underbrace{\frac{\partial o}{\partial w_j^{(3)}}}_{\text{circled}} \longrightarrow \frac{\partial l}{\partial w_j^{(3)}} = \frac{\partial l}{\partial o} \cdot \underbrace{\frac{\partial o}{\partial z_j^{(3)}}}_{\text{circled}} \cdot \frac{\partial z_j^{(3)}}{\partial w_j^{(3)}}$$

$$\frac{\partial l}{\partial w_j^{(2)}} = \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_j^{(2)}} \cdot \frac{\partial a_j^{(2)}}{\partial w_j^{(2)}} \longrightarrow \dots$$

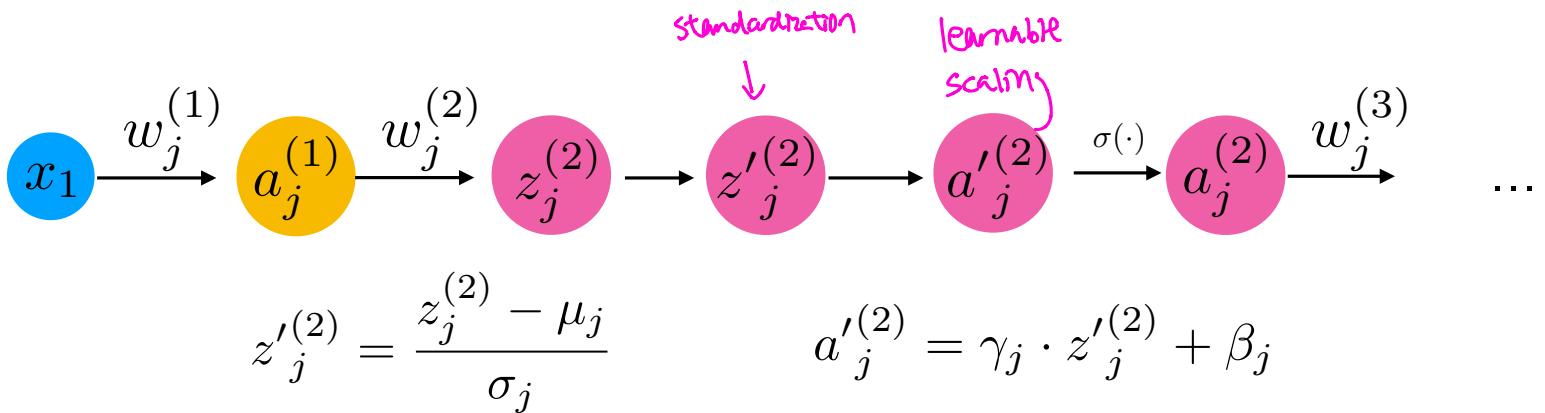
$$\frac{\partial l}{\partial w_j^{(1)}} = \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_j^{(2)}} \cdot \frac{\partial a_j^{(2)}}{\partial a_j^{(1)}} \cdot \frac{\partial a_j^{(1)}}{\partial w_j^{(1)}} \longrightarrow \dots$$

$\therefore z = w^T x + b$   
 $\uparrow$   
 $a^{(j-1)}$   
 $\partial z = x$

(previously, we didn't write the net input explicitly in the comp. graph)

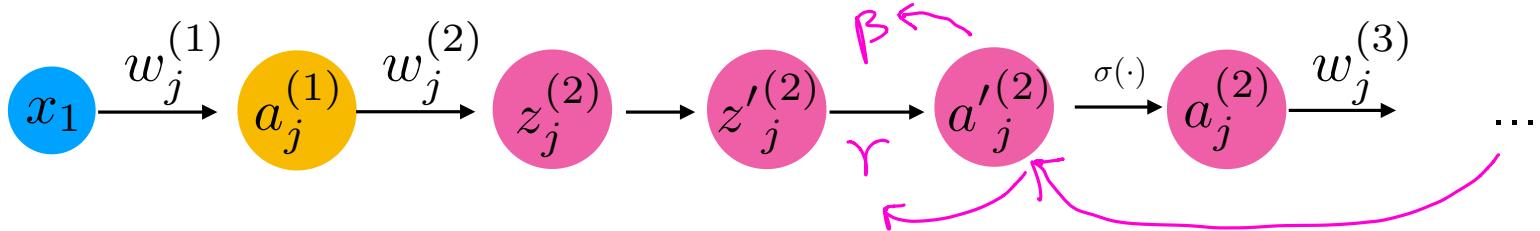


↓  
Adding a  
BatchNorm layer ...



# Backprop for BatchNorm Parameters

$$z_j^{(2)} = \frac{z_j^{(2)} - \mu_j}{\sigma_j} \quad a_j'^{(2)} = \gamma_j \cdot z_j'^{(2)} + \beta_j$$



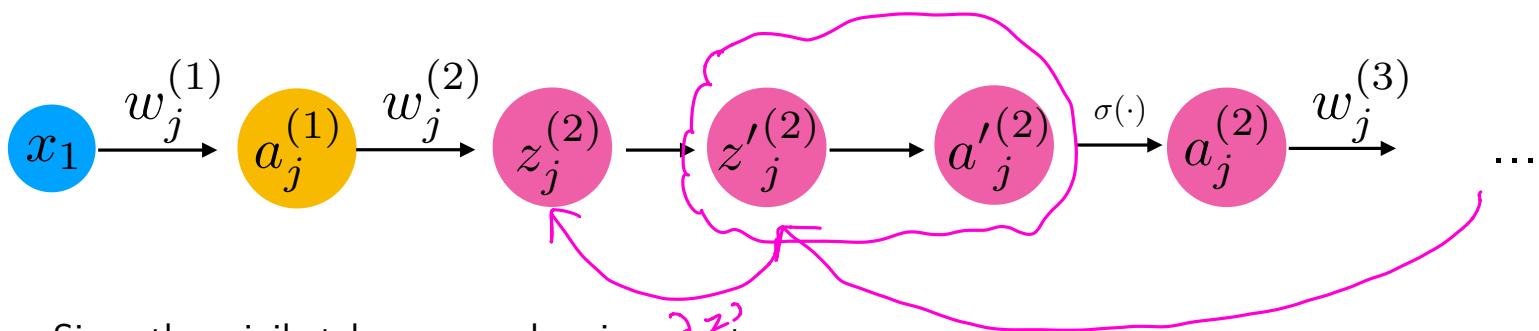
$$\frac{\partial l}{\partial \beta_j} = \sum_{i=1}^n \frac{\partial l}{\partial a_j'^{(2)[i]}} \cdot \cancel{\frac{\partial a_j'^{(2)[i]}}{\partial \beta_j}}$$

# in minibatch

*4x + 5*

$$\frac{\partial l}{\partial \gamma_j} = \sum_{i=1}^n \frac{\partial l}{\partial a_j'^{(2)[i]}} \cdot \frac{\partial a_j'^{(2)[i]}}{\partial \gamma_j} = \sum_{i=1}^n \frac{\partial l}{\partial a_j'^{(2)[i]}} \cdot z_j'^{(2)[i]}$$

# Backprop Beyond the BatchNorm Layer



Since the minibatch mean and variance act as parameters, we can/have to apply the multivariable chain rule

$$\begin{aligned} \frac{\partial l}{\partial z_j^{(2)[i]}} &= \left\{ \frac{\partial l}{\partial z'^{(2)[i]}} \cdot \frac{\partial z'^{(2)[i]}}{\partial z_j^{(2)[i]}} + \frac{\partial l}{\partial \mu_j} \cdot \frac{\partial \mu_j}{\partial z_j^{(2)[i]}} + \frac{\partial l}{\partial \sigma_j^2} \cdot \frac{\partial \sigma_j^2}{\partial z_j^{(2)[i]}} \right\} \\ &= \frac{\partial l}{\partial z'^{(2)[i]}} \cdot \frac{1}{\sigma_j} + \frac{\partial l}{\partial \mu_j} \cdot \frac{1}{n} + \frac{\partial l}{\partial \sigma_j^2} \cdot \frac{2(z_j^{(2)} - \mu_j)}{n} \end{aligned}$$

$\frac{\partial a^2}{\partial z'} = r$

# Backprop for BatchNorm Parameters

$$\begin{aligned}\frac{\partial l}{\partial z_j^{(2)[i]}} &= \frac{\partial l}{\partial z'_j^{(2)[i]}} \cdot \frac{\partial z'_j^{(2)[i]}}{\partial z_j^{(2)[i]}} + \frac{\partial l}{\partial \mu_j} \cdot \frac{\partial \mu_j}{\partial z_j^{(2)[i]}} + \frac{\partial l}{\partial \sigma_j^2} \cdot \frac{\partial \sigma_j^2}{\partial z_j^{(2)[i]}} \\ &= \boxed{\frac{\partial l}{\partial z'_j^{(2)[i]}} \cdot \frac{1}{\sigma_j}} + \boxed{\frac{\partial l}{\partial \mu_j} \cdot \frac{1}{n}} + \boxed{\frac{\partial l}{\partial \sigma_j^2} \cdot \frac{2(z_j^{(2)} - \mu_j)}{n}}\end{aligned}$$

Annotations:

- $\bar{z} = \frac{z - \mu}{\sigma}$
- $\frac{\partial \bar{z}}{\partial z} = \frac{1}{\sigma}$
- $\frac{1}{n} \sum (z - \mu)^2$

If you like math & engineering, you can solve the remaining terms  
as an ungraded HW exercise ;)

# BatchNorm in PyTorch

```
class MultilayerPerceptron(torch.nn.Module):  
  
    def __init__(self, num_features, num_classes):  
        super(MultilayerPerceptron, self).__init__()  
  
        ### 1st hidden layer  
        self.linear_1 = torch.nn.Linear(num_features, num_hidden_1)  
        self.linear_1_bn = torch.nn.BatchNorm1d(num_hidden_1) → Batch norm layer  
                                         vector  
        ### 2nd hidden layer  
        self.linear_2 = torch.nn.Linear(num_hidden_1, num_hidden_2)  
        self.linear_2_bn = torch.nn.BatchNorm1d(num_hidden_2)  
  
        ### Output layer  
        self.linear_out = torch.nn.Linear(num_hidden_2, num_classes)  
  
    def forward(self, x):  
        out = self.linear_1(x)  
        # note that batchnorm is in the classic  
        # sense placed before the activation  
        out = self.linear_1_bn(out) ←  
        out = F.relu(out)  
  
        out = self.linear_2(out)  
        out = self.linear_2_bn(out)  
        out = F.relu(out)  
  
        logits = self.linear_out(out)  
        probas = F.softmax(logits, dim=1)  
        return logits, probas
```

# BatchNorm in PyTorch

```
class MultilayerPerceptron(torch.nn.Module):  
  
    def __init__(self, num_features, num_classes):  
        super(MultilayerPerceptron, self).__init__()  
  
        ### 1st hidden layer  
        self.linear_1 = torch.nn.Linear(num_features, num_hidden_1)  
        self.linear_1_bn = torch.nn.BatchNorm1d(num_hidden_1)  
  
        ### 2nd hidden layer  
        self.linear_2 = torch.nn.Linear(num_hidden_1, num_hidden_2)  
        self.linear_2_bn = torch.nn.BatchNorm1d(num_hidden_2)  
  
        ### Output layer  
        self.linear_out = torch.nn.Linear(num_hidden_2, num_classes)  
  
    def forward(self, x):  
        out = self.linear_1(x)  
        # note that batchnorm is in the classic  
        # sense placed before the activation  
        out = self.linear_1_bn(out)  
        out = F.relu(out)  
  
        out = self.linear_2(out)  
        out = self.linear_2_bn(out)  
        out = F.relu(out)  
  
        logits = self.linear_out(out)  
        probas = F.softmax(logits, dim=1)  
        return logits, probas
```

don't forget `model.train()`  
and `model.eval()`  
in training and test loops

# BatchNorm During Prediction ("Inference")

- Use exponentially weighted average (moving average) of mean and variance

*If testing has a data point, can't compute mean*

*running\_mean = momentum \* running\_mean + (1 - momentum) \* sample\_mean*

*running\_mean* *saved for testing during training* *updated* *For current mini batch*

*sample\_mean*

(where momentum is typically ~0.1; and same for variance)

- Alternatively, can also use global training set mean and variance

# BatchNorm and Internal Covariate Shift

Ioffe, S., & Szegedy, C. (2015, June). Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In *International Conference on Machine Learning* (pp. 448-456).

<http://proceedings.mlr.press/v37/ioffe15.html>

Internal Covariate Shift is jargon for saying that the layer input distribution changes ("feature shift" in hidden layers )

But there is actually no guarantee or evidence that BatchNorm helps with covariate shift

In my opinion, BatchNorm just provides additional parameters that will help layers to learn a little bit more independently  $\gamma, \beta$

---

# How Does Batch Normalization Help Optimization?

---

**Shibani Santurkar\***

MIT

shibani@mit.edu

**Dimitris Tsipras\***

MIT

tsipras@mit.edu

**Andrew Ilyas\***

MIT

ailyas@mit.edu

**Aleksander Mądry**

MIT

madry@mit.edu

## Abstract

Batch Normalization (BatchNorm) is a widely adopted technique that enables faster and more stable training of deep neural networks (DNNs). Despite its pervasiveness, the exact reasons for BatchNorm’s effectiveness are still poorly understood. The popular belief is that this effectiveness stems from controlling the change of the layers’ input distributions during training to reduce the so-called “internal covariate shift”. In this work, we demonstrate that such distributional stability of layer inputs has little to do with the success of BatchNorm. Instead, we uncover a more fundamental impact of BatchNorm on the training process: it makes the optimization landscape significantly smoother. This smoothness induces a more predictive and stable behavior of the gradients, allowing for faster training.

# BatchNorm Enables Faster Convergence By

① Allowing Larger Learning Rates

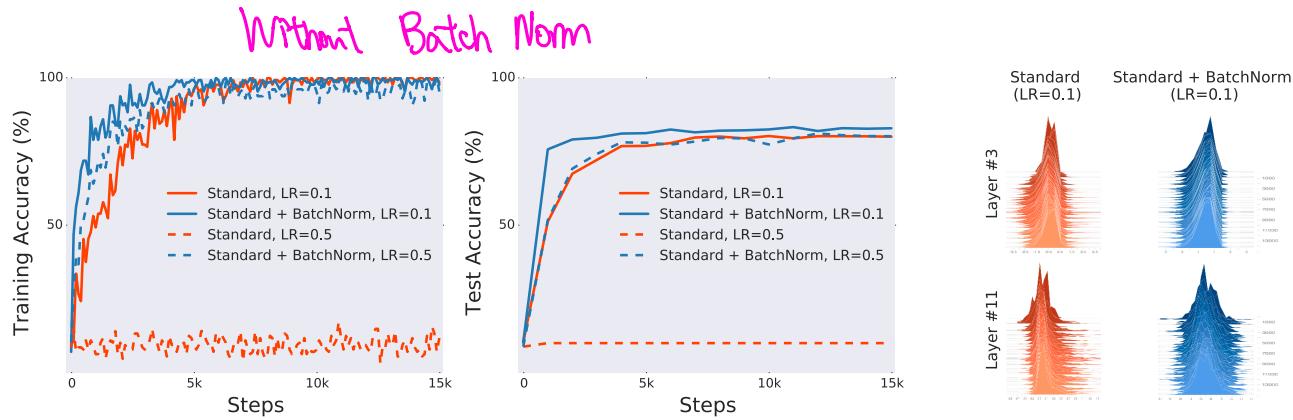


Figure 1: Comparison of (a) training (optimization) and (b) test (generalization) performance of a standard VGG network trained on CIFAR-10 with and without BatchNorm (details in Appendix A). There is a consistent gain in training speed in models with BatchNorm layers. (c) Even though the gap between the performance of the BatchNorm and non-BatchNorm networks is clear, the difference in the evolution of layer input distributions seems to be much less pronounced. (Here, we sampled activations of a given layer and visualized their distribution over training steps.)

Santurkar, S., Tsipras, D., Ilyas, A., & Madry, A. (2018). How does batch normalization help optimization?. In *Advances in Neural Information Processing Systems* (pp. 2488-2498).

# Good Performance of BatchNorm Seems Unrelated to Covariate Shift Prevention

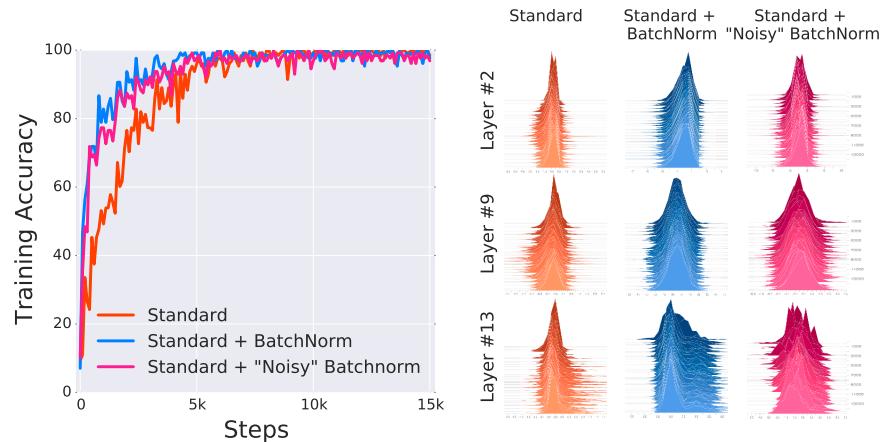


Figure 2: Connections between distributional stability and BatchNorm performance: We compare VGG networks trained without BatchNorm (Standard), with BatchNorm (Standard + BatchNorm) and with explicit “covariate shift” added to BatchNorm layers (Standard + “Noisy” BatchNorm). In the later case, we induce distributional instability by adding *time-varying, non-zero* mean and *non-unit* variance noise independently to each batch normalized activation. The “noisy” BatchNorm model nearly matches the performance of standard BatchNorm model, despite complete distributional instability. We sampled activations of a given layer and visualized their distributions (also cf. Figure 7).

Santurkar, S., Tsipras, D., Ilyas, A., & Madry, A. (2018). How does batch normalization help optimization?. In *Advances in Neural Information Processing Systems* (pp. 2488-2498).

# How Does BatchNorm Work?

## Why Does BatchNorm Help?

2015:

Reduces covariate shift.

Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*.

2018:

Networks with BatchNorm train well with or without ICS.

Hypothesis is that BatchNorm makes the optimization landscape smoother.

Santurkar, S., Tsipras, D., Ilyas, A., & Madry, A. (2018). How does batch normalization help optimization? In *Advances in Neural Information Processing Systems* (pp. 2483-2493).

*depend on single unit between layers*

2018:

"Batch normalization implicitly **discourages single direction reliance**" (here, "single direction reliance" means that an input influences only a single unit or linear combination of single units)

Morcos, A. S., Barrett, D. G., Rabinowitz, N. C., & Botvinick, M. (2018). On the importance of single directions for generalization. *arXiv preprint arXiv:1803.06959*.

# How Does BatchNorm Work?

# Why Does BatchNorm Help?

2018:

*regularization technique*

BatchNorm acts as an implicit regularizer and improves generalization accuracy

Luo, P., Wang, X., Shao, W., & Peng, Z. (2018). Towards understanding regularization in batch normalization. *arXiv preprint arXiv:1809.00846*.

2019:

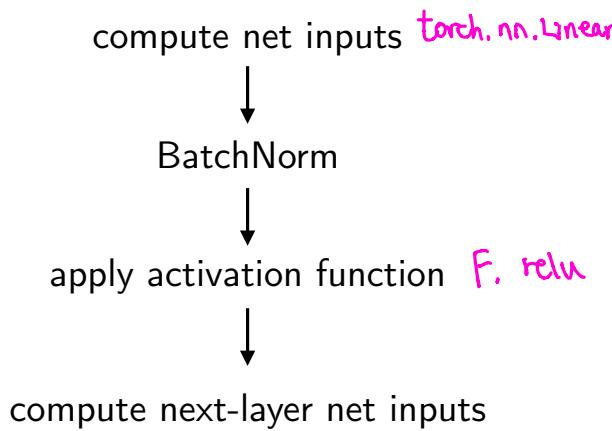
BatchNorm causes exploding gradients, requiring careful tuning when training deep neural nets without skip connections (more about skip connections soon)

Yang, G., Pennington, J., Rao, V., Sohl-Dickstein, J., & Schoenholz, S. S. (2019). A mean field theory of batch normalization. *arXiv preprint arXiv:1902.08129*.

# BatchNorm Variants

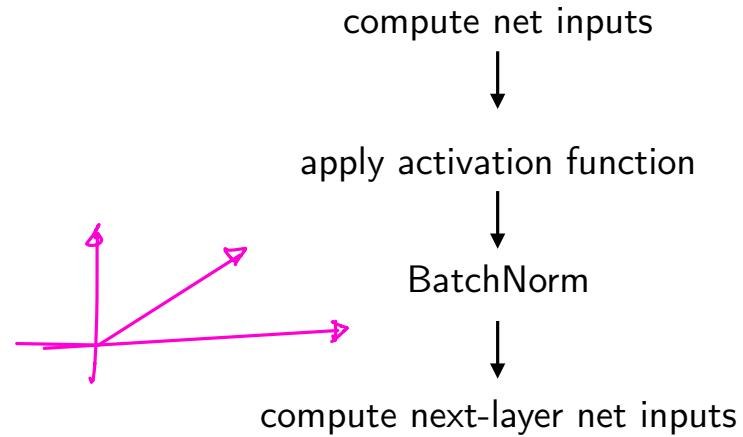
## Pre-Activation

"Original" version  
as discussed in  
previous slides



## Post-Activation

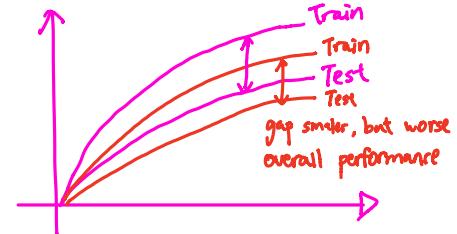
May make more sense,  
but less common



# Some Benchmarks

<https://github.com/ducha-aiki/caffenet-benchmark/blob/master/batchnorm.md#bn---before-or-after-relu>

BN -- before or after ReLU?



Name	Accuracy	LogLoss	Comments
Before	0.474	2.35	As in paper
Before + scale&bias layer	0.478	2.33	As in paper
After	<b>0.499</b>	<b>2.21</b>	
After + scale&bias layer	0.493	2.24	

# Some Benchmarks

<https://github.com/ducha-aiki/caffenet-benchmark/blob/master/batchnorm.md#bn---before-or-after-relu>

## BN and activations

Name	Accuracy	LogLoss	Comments
ReLU	0.499	2.21	
RReLU	0.500	2.20	
PReLU	<b>0.503</b>	<b>2.19</b>	
ELU	0.498	2.23	
Maxout	0.487	2.28	
Sigmoid	0.475	2.35	
TanH	0.448	2.50	
No	0.384	2.96	

# Some Benchmarks

<https://github.com/ducha-aiki/caffenet-benchmark/blob/master/batchnorm.md#bn---before-or-after-relu>

## BN and dropout

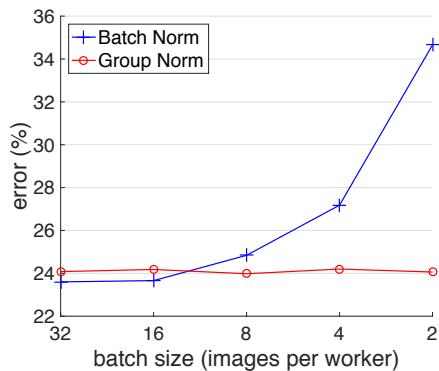
ReLU non-linearity, fc6 and fc7 layer only

Name	Accuracy	LogLoss	Comments
Dropout = 0.5	0.499	2.21	
Dropout = 0.2	<b>0.527</b>	<b>2.09</b>	
Dropout = 0	0.513	2.19	

# Practical Consideration

BatchNorm become more stable with larger minibatch sizes

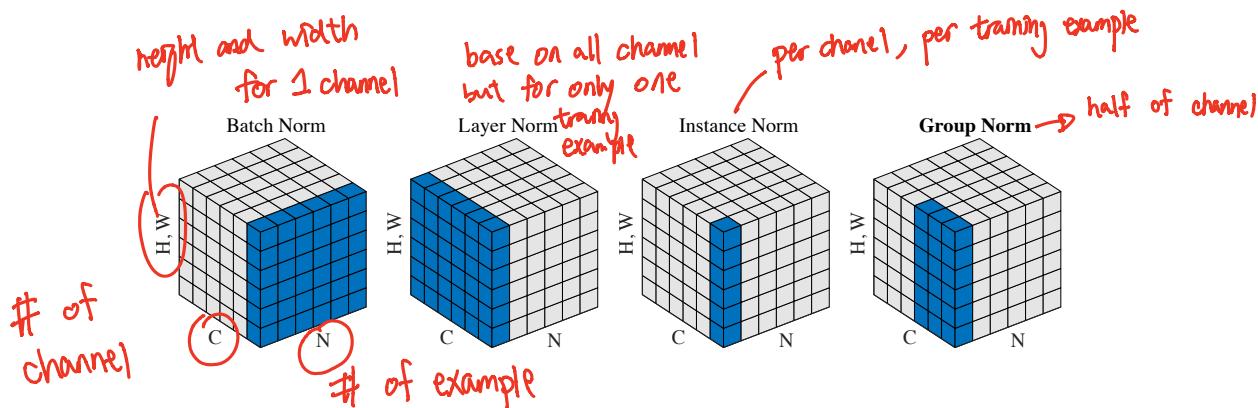
$\sigma, \mu$  will be too noisy if I use small minibatch size, 'compute them from minibatch'



**Figure 1. ImageNet classification error vs. batch sizes.** The model is ResNet-50 trained in the ImageNet training set using 8 workers (GPUs) and evaluated in the validation set. BN's error increases rapidly when reducing the batch size. GN's computation is independent of batch sizes, and its error rate is stable despite the batch size changes. GN has substantially lower error (by 10%) than BN with a batch size of 2.

Wu, Y., & He, K. (2018). Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)* (pp. 3-19).

# Other Normalization Methods for Hidden Activations



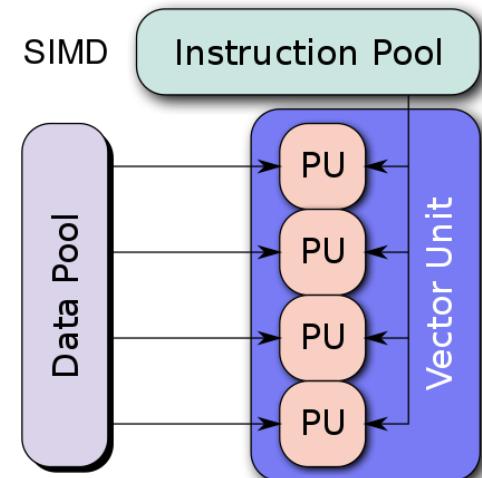
**Figure 2. Normalization methods.** Each subplot shows a feature map tensor. The pixels in blue are normalized by the same mean and variance, computed by aggregating the values of these pixels. Group Norm is illustrated using a group number of 2.

Wu, Y., & He, K. (2018). Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)* (pp. 3-19).

(will revisit after introducing Convolutional Neural Networks)

# Why Minibatch Sizes as Powers of 2?

- Related to SIMD - Single Instruction Multiple Data - paradigm used by CPUs/GPUs
- Comes from mapping the computations (e.g., dot products) to physical processing cores on the GPU, where the number of processing cores is usually a power of 2
- E.g., if we have 48 columns in a matrix, we can map 3 dot products to each processing core if we have 16 processing cores (GPUs usually have many, many more processing cores)



(It might be one of the archaic DL conventions/traditions, and I don't think this matters much anymore for modern frameworks)

Source: <https://upload.wikimedia.org/wikipedia/commons/thumb/c/ce/SIMD2.svg/440px-SIMD2.svg.png>

# Reading Assignments (Optional)

Ioffe, S., & Szegedy, C. (2015, June). Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In *International Conference on Machine Learning* (pp. 448-456).

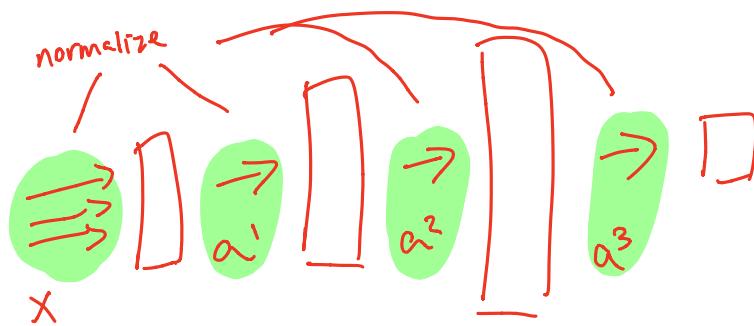
<http://proceedings.mlr.press/v37/ioffe15.html>

update our network more  
efficiently

# Part 2: Weight Initialization

Batch Norm:

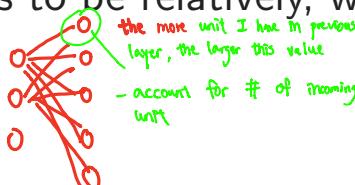
normalize input to NN, and each hidden layer



# Weight Initialization

only happen once

- We previously discussed that we want to initialize weight to small, random numbers to break symmetry
- Also, we want the weights to be relatively, why?

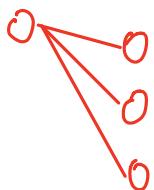


$$Z = W^T X + b$$

initializing this

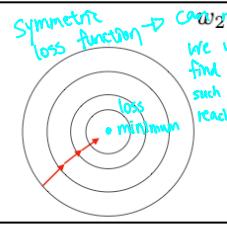
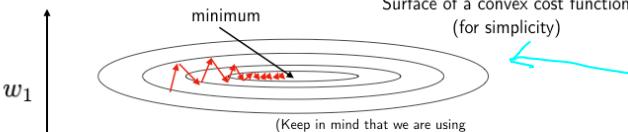
Tip (from an earlier slide):

we would create symmetry in the network if we initialize weight to 0,



net input to hidden layer will all be the same, so if

$$\begin{aligned} w_{1,1} \rightarrow z_1 &= x_1 w_{1,1} + b_1 \\ w_{2,1} \rightarrow z_2 &= x_2 w_{2,1} + b_2 \\ w_{3,1} \rightarrow z_3 &= x_3 w_{3,1} + b_3 \end{aligned}$$



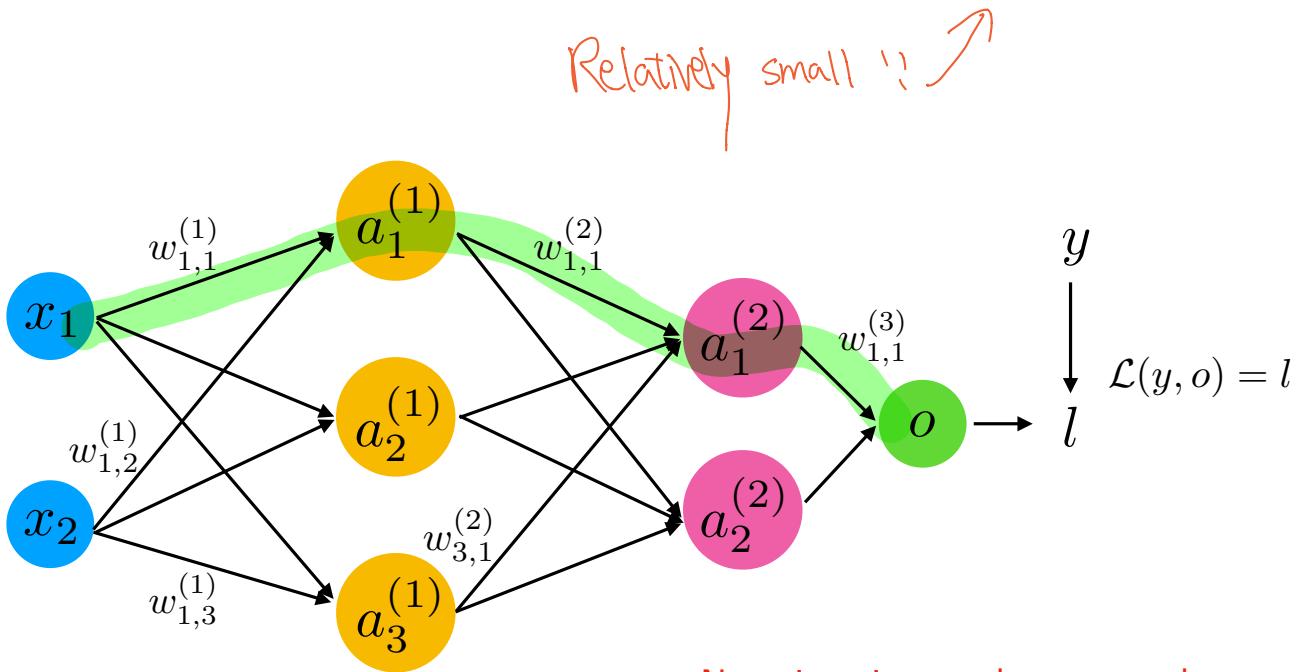
$$x_j^{[i]} = \frac{x_j^{[i]} - \mu_j}{\sigma_j}$$

(scaled feature will have zero mean, unit variance)

- symmetry create problem
- backpropagation would update each unit the same weight

We want weights roughly on the same scale to avoid this problem

# Sidenote: Vanishing/Exploding Gradient Problems



Now, imagine, we have many layers and sigmoid activations ...

$$\begin{aligned}\frac{\partial l}{\partial w_{1,1}^{(1)}} &= \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_1^{(2)}} \cdot \boxed{\frac{\partial a_1^{(2)}}{\partial a_1^{(1)}}} \cdot \frac{\partial a_1^{(1)}}{\partial w_{1,1}^{(1)}} \\ &\quad + \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_2^{(2)}} \cdot \frac{\partial a_2^{(2)}}{\partial a_1^{(1)}} \cdot \frac{\partial a_1^{(1)}}{\partial w_{1,1}^{(1)}}\end{aligned}$$

# Sidenote: Vanishing/Exploding Gradient problems

$$\begin{aligned}\frac{\partial l}{\partial w_{1,1}^{(1)}} &= \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_1^{(2)}} \cdot \boxed{\frac{\partial a_1^{(2)}}{\partial a_1^{(1)}}} \cdot \frac{\partial a_1^{(1)}}{\partial w_{1,1}^{(1)}} \\ &\quad + \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_2^{(2)}} \cdot \frac{\partial a_2^{(2)}}{\partial a_1^{(1)}} \cdot \frac{\partial a_1^{(1)}}{\partial w_{1,1}^{(1)}}\end{aligned}$$

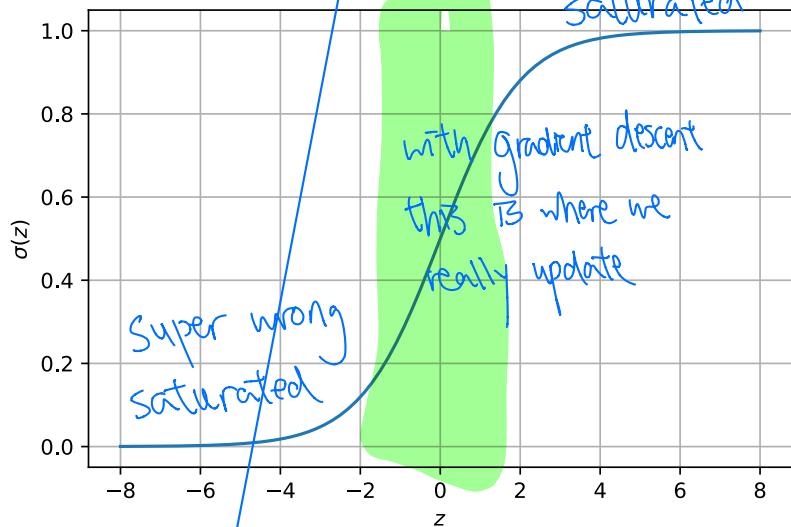
Now, imagine, we have many layers and logistic sigmoid activations ...

$$\sigma'(z^{[i]}) = \sigma(z^{[i]}) \cdot (1 - \sigma(z^{[i]}))$$

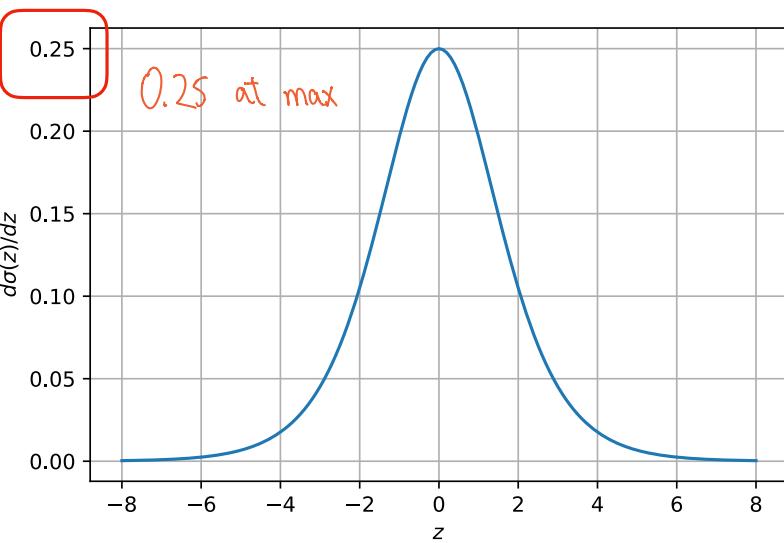
# Sidenote: Vanishing/Exploding Gradient Problems

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

super right  
saturated



$$\frac{d}{dz} \sigma(z) = \frac{e^{-z}}{(1 + e^{-z})^2} = \sigma(z)(1 - \sigma(z))$$



## Sidenote: Vanishing/Exploding Gradient problems



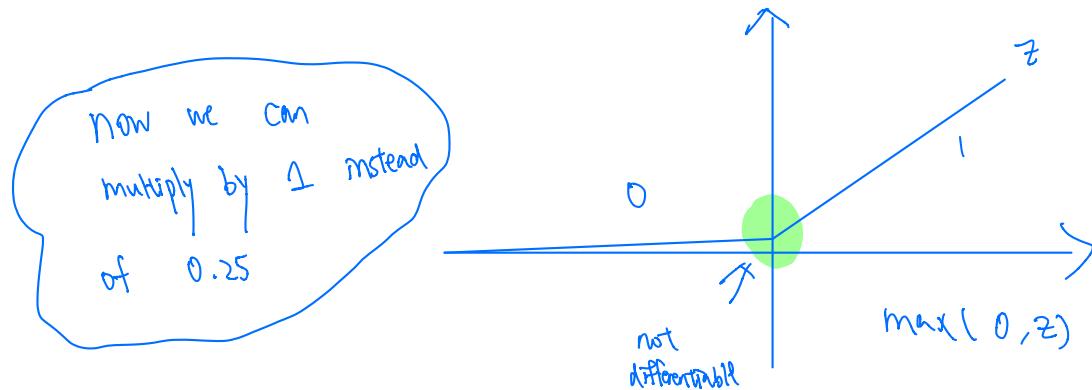
Assume, we have the largest gradient:

$$\frac{d}{dz} \sigma(0.0) = \sigma(0.0)(1 - \sigma(0.0)) = 0.25$$

Even then, for, e.g., 10 layers, we degrade the other gradients substantially!

$$0.25^{10} \approx 10^{-6}$$

# Sidenote: Vanishing/Exploding Gradient problems

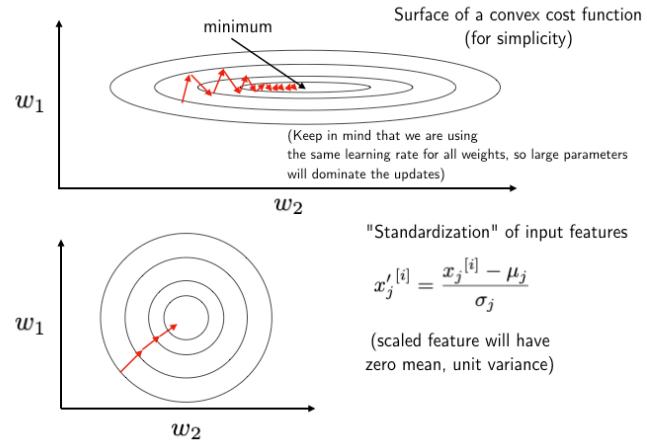


How, do you think, does ReLU behave?

# Weight Initialization

- Traditionally, we can initialize weights by sampling from a random uniform distribution in range [0, 1], or better, [-0.5, 0.5] *more prob. of getting a certain output (more expressive power)*
- Or, we could sample from a Gaussian distribution with mean 0 and small variance (e.g., 0.1 or 0.01)
- When would you choose which?

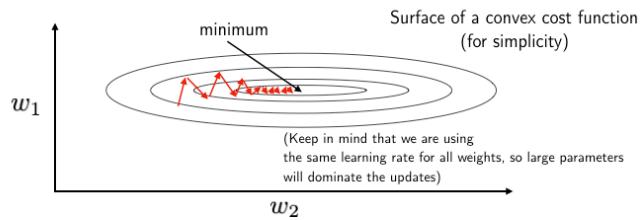
Tip (from an earlier slide):



# Weight Initialization

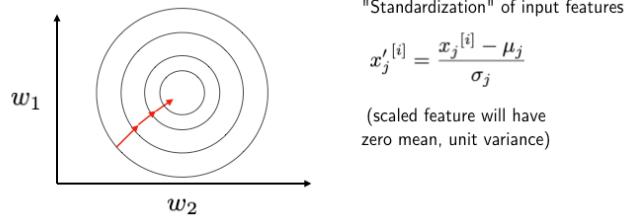
- Traditionally, we can initialize weights by sampling from a random uniform distribution in range [0, 1], or better, [-0.5, 0.5]
- Or, we could sample from a Gaussian distribution with mean 0 and small variance (e.g., 0.1 or 0.01)
- When would you choose which?

Tip (from an earlier slide):



Sidenote: You can initialize the bias units to all zeros

'already break symmetry'



# Custom Weight Initialization in PyTorch

```
class MLP(torch.nn.Module):

    def __init__(self, num_features, num_hidden, num_classes):
        super(MLP, self).__init__()

        self.num_classes = num_classes

        ### 1st hidden layer
        self.linear_1 = torch.nn.Linear(num_features, num_hidden)
        self.linear_1.weight.detach().normal_(0.0, 0.1)          replace
        self.linear_1.bias.detach().zero_()                      mean s.d.

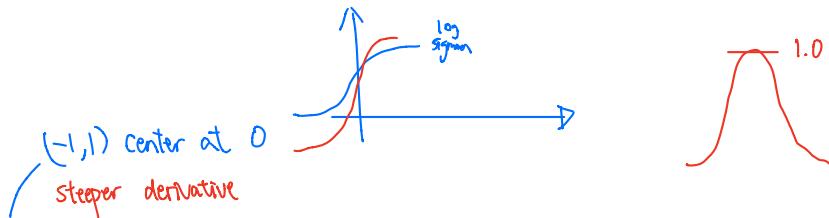
        ### Output layer
        self.linear_out = torch.nn.Linear(num_hidden, num_classes)
        self.linear_out.weight.detach().normal_(0.0, 0.1)
        self.linear_out.bias.detach().zero_() → initialize bias to 0

    def forward(self, x):
        out = self.linear_1(x)
        out = torch.sigmoid(out)
        logits = self.linear_out(out)
        probas = torch.sigmoid(logits)
        return logits, probas
```

*Self. linear. weight = normal*

# Weight Initialization -- Xavier Initialization

Glorot, Xavier, and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks." *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010.



- TanH is a bit more robust regarding vanishing gradients (compared to logistic sigmoid)
- It still has the problem of saturation (near zero gradients if inputs are very large, positive or negative values)
- Xavier initialization is a small improvement for initializing weights for tanH

# Weight Initialization -- Xavier Initialization

Glorot, Xavier, and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks." *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010.

## Method:

Step 1: Initialize weights from Gaussian or uniform distribution with (previous slide)

Step 2: Scale the weights proportional to the number of inputs to the layer

(For the first hidden layer, that is the number of features in the dataset; for the second hidden layer, that is the number of units in the 1st hidden layer etc.)

# Weight Initialization -- Xavier Initialization

Glorot, Xavier, and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks." *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010.

## Method:

Scale the weights proportional to the number of inputs to the layer

In particular, scale as follows:

$$\mathbf{W}^{(l)} := \mathbf{W}^{(l)} \cdot \sqrt{\frac{1}{m^{(l-1)}}}$$

where  $m$  is the number of  
input units to the next  
layer

*layer index*

e.g.,

$$W_{i,j}^{(l)} \sim N(\mu = 0, \sigma^2 = 0.01)$$

(or uniform distr. in a fixed interval, as in the original paper)

# Weight Initialization -- Xavier Initialization

Glorot, Xavier, and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks." *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010.

Sidenote: If you didn't initialize the bias units to all zeros, also include those in the scaling.

$$\mathbf{W}^{(l)} := \mathbf{W}^{(l)} \cdot \sqrt{\frac{1}{m^{(l-1)}}}$$

e.g.,

$$W_{i,j}^{(l)} \sim N(\mu = 0, \sigma^2 = 0.01)$$

(or uniform distr. in a fixed interval, as in the original paper)

# Weight Initialization -- Xavier Initialization

Glorot, Xavier, and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks." *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010.

## Rationale behind this scaling:

Variance of the sample (between data points, not variance of the mean) linearly increases as the sample size increases (variance of the sum of independent variables is the sum of the variances); square root for standard deviation

Want to have same variance between different layers,  $\uparrow$  Variance depends on the # of input

$$\mathbf{W}^{(l)} := \mathbf{W}^{(l)} \cdot \sqrt{\frac{1}{m^{(l-1)}}}$$

e.g.,

$$W_{i,j}^{(l)} \sim N(\mu = 0, \sigma^2 = 0.01)$$

(or uniform distr. in a fixed interval, as in the original paper)

# Weight Initialization -- Xavier Initialization

Glorot, Xavier, and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks." *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010.

Rationale behind this scaling:

Variance of the sample (between data points, not variance of the mean)

linearly increases as the sample size increases (variance of the sum of independent variables is the sum of the variances); square root for standard deviation

Variance scale with the # of input

$$\begin{aligned}\text{Var} \left( a^{(l)} \right) &\approx \text{Var} \left( z^{(l)} \right) = \text{Var} \left( z_j^{(l)} \right) = \text{Var} \left( \sum_{i=1}^{m^{(l-1)}} W_{jk}^{(l)} a_k^{(l-1)} \right) \\ &= \sum_{i=1}^{m^{(l-1)}} \text{Var} \left[ W_{jk}^{(l)} a_k^{(l-1)} \right] \xrightarrow{\text{!} \text{ Independent}} \sum_{i=1}^{m^{(l-1)}} \text{Var} \left[ W_{jk}^{(l)} \right] \text{Var} \left[ a_k^{(l-1)} \right] \\ &= \sum_{i=1}^{m^{(l-1)}} \text{Var} \left[ W^{(l)} \right] \text{Var} \left[ a^{(l-1)} \right] = m^{(l-1)} \text{Var} \left[ W^{(l)} \right] \text{Var} \left[ a^{(l-1)} \right]\end{aligned}$$

make sure the activation function are not saturated and that we have a gradient

# Xavier Initialization in PyTorch

## Semi-Automatic:

```
...
self.linear = torch.nn.Linear(...) # init the layer
torch.nn.init.xavier_uniform_(conv1.weight)
...
...
```

More conveniently for all weights in e.g., fully-connected layers:

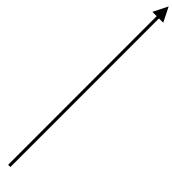
```
...
def weights_init(m):
    if isinstance(m, nn.Linear):
        torch.nn.init.xavier_uniform_(m.weight)
        torch.nn.init.xavier_uniform_(m.bias)
    nn.module
model.apply(weights_init)
...
...
```

if a module is a fully connected layer

nn.module

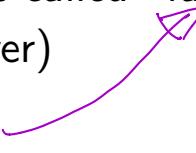
# Weight Initialization -- Xavier Initialization

$$\mathbf{W}^{(l)} := \mathbf{W}^{(l)} \cdot \sqrt{\frac{1}{m^{[l-1]}}}$$



Again, some DL jargon: This is sometimes called "fan in"  
(= number of inputs to a layer)

# of input to ths h.L



# Weight Initialization -- Xavier Initialization

$$\mathbf{W}^{(l)} := \mathbf{W}^{(l)} \cdot \sqrt{\frac{1}{m^{[l-1]}}}$$



From the original paper:

We initialized the biases to be 0 and the weights  $W_{ij}$  at each layer with the following commonly used heuristic:

$$W_{ij} \sim U\left[-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\right], \quad (1)$$

where  $U[-a, a]$  is the uniform distribution in the interval  $(-a, a)$  and  $n$  is the size of the previous layer (the number of columns of  $W$ ).

However, in practice, some people also use "fan in" + "fan out" in the denominator, and it works fine

Glorot, Xavier, and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks." *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010.

# Weight Initialization -- Xavier Initialization

$$\mathbf{W}^{(l)} := \mathbf{W}^{(l)} \cdot \sqrt{\frac{1}{m^{[l-1]}}}$$

However, in practice, some people also use "fan in" + "fan out" in the denominator, and it works fine

From the original paper:

We initialized the biases to be 0 and the weights  $W_{ij}$  at each layer with the following commonly used heuristic:

$$W_{ij} \sim U\left[-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\right], \quad (1)$$

where  $U[-a, a]$  is the uniform distribution in the interval  $(-a, a)$  and  $n$  is the size of the previous layer (the number of columns of  $W$ ).

Also from the original paper:

The normalization factor may therefore be important when initializing deep networks because of the multiplicative effect through layers, and we suggest the following initialization procedure to approximately satisfy our objectives of maintaining activation variances and back-propagated gradients variance as one moves up or down the network. We call it the **normalized initialization**:

$$W \sim U\left[-\frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}\right] \quad (16)$$

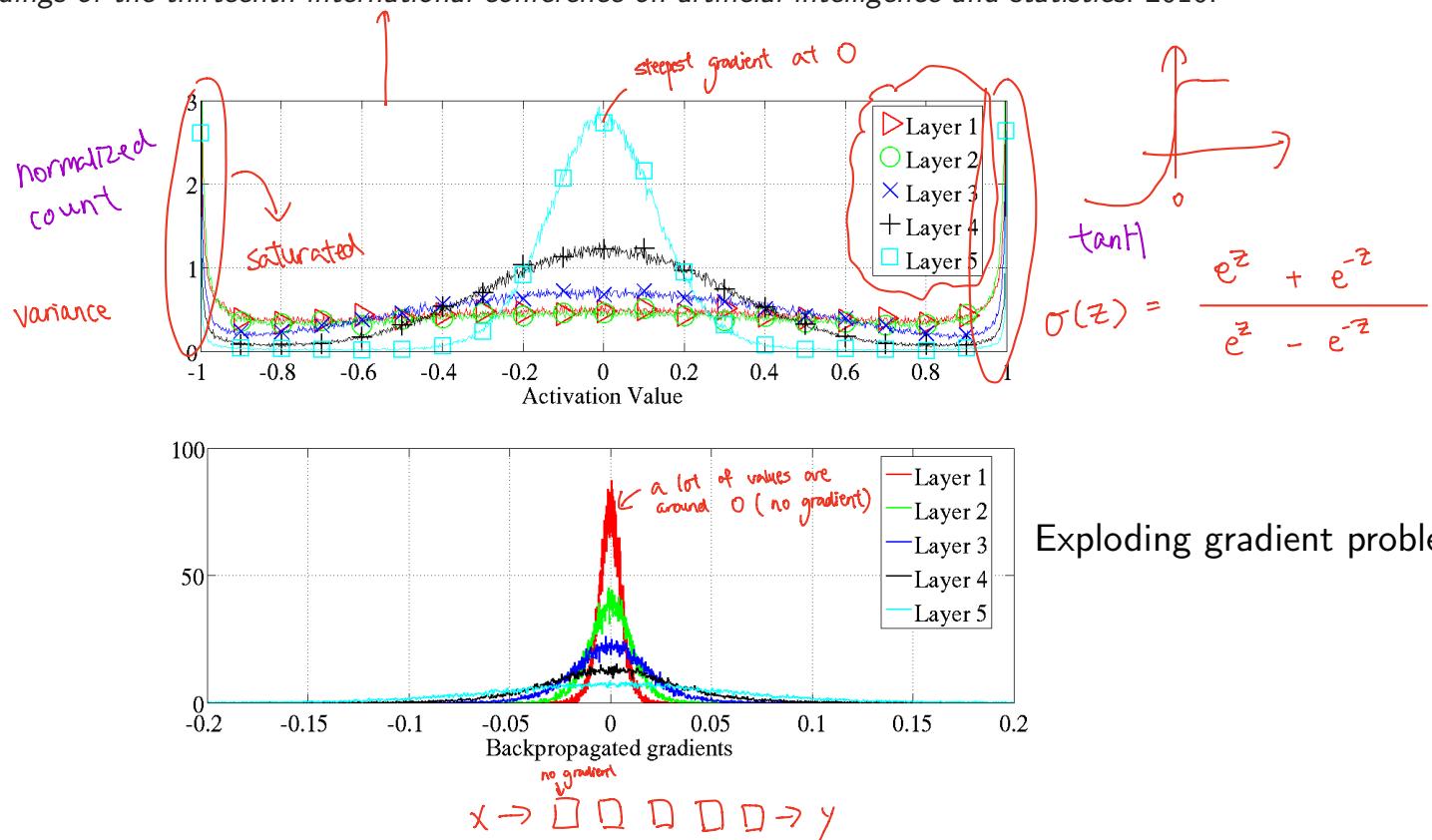
fan in ↑      fan out ↓

Glorot, Xavier, and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks." *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010.

not using standard initialization

# Weight Initialization -- Xavier Initialization

Glorot, Xavier, and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks." *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010.



# Weight Initialization -- Xavier Initialization

Glorot, Xavier, and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks." *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010.

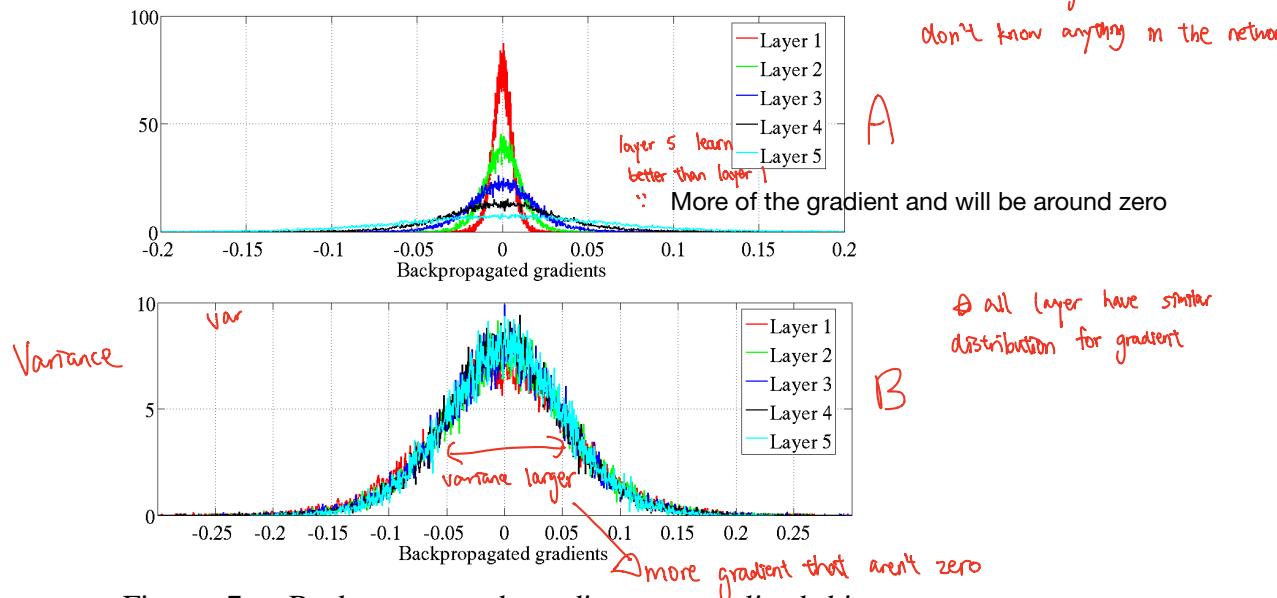


Figure 7: Back-propagated gradients normalized histograms with hyperbolic tangent activation, with standard (top) vs normalized (bottom) initialization. Top: 0-peak decreases for higher layers.

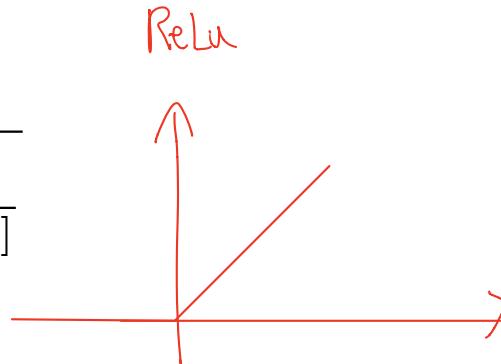
# Weight Initialization -- He Initialization

He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification." In *Proceedings of the IEEE international conference on computer vision*, pp. 1026-1034. 2015.

designed for ReLU

- Assuming activations with mean 0, which is reasonable, Xavier Initialization assumes a derivative of 1 for the activation function (which is reasonable for tanH)
- For ReLU, this is different, as the activations are not centered at zero anymore
- He initialization takes this into account (to see that worked out in math, see the paper)
- The result is that we add a scaling factor of  $2^{0.5}$

$$\mathbf{W}^{(l)} := \mathbf{W}^{(l)} \cdot \sqrt{\frac{2}{m^{[l-1]}}}$$



# Weight Initialization -- He Initialization

He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification." In *Proceedings of the IEEE international conference on computer vision*, pp. 1026-1034. 2015.

- Assuming activations with mean 0, which is reasonable, Xavier Initialization assumes a derivative of 1 for the activation function (which is reasonable for tanH)
- For ReLU, this is different, as the activations are not centered at zero anymore
- He initialization takes this into account (to see that worked out in math, see the paper)
- The result is that we add a scaling factor of  $2^{0.5}$

$$\sigma(x) = \max(0, x)$$

For Leaky ReLU with negative slope alpha:

$$\mathbf{W}^{(l)} := \mathbf{W}^{(l)} \cdot \sqrt{\frac{2}{m^{[l-1]}}}$$

$$\mathbf{W}^{(l)} := \mathbf{W}^{(l)} \cdot \sqrt{\frac{2}{(1 + \alpha^2) \cdot m^{[l-1]}}}$$

# PyTorch Default Weights

PyTorch uses the following scheme by default, which is somewhat similar to Xavier initialization, and works ok in practice most of the time

```
def reset_parameters(self):
    bound = 1 / math.sqrt(self.weight.size(1)) Scaling factor
    init.uniform_(self.weight, -bound, bound)
    if self.bias is not None: [ -1, 1 ]
        init.uniform_(self.bias, -bound, bound)
```

<https://github.com/pytorch/pytorch/blob/master/torch/nn/modules/linear.py#L148>

# PyTorch Default Weights

However, note that different layers have different defaults

He

```
55     def reset_parameters(self):
56         init.kaiming_uniform_(self.weight, a=math.sqrt(5))
57         if self.bias is not None:
58             fan_in, _ = init._calculate_fan_in_and_fan_out(self.weight)
59             bound = 1 / math.sqrt(fan_in)
60             init.uniform_(self.bias, -bound, bound)
61
```

<https://github.com/pytorch/pytorch/blob/master/torch/nn/modules/conv.py#L55>

**Note that if BatchNorm is used,  
initial feature weight choice is less important anyway**

# Reading Assignments (Optional)

Glorot, Xavier, and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks." *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010.

<http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>

He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification." In *Proceedings of the IEEE international conference on computer vision*, pp. 1026-1034. 2015.

<https://arxiv.org/abs/1502.01852>