

ENGR-E 533

Deep Learning Systems

Module 03

Adult Optimization

Minje Kim

Department of Intelligent Systems Engineering

Email: minje@indiana.edu

Website: <http://minjekim.com>

Research Group: <http://saige.sice.indiana.edu>

Meeting Request: <http://doodle.com/minje>



INDIANA UNIVERSITY

**SCHOOL OF INFORMATICS,
COMPUTING, AND ENGINEERING**

Initialization Methods, Activation Functions, and Batch Normalization

Yes, they are related

Initialization

-Why random?

- I remember a fellow student who didn't randomly initialize in his machine learning class a while ago
- What he did was to try out a bunch of different random numbers and pick up the best one
 - Best in the sense of test accuracy
- He was still using some MATLAB functions that generate random numbers
 - Why is not the random initialization?
 - What would be the correct random initialization?
- Random initialization means you generate random numbers to initialize your parameters
 - Therefore, there could be some good choices and bad choices depending on your luck
 - You never know until you see the results
- Eventually there's a technique called ensemble, but for now let's just bear with this uncertainty
- Remaining questions
 - What would be the right p.d.f. to sample from
 - How do we define the parameters of the p.d.f.?



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

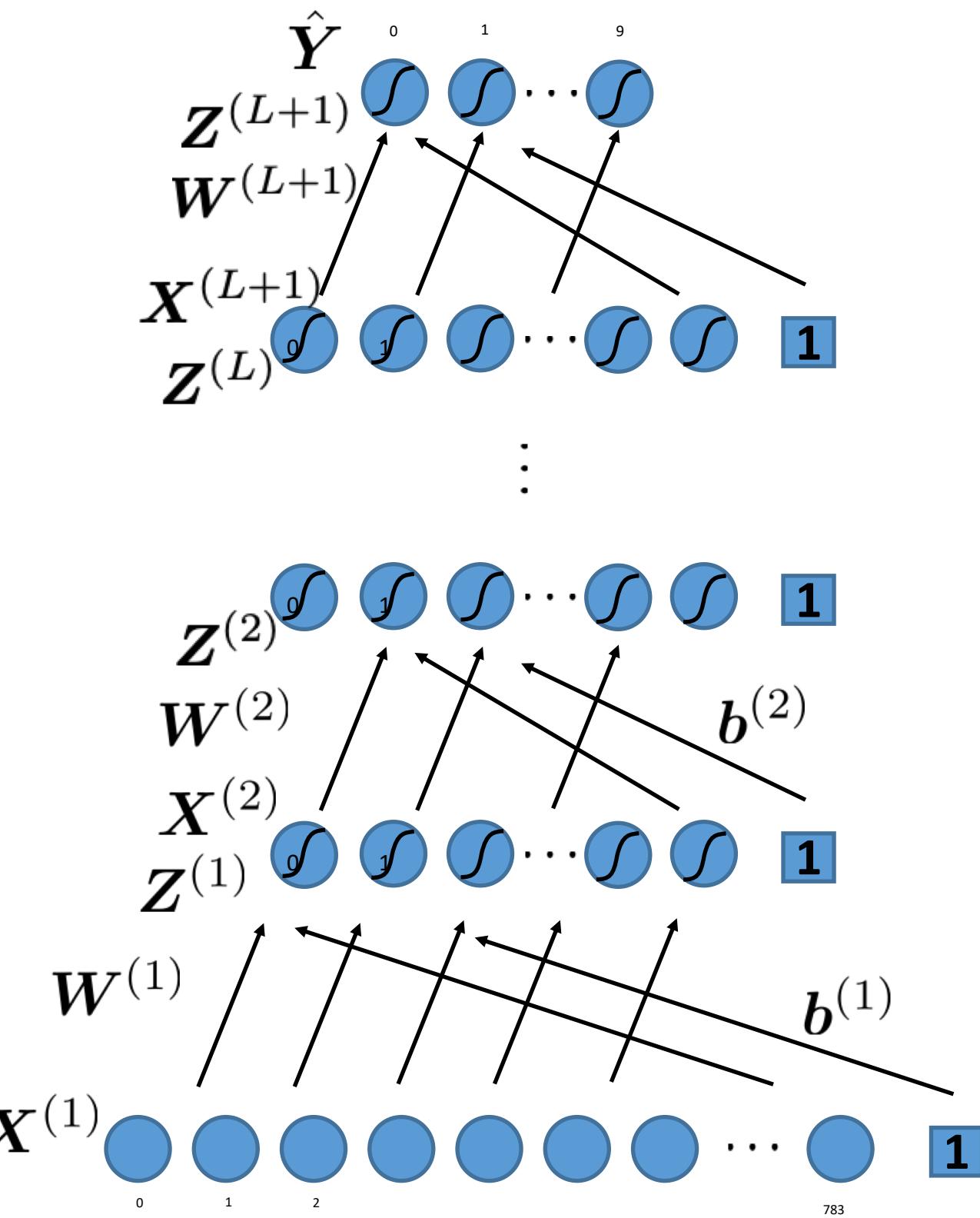
Initialization

-Choosing the best p.d.f.

- Optimal p.d.f. to sample your initial parameters from
 - Suppose the set of weights that give you the global optimum of the objective function
 - e.g. The weights that predict the test samples with the maximum possible accuracy
 - If we find a p.d.f. that's very similar to the sample distribution of these optimal weights
 - I have no theory behind this, but that would be the best one to start from
- Normally we should use some other distributions
 - Maybe some mismatching ones
 - What would be your choice?
 - Normal distribution $\mathcal{N}(\mu, \sigma)$
- Many people also like to use uniform distributions, too
- What would be the mean and variance?
 - It might be safer to use zero mean
 - How about the variance?
 - Let's take a look at them
- Hold on.. You need to know which part of the network to look at:

$$\frac{\partial \mathcal{E}}{\partial \mathbf{W}^{(1)}} = \frac{\partial \mathcal{E}}{\partial \hat{\mathbf{Y}}} \frac{\partial \hat{\mathbf{Y}}}{\partial \mathbf{Z}^{(l+1)}} \prod_{l=1}^L \left[\frac{\partial \mathbf{Z}^{(l+1)}}{\partial \mathbf{X}^{(l+1)}} \frac{\partial \mathbf{X}^{(l+1)}}{\partial \mathbf{Z}^{(l)}} \right] \frac{\partial \mathbf{Z}^{(1)}}{\partial \mathbf{W}^{(1)}} = \left(\left(\mathbf{W}^{(2)\top} \left(\dots \left(\mathbf{W}^{(L+1)\top} (\hat{\mathbf{Y}} - \mathbf{Y}) \right) \odot \sigma'(\checkmark \mathbf{Z}^{(L)}) \dots \right) \right) \odot \sigma'(\checkmark \mathbf{Z}^{(1)}) \right) \mathbf{X}^{(1)\top}$$

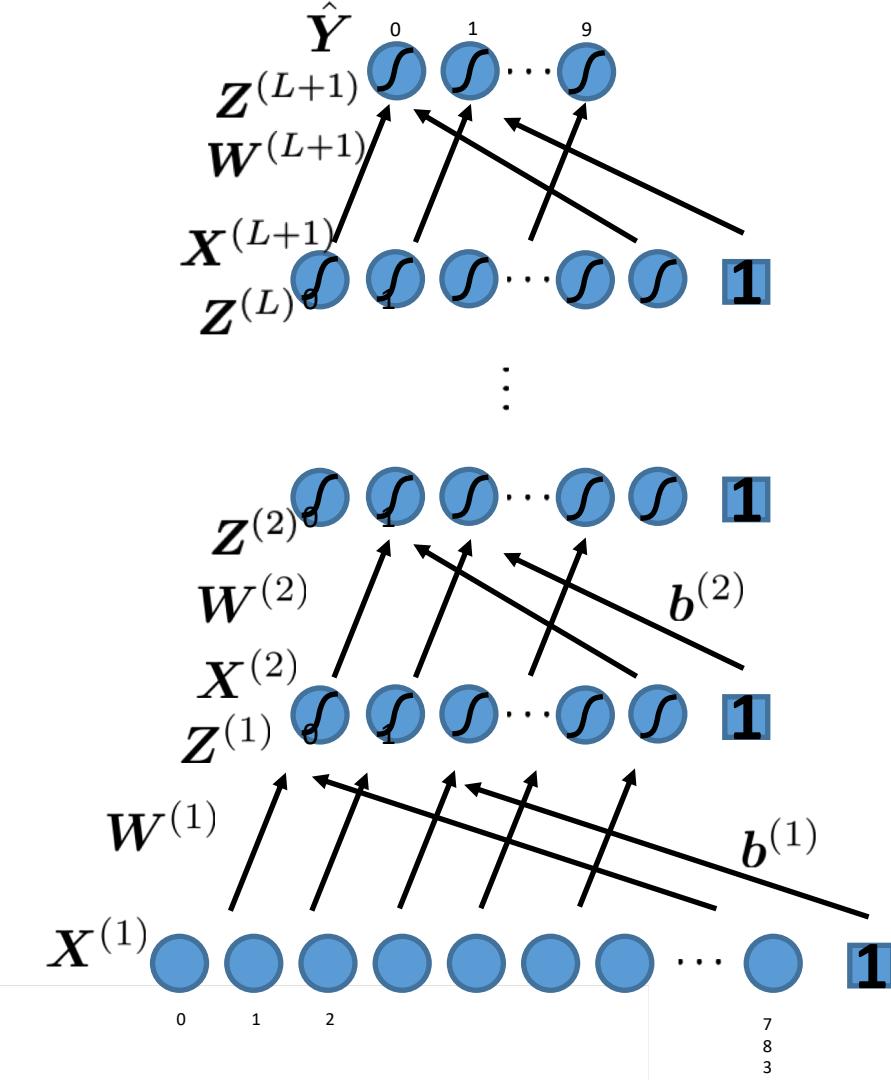
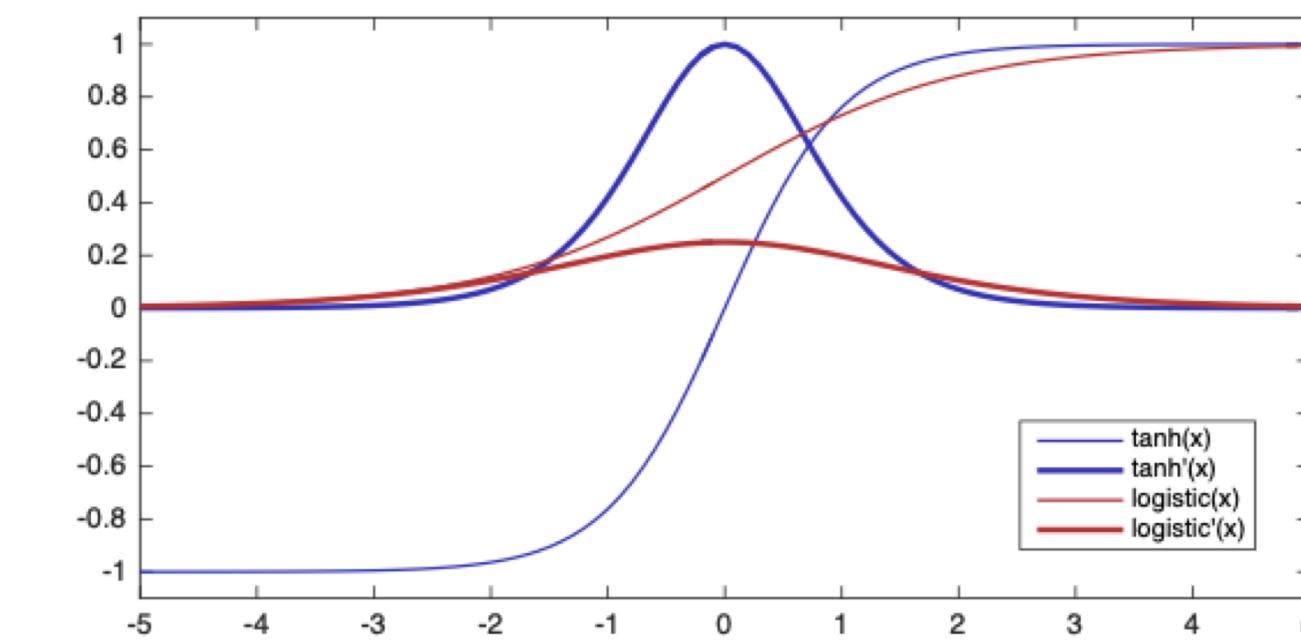
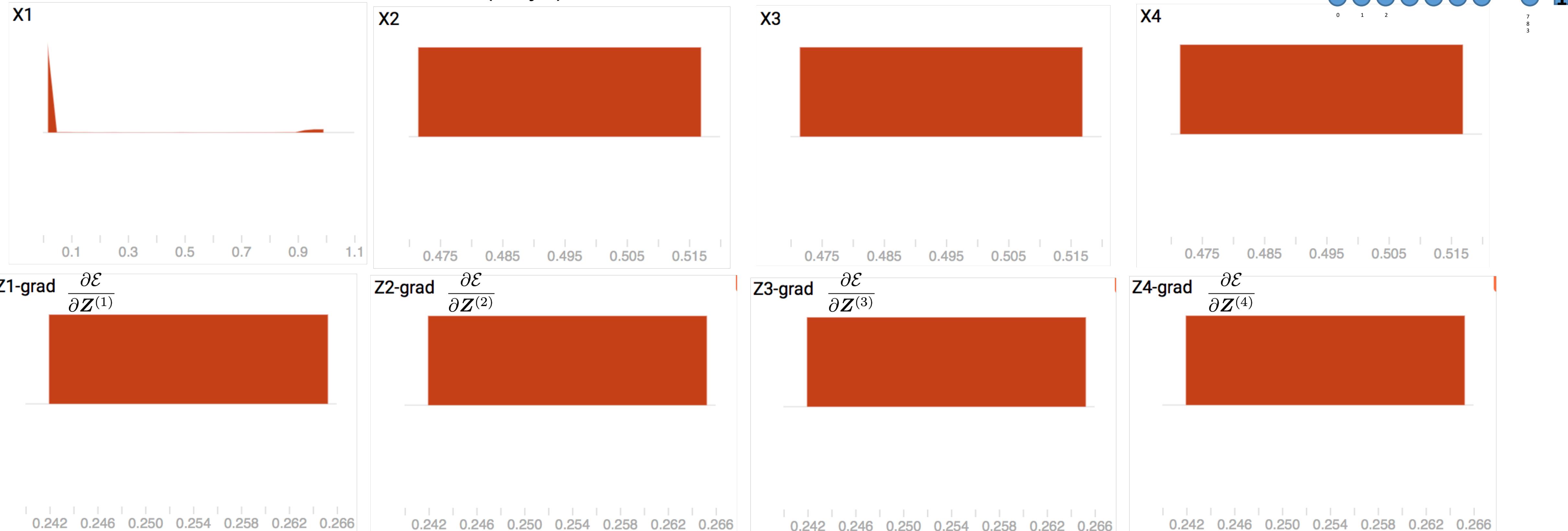
L speed bumps



Initialization

-Choosing the best p.d.f.

- Too small variance $\sigma = 0.0001$
 - All the hidden unit activations (features) are centered around a value
 - Not healthy (features should be distinctive)
 - Even worse with the tanh activations (why?)



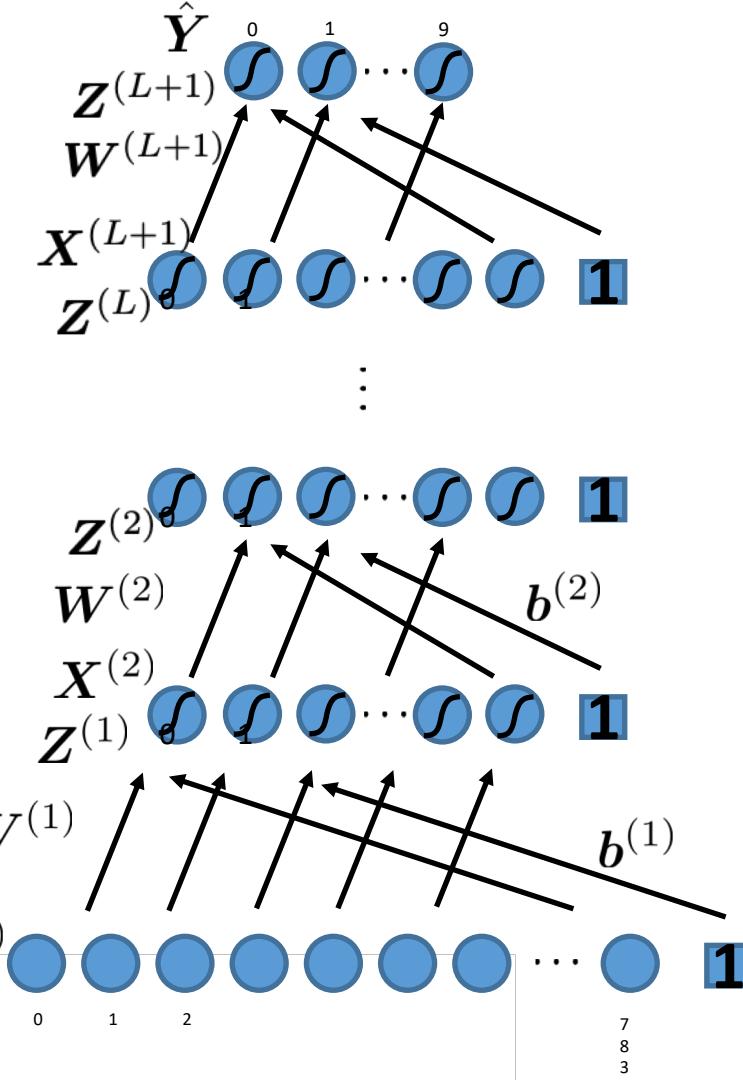
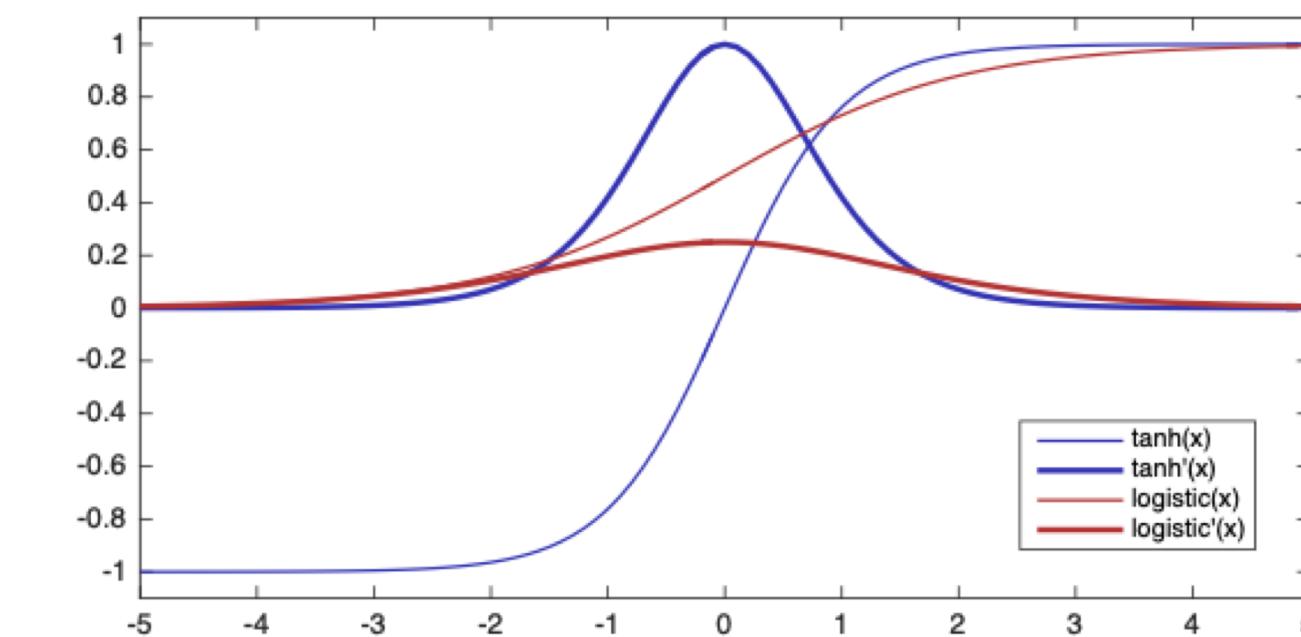
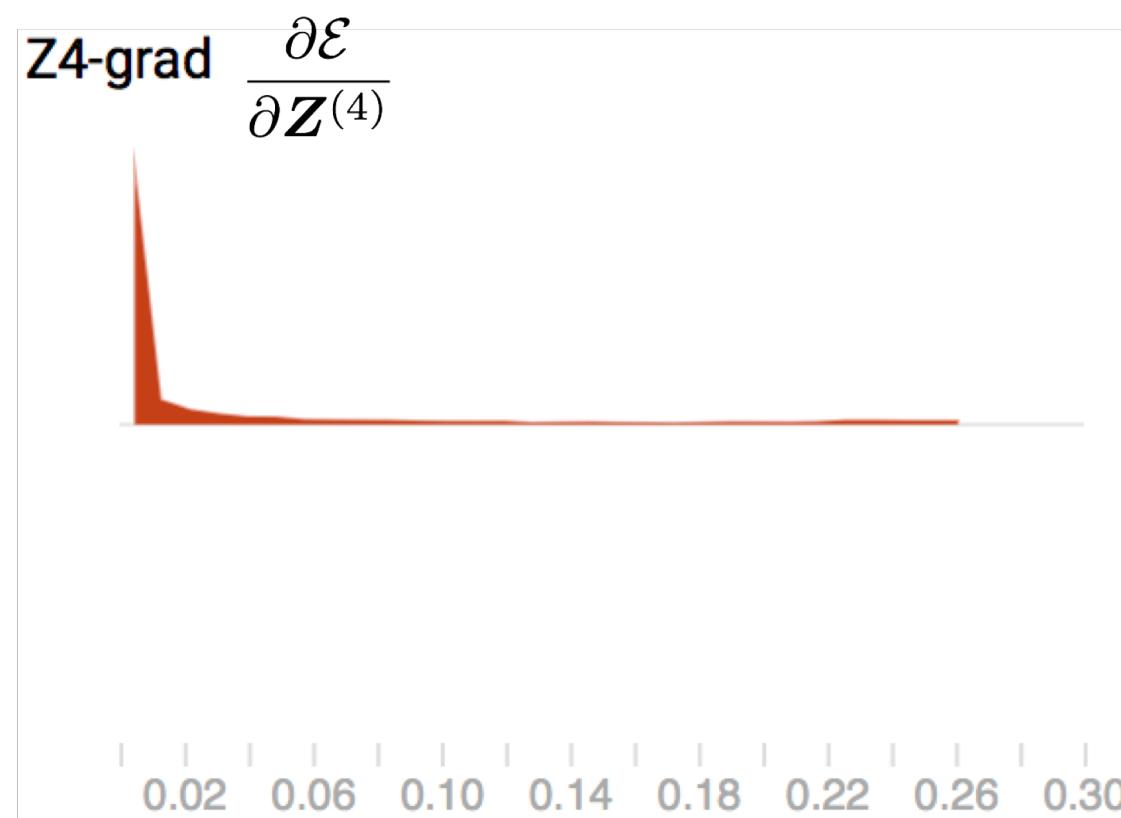
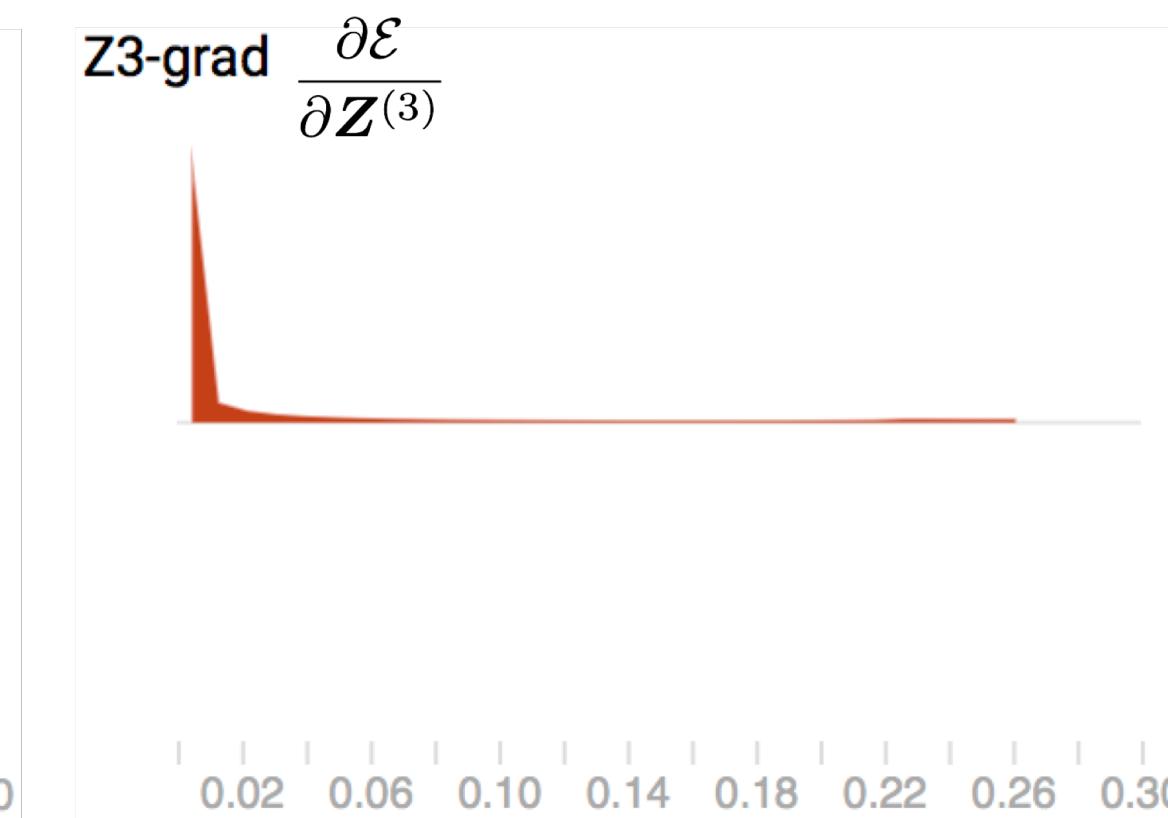
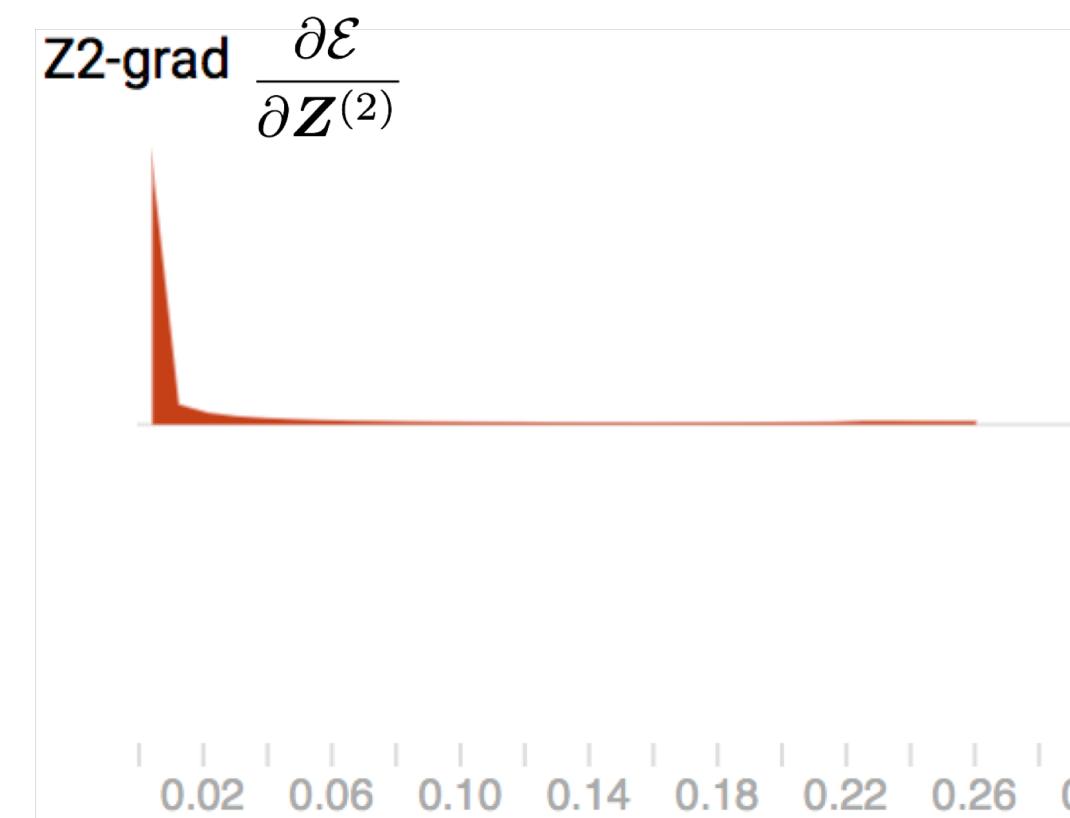
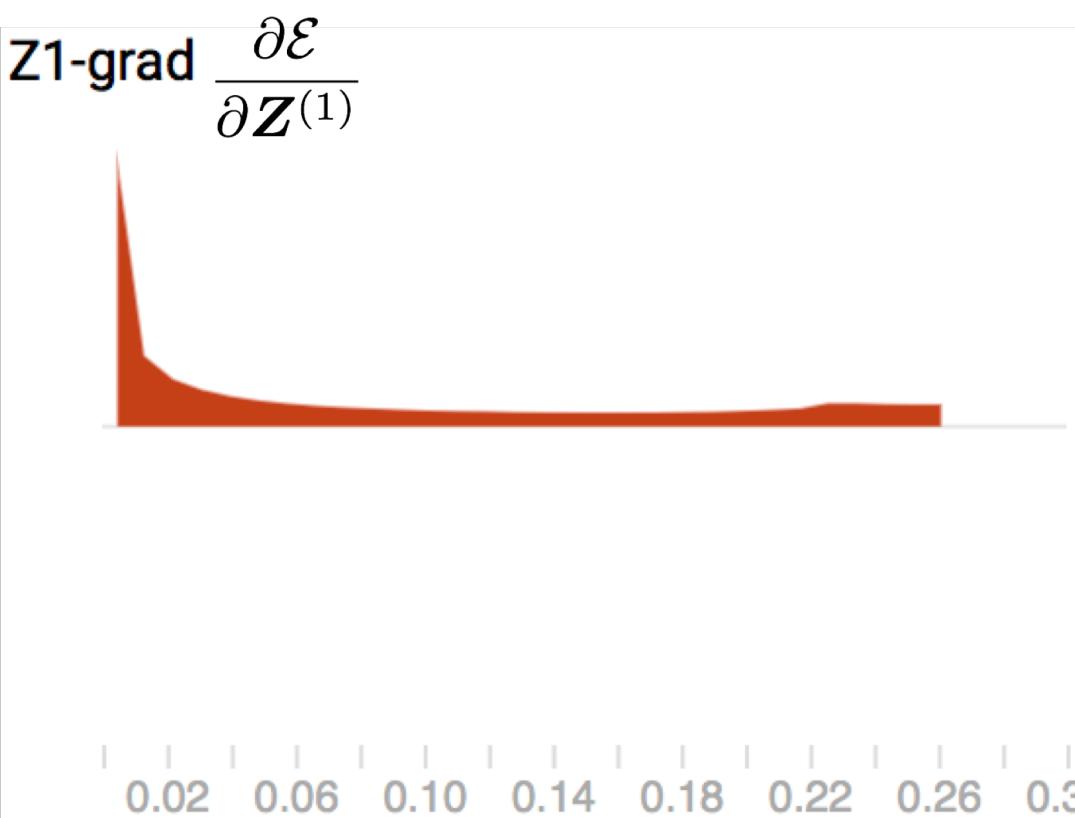
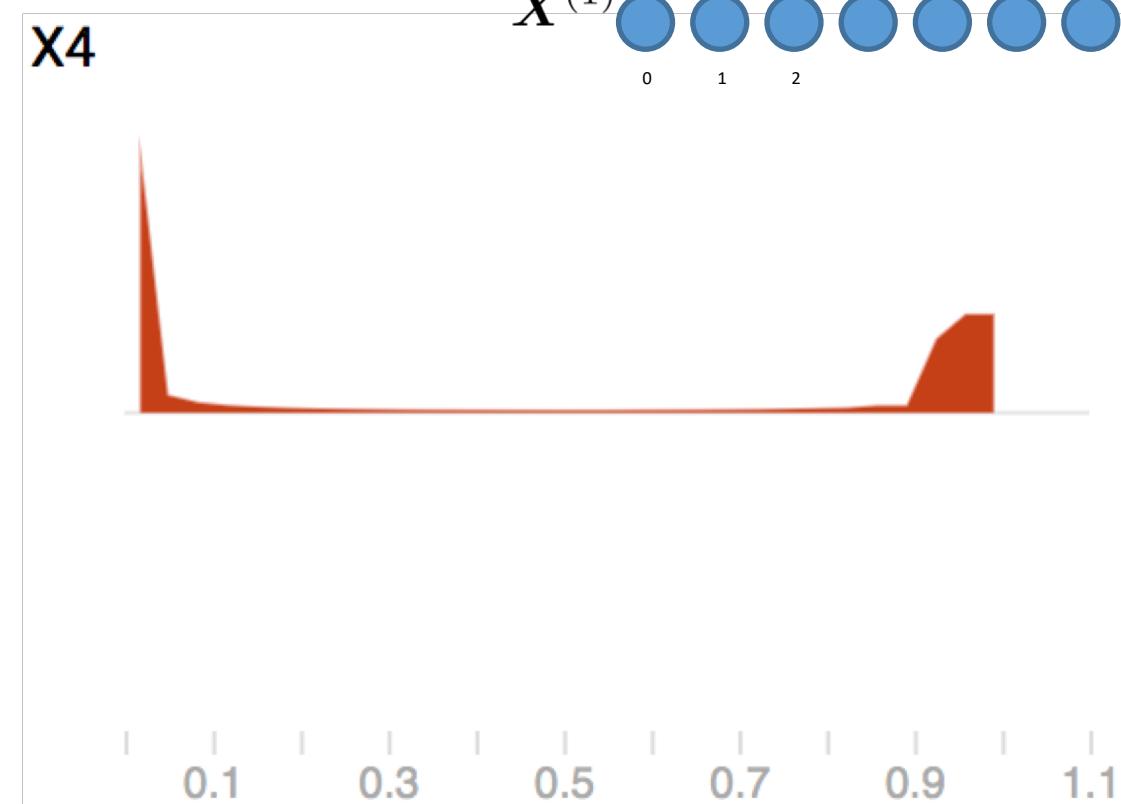
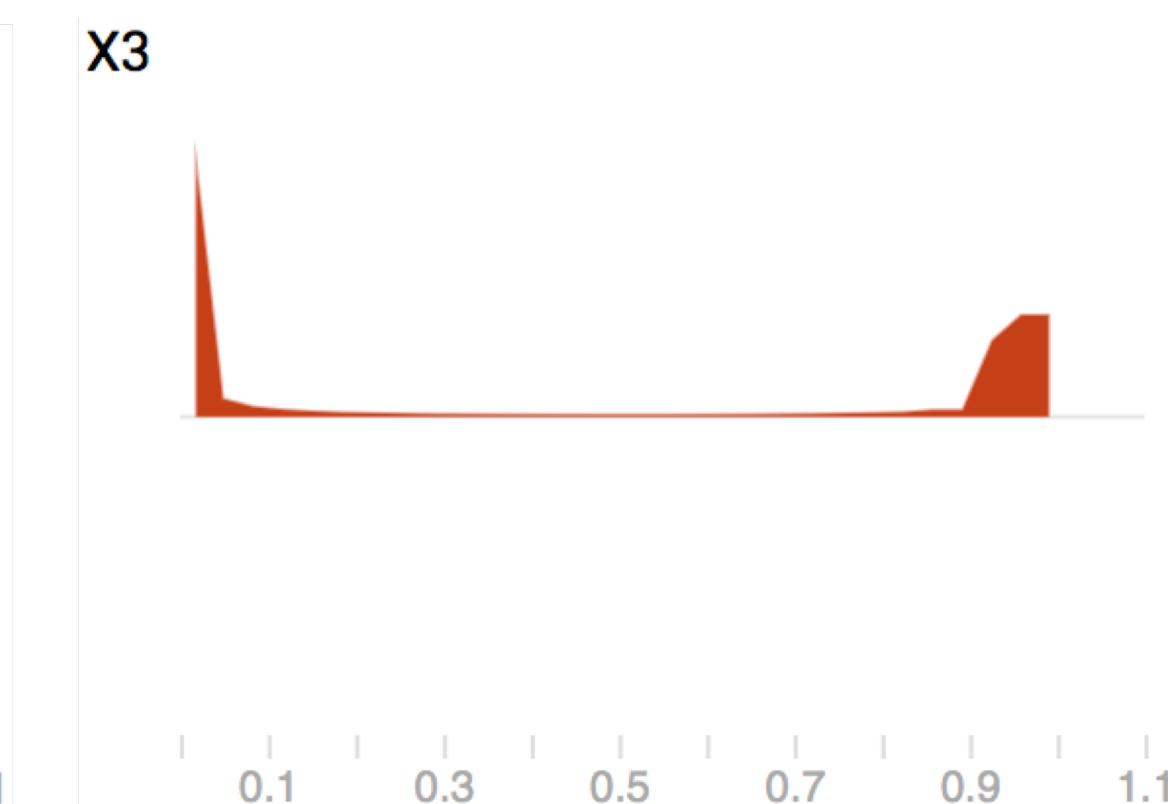
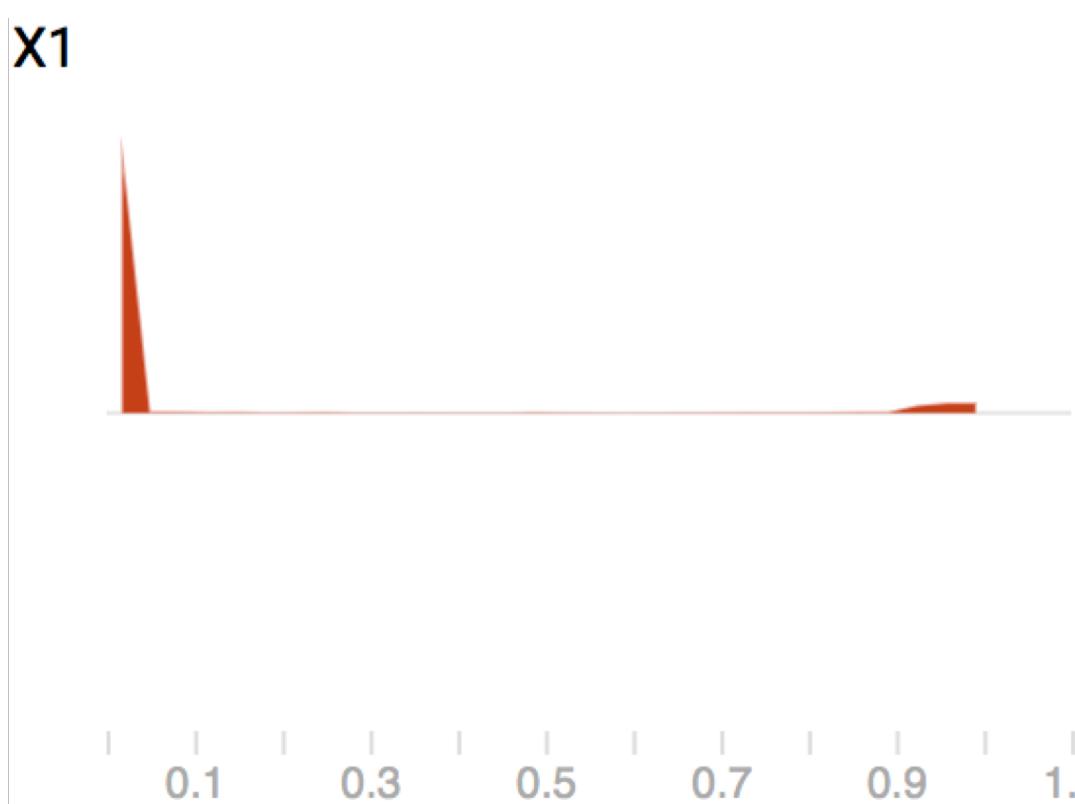
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Initialization

-Choosing the best p.d.f.

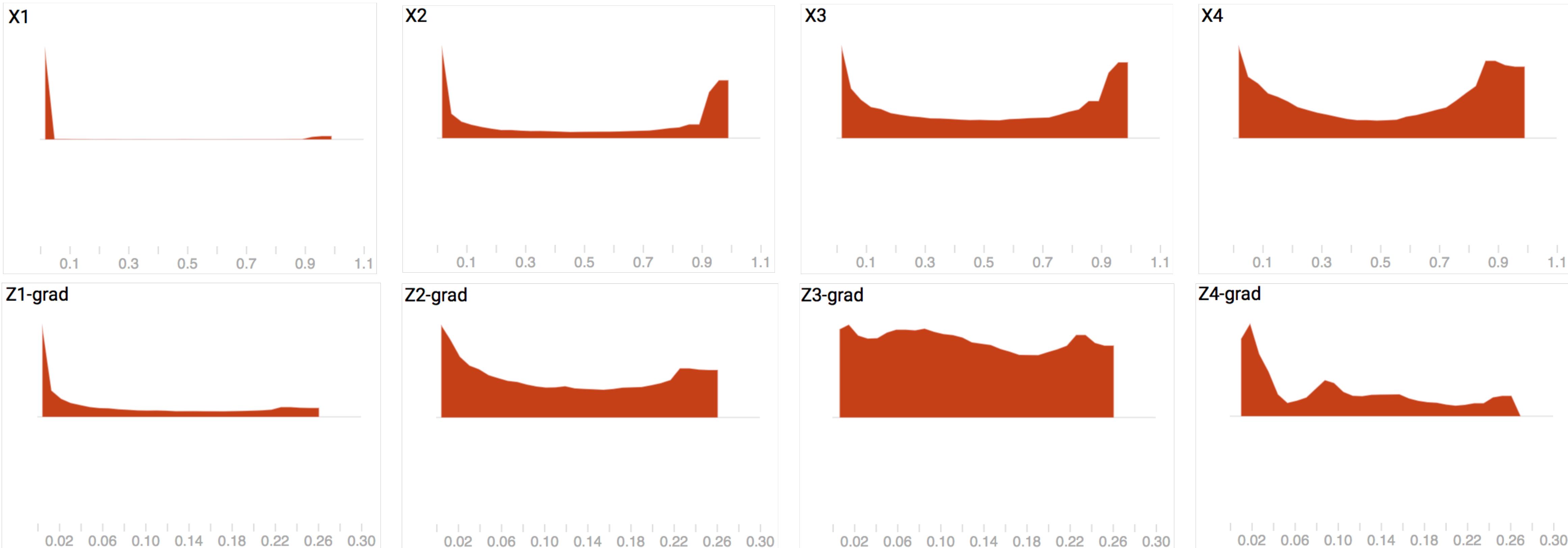
- Too large variance $\sigma = 0.5$
 - Hidden unit outputs are too extreme (why?)
 - Gradient vanishes



Initialization

-Dimension matters

- Still too large variance $\sigma = 0.5$, but now are with only 50 hidden units per layer (previously it was 1,000)
 - The activations are not too extreme
 - Gradient doesn't vanish too much
 - So, what?



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Frishman, Fred. "On the arithmetic means and variances of products and ratios of random variables." A Modern Course on Statistical Distributions in Scientific Work. Springer, Dordrecht, 1975. 401-406.

Initialization

-Xavier Initialization

- Variance of the linear combination $Z_{i,t} = \sum_{j=1}^{N_j} W_{i,j} X_{j,t}$
- To match the variance? $\text{var}(W_{i,j}) = \frac{1}{N_j}$

- We forgot about the other linear combination

$$\frac{\partial \mathcal{E}}{\partial \mathbf{W}^{(1)}} = \frac{\partial \mathcal{E}}{\partial \hat{\mathbf{Y}}} \frac{\partial \hat{\mathbf{Y}}}{\partial \mathbf{Z}^{(l+1)}} \prod_{l=1}^L \left[\frac{\partial \mathbf{Z}^{(l+1)}}{\partial \mathbf{X}^{(l+1)}} \frac{\partial \mathbf{X}^{(l+1)}}{\partial \mathbf{Z}^{(l)}} \right] \frac{\partial \mathbf{Z}^{(1)}}{\partial \mathbf{W}^{(1)}} = \left(\left(\mathbf{W}^{(2)\top} \left(\dots \left(\mathbf{W}^{(L+1)\top} (\hat{\mathbf{Y}} - \mathbf{Y}) \right) \odot \sigma'(\mathbf{Z}^{(L)}) \dots \right) \odot \sigma'(\mathbf{Z}^{(1)}) \right) \mathbf{X}^{(1)\top} \right.$$

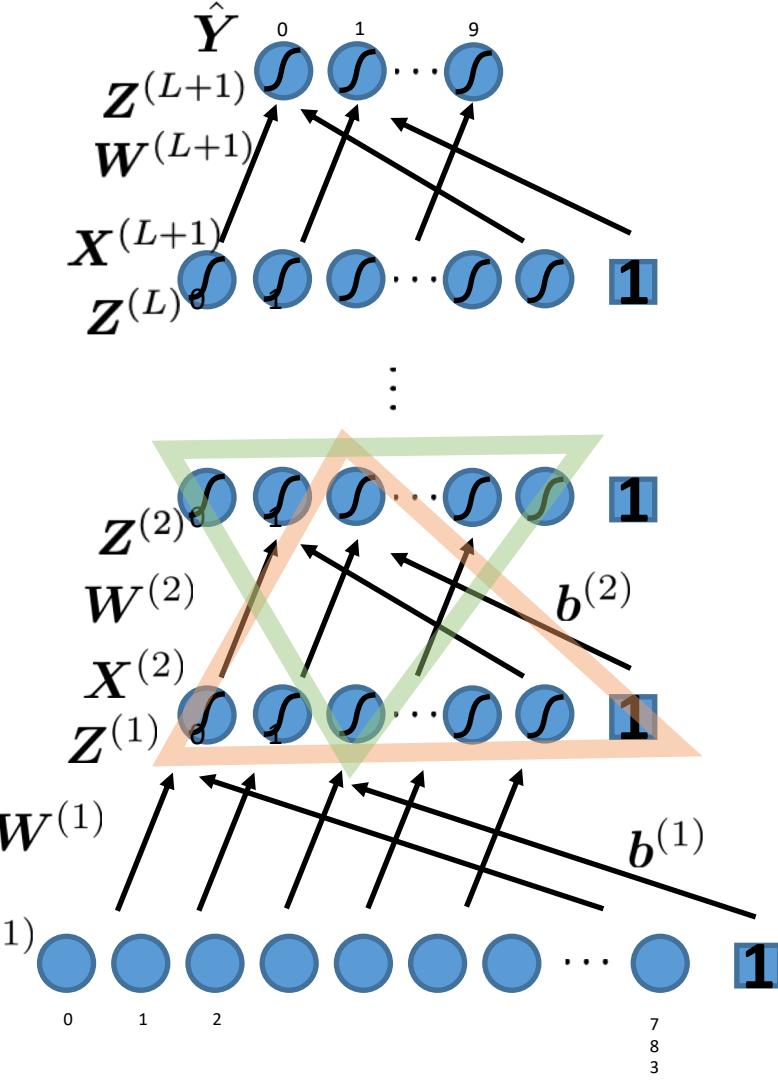
$$\frac{\partial \mathcal{E}}{\partial \mathbf{X}^{(l)}} = \mathbf{W}^{(l)\top} \Delta^{(l)}$$

$$\text{var} \left(\left[\frac{\partial \mathcal{E}}{\partial \mathbf{X}^{(l)}} \right]_{j,t} \right) = \text{var} \left(\sum_{i=1}^{N_i} \mathbf{W}_{i,j}^{(l)} \Delta_{i,t}^{(l)} \right)$$

- To match the variance? $\text{var}(\mathbf{W}_{i,j}) = \frac{1}{N_i}$
- Which one to choose?
 - Average number of units: $(N_j + N_i)/2$
 - Hence, $\text{var}(\mathbf{W}_{i,j}) = \frac{2}{N_i + N_j}$

$$\text{var}(XY) = \text{var}(X)\text{var}(Y)$$

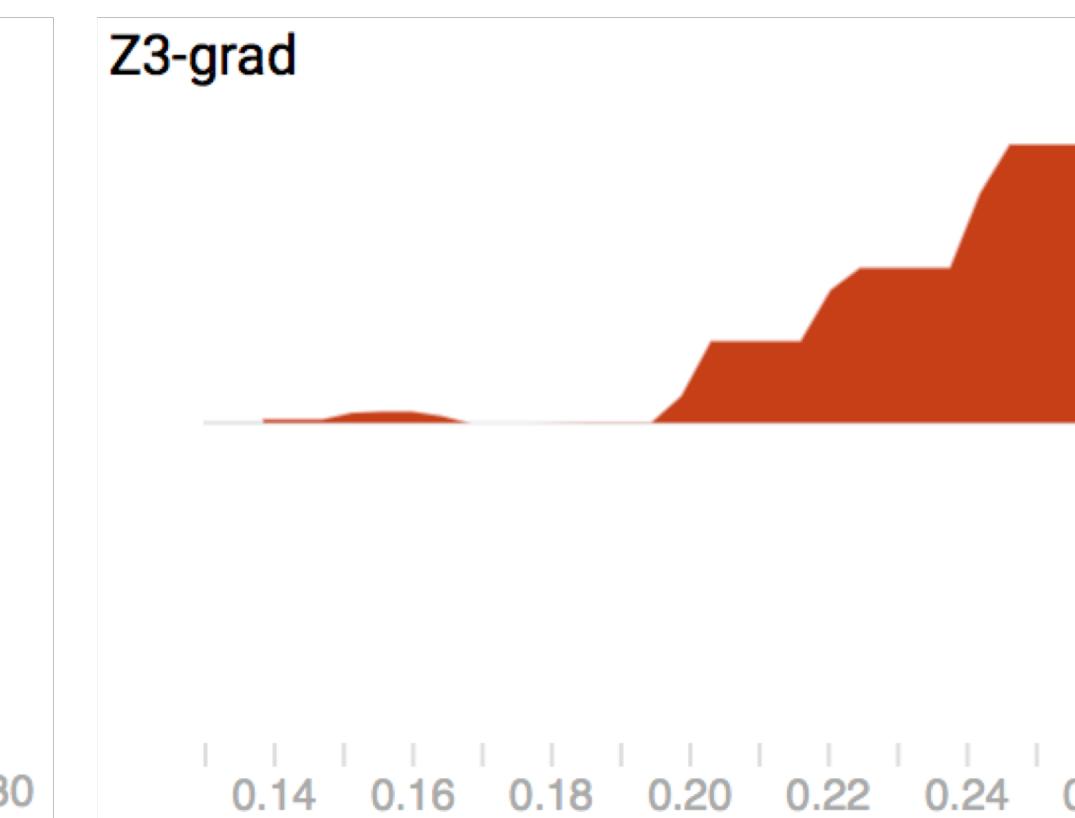
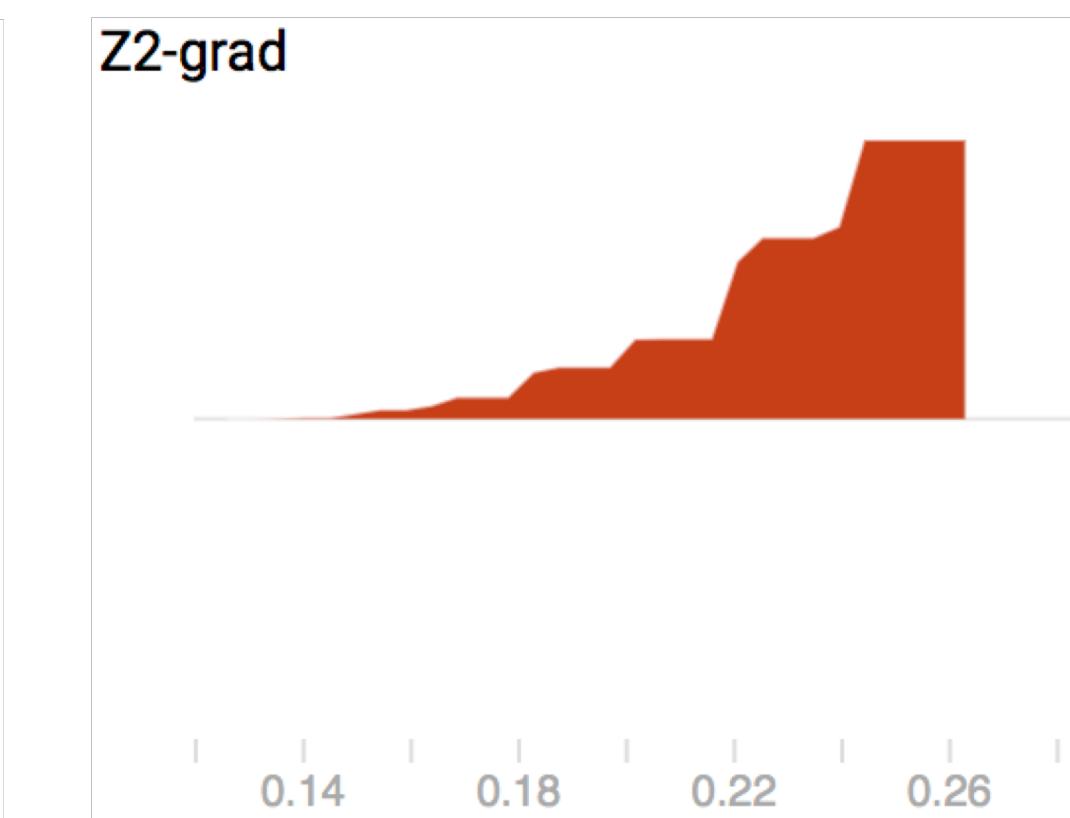
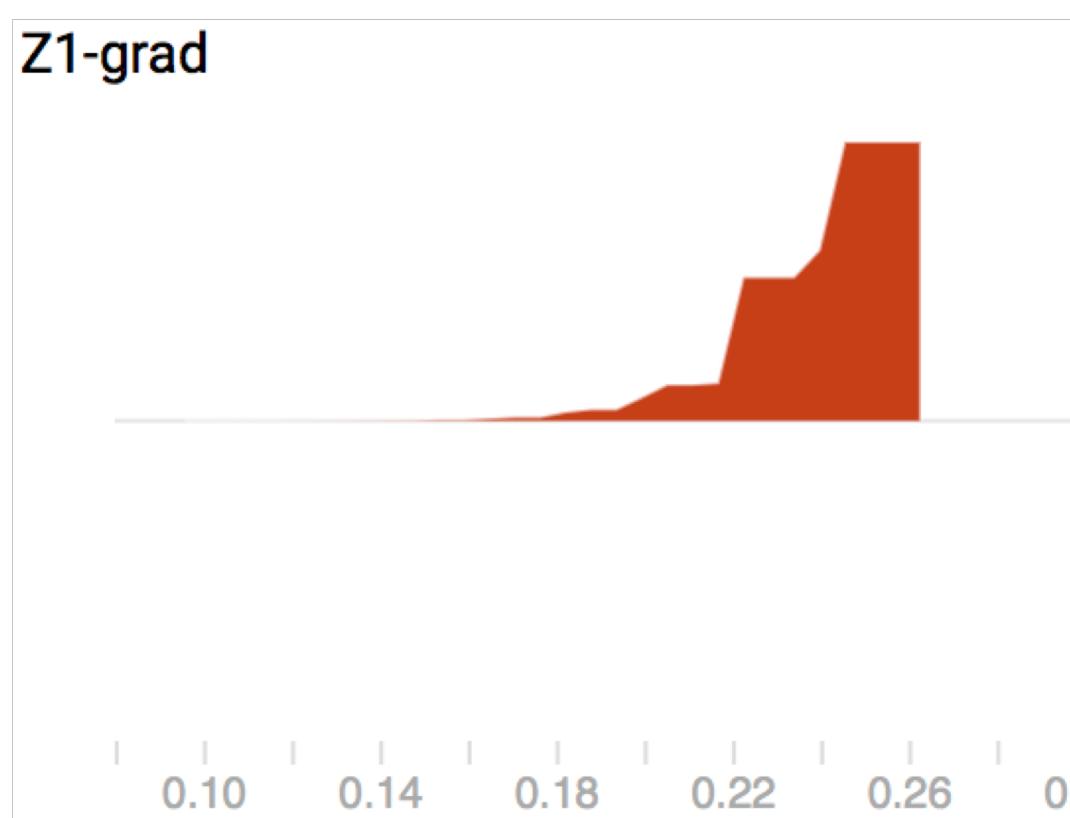
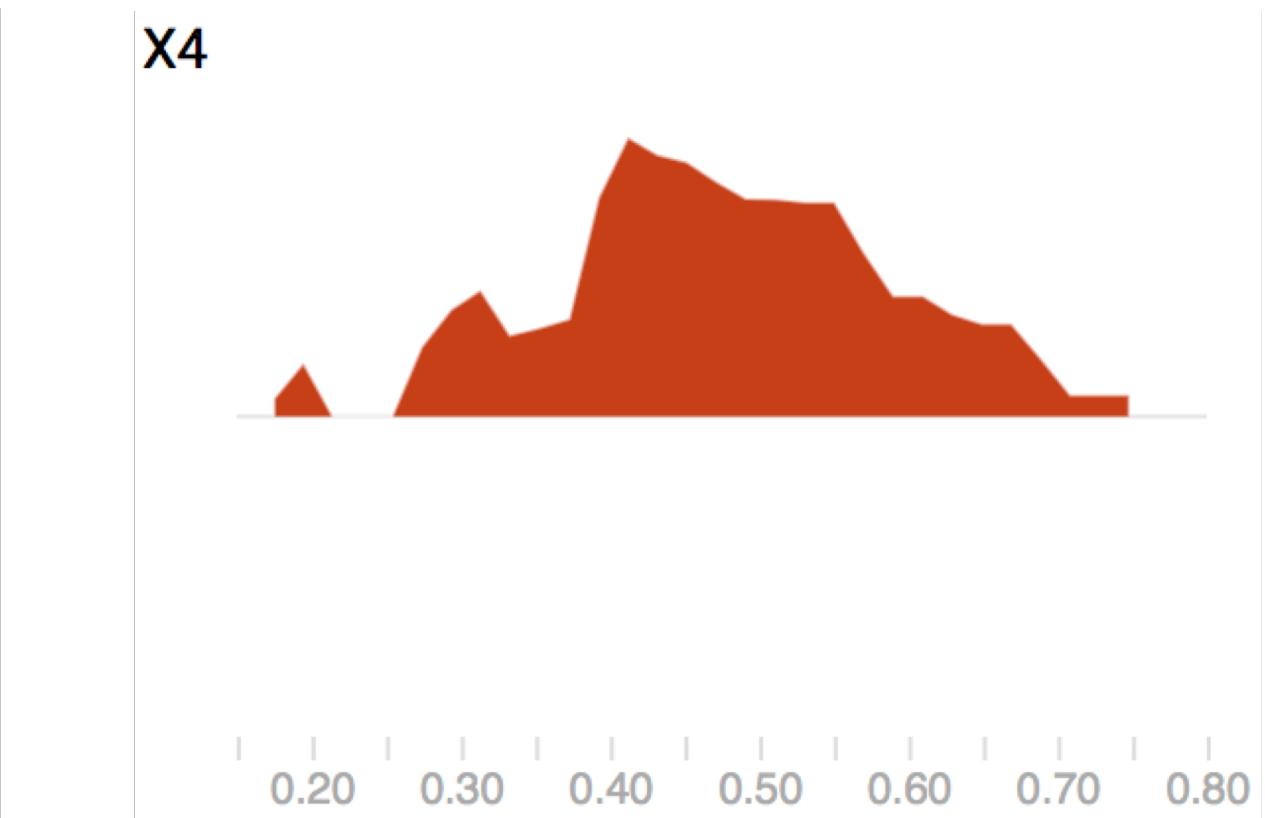
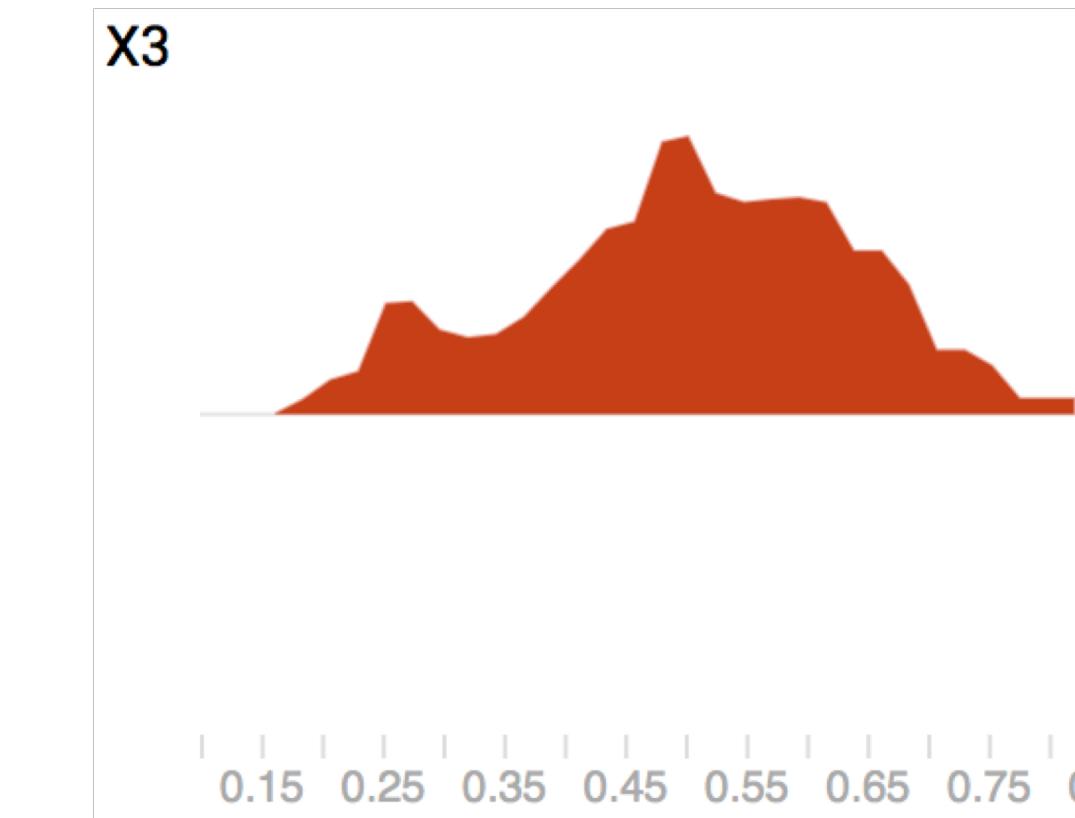
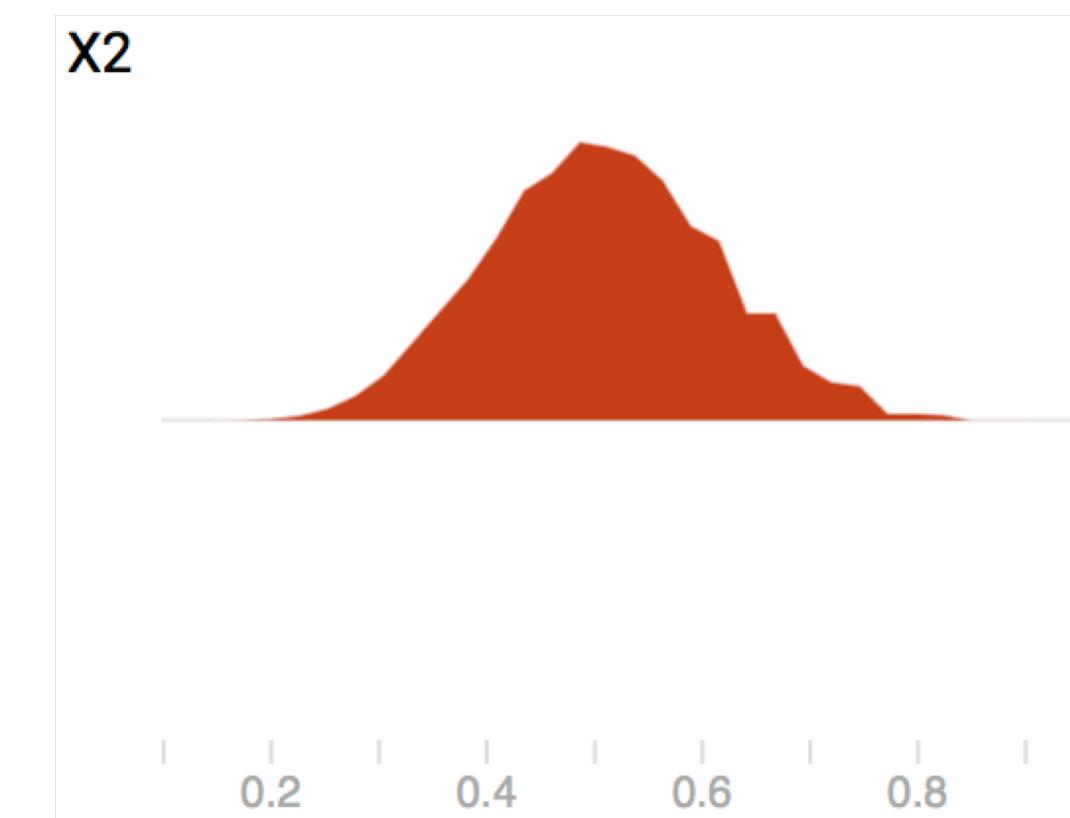
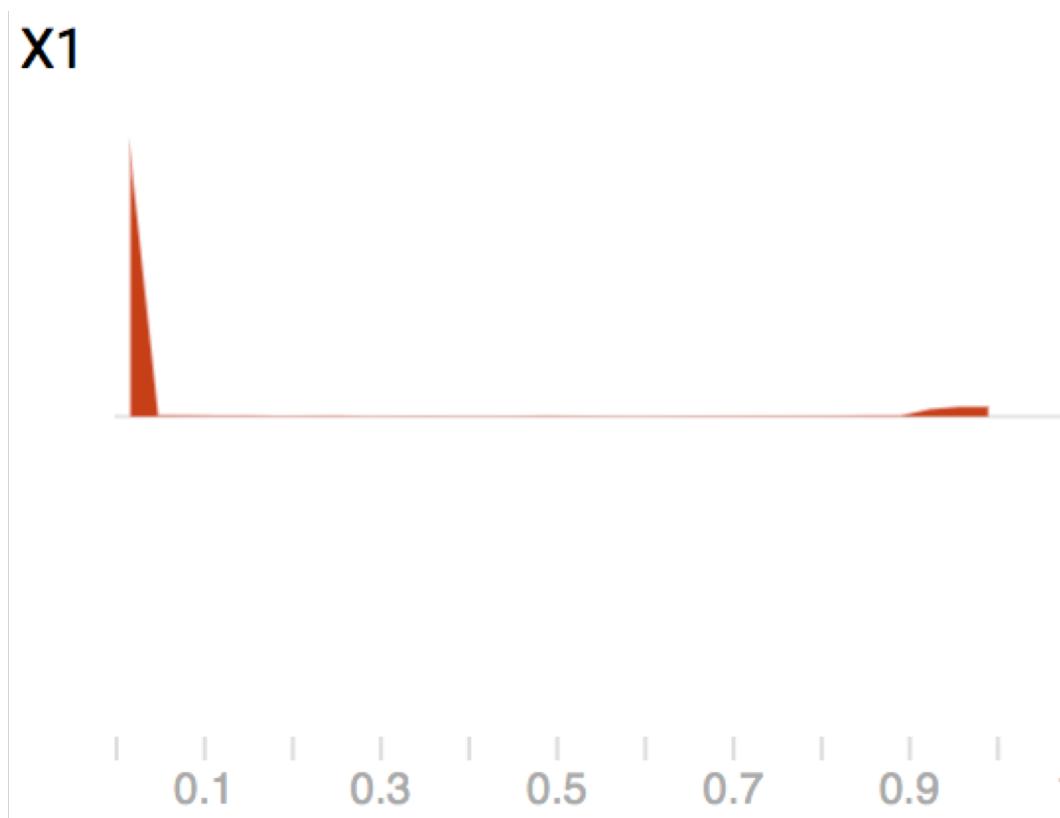
(only if they are independent and zero mean)



Initialization

-Xavier Initialization

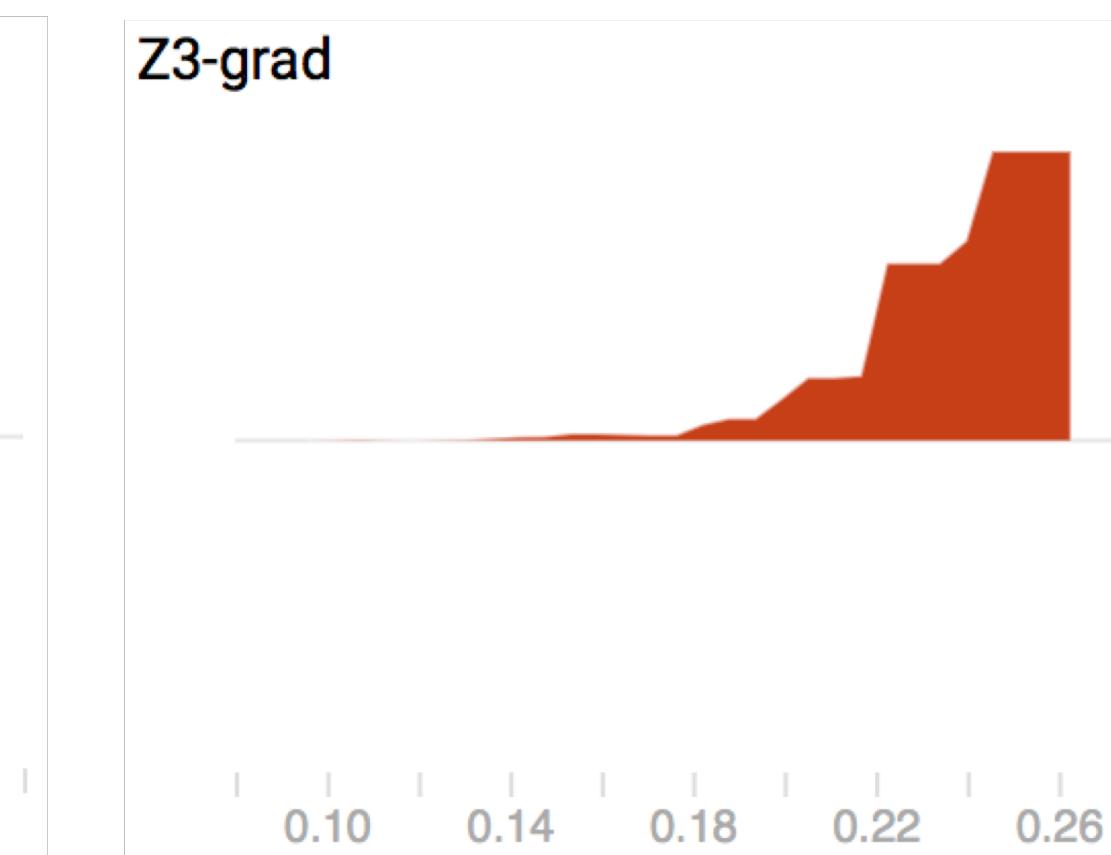
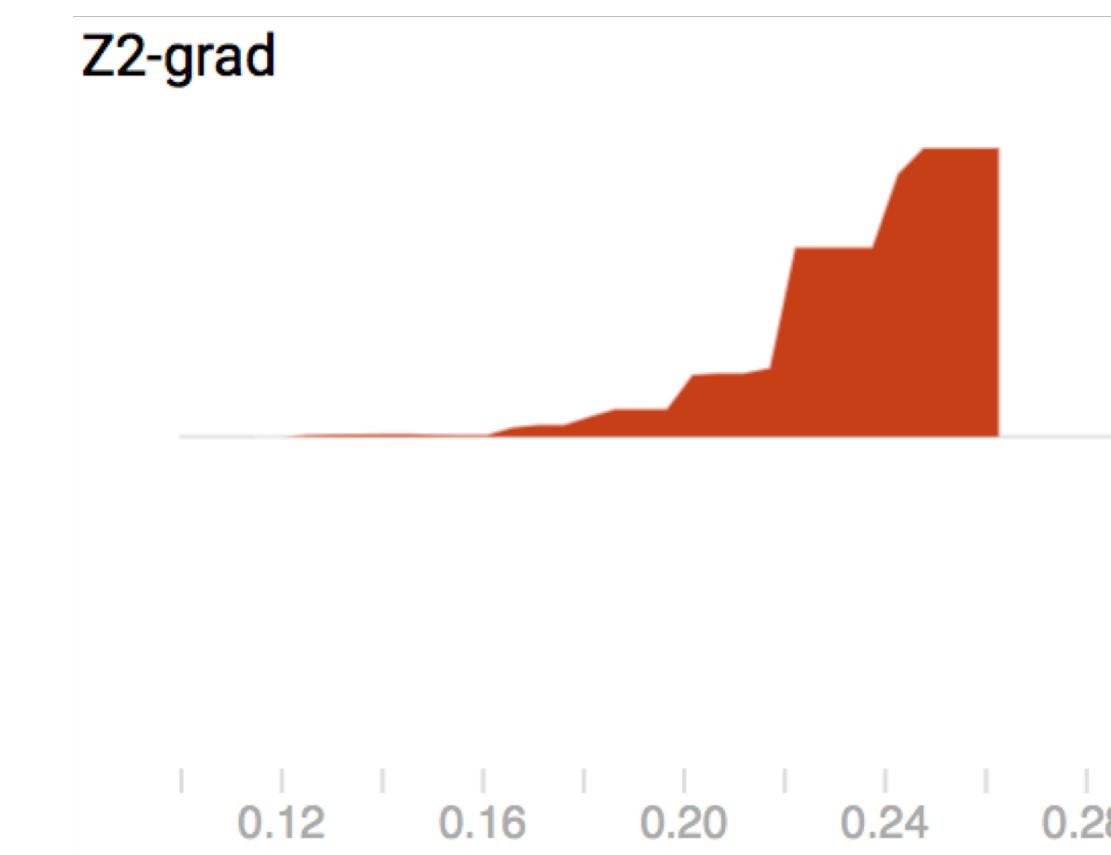
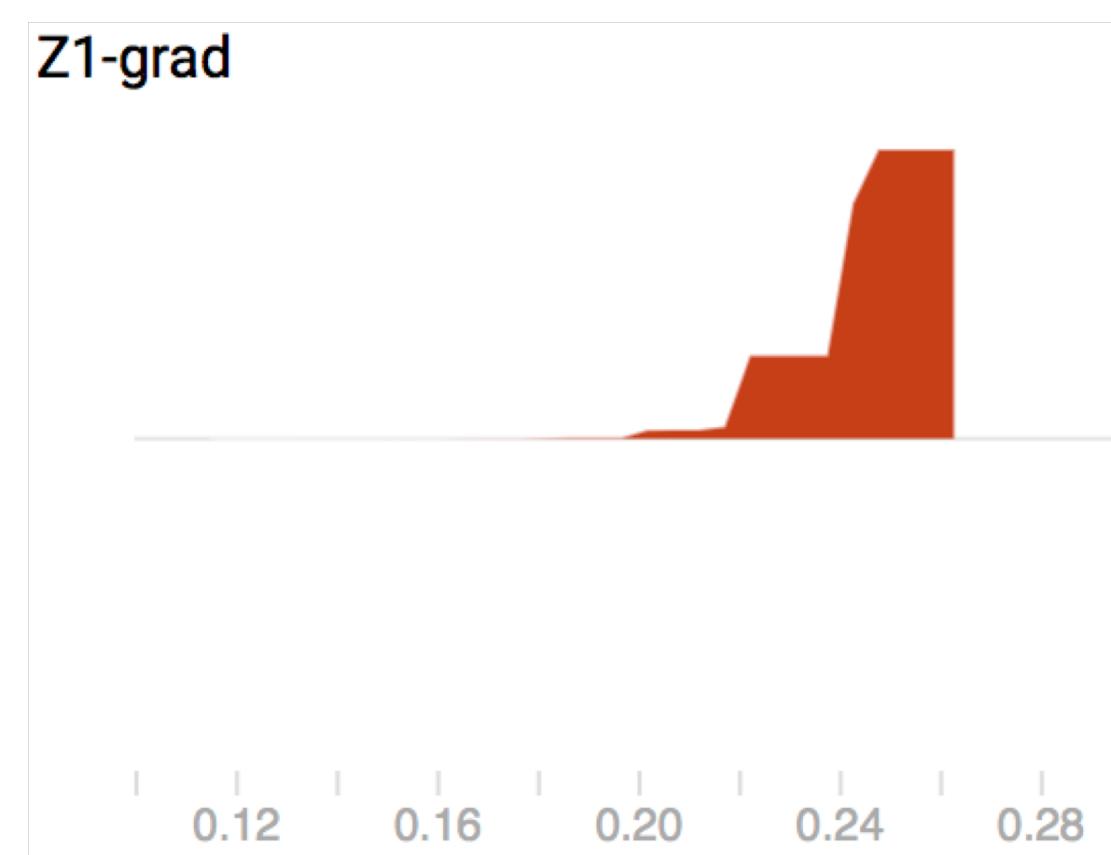
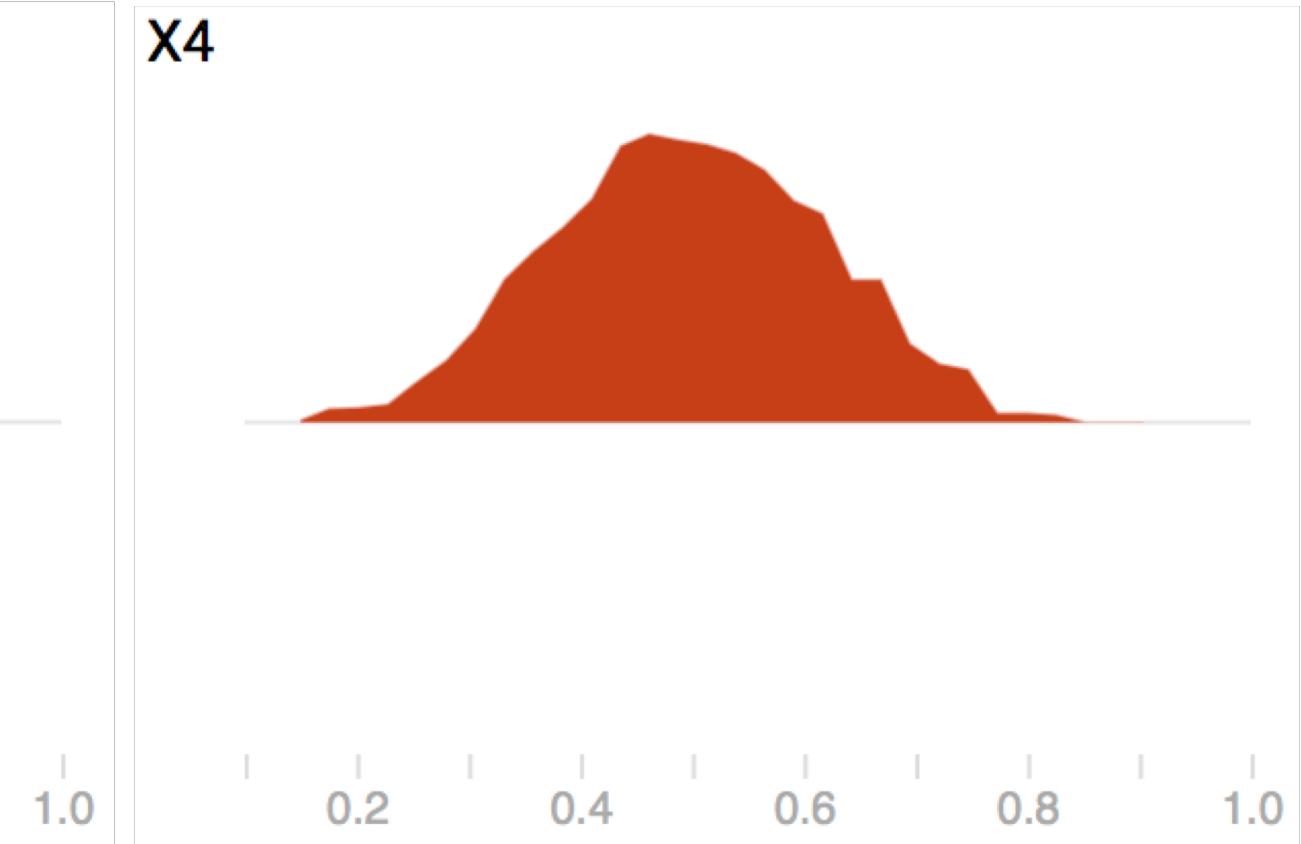
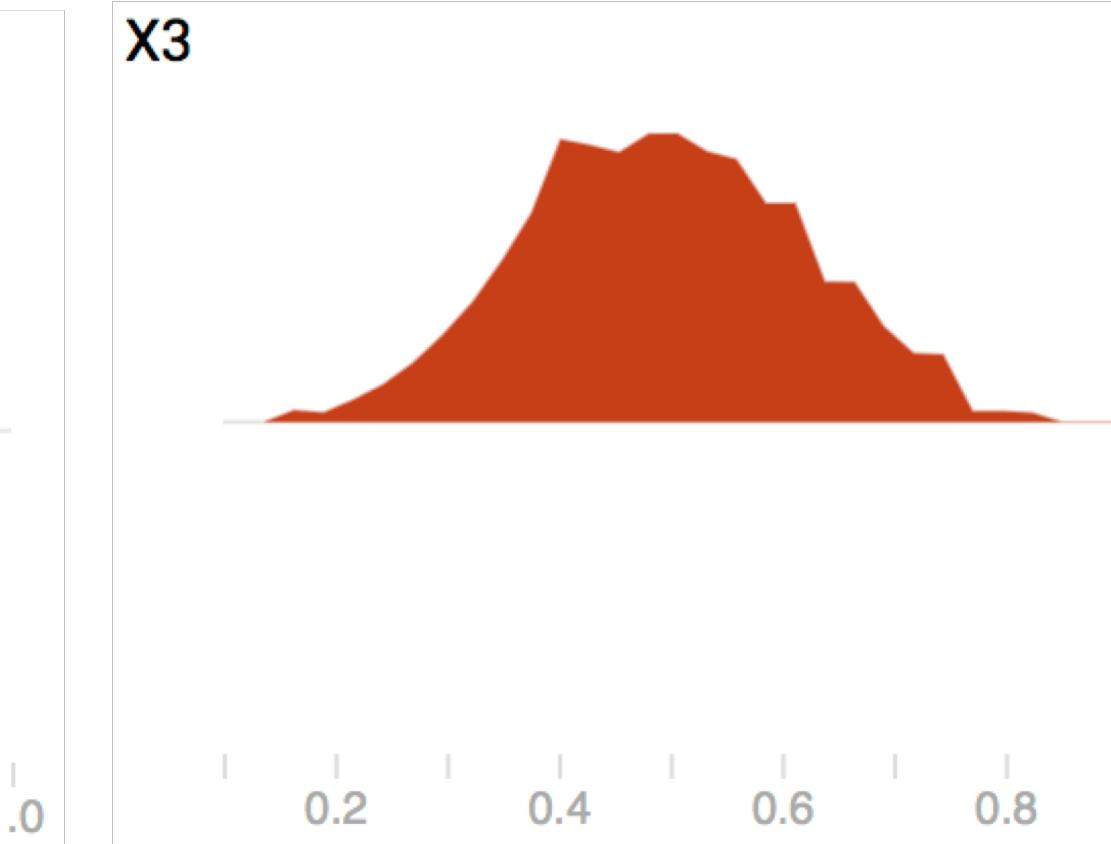
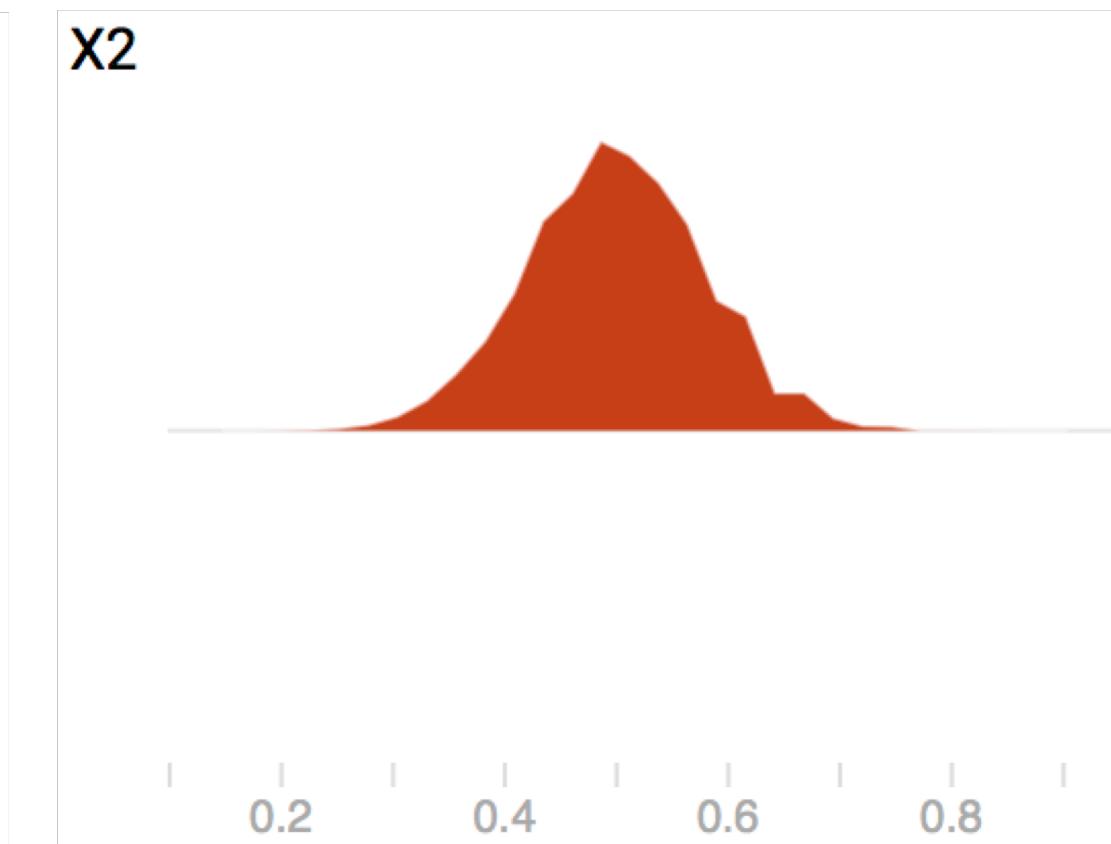
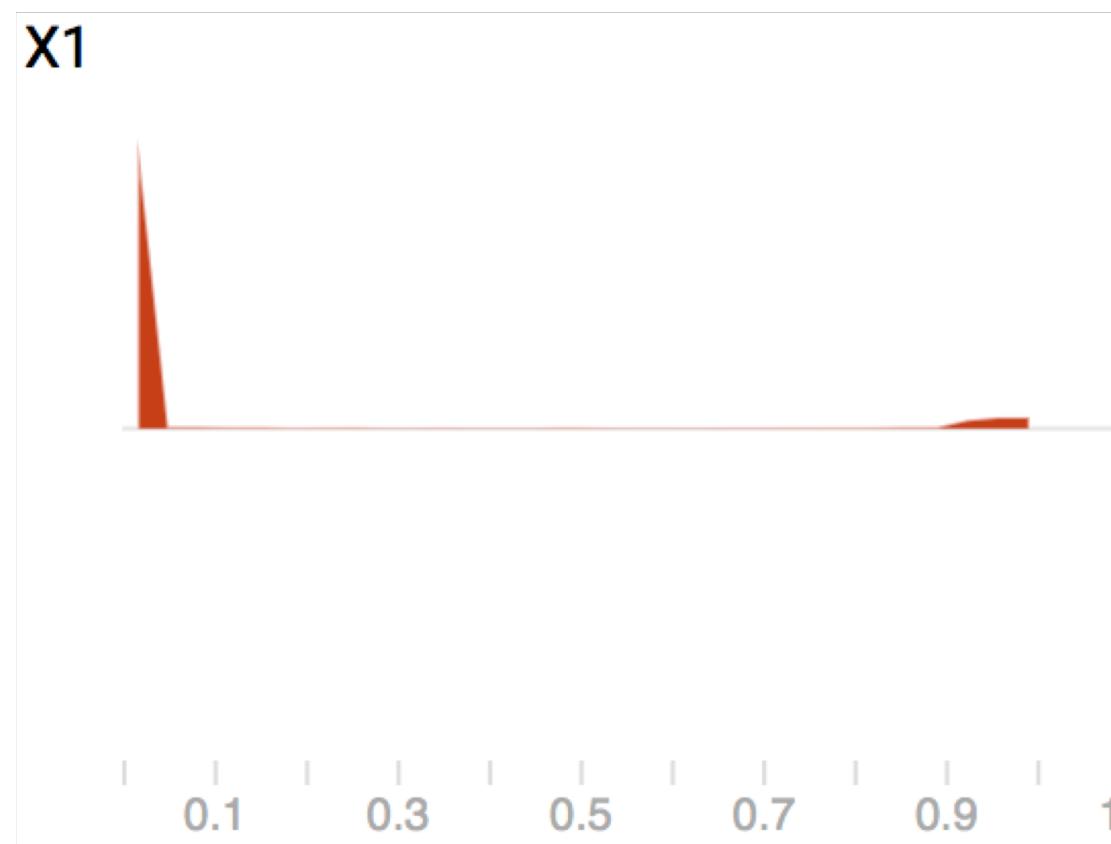
- Xavier initialization with 100 hidden units



Initialization

-Xavier Initialization

- Xavier initialization with 1000 hidden units
 - Properly adapts to the change of the network topology



Activation Functions

-Rectified Linear Units (ReLU)

- We all know that sigmoid functions are not a good choice for deep learning

- Their derivatives block the flow of the BP error

- ReLU doesn't (entirely) block the BP error

- It actually does block about half of them, but the other half survive
 - It actually blocks about half of the FP (forwardpropagation), too, but the other half survive

- Blocking FP path is actually good

- A natural way to enforce sparsity

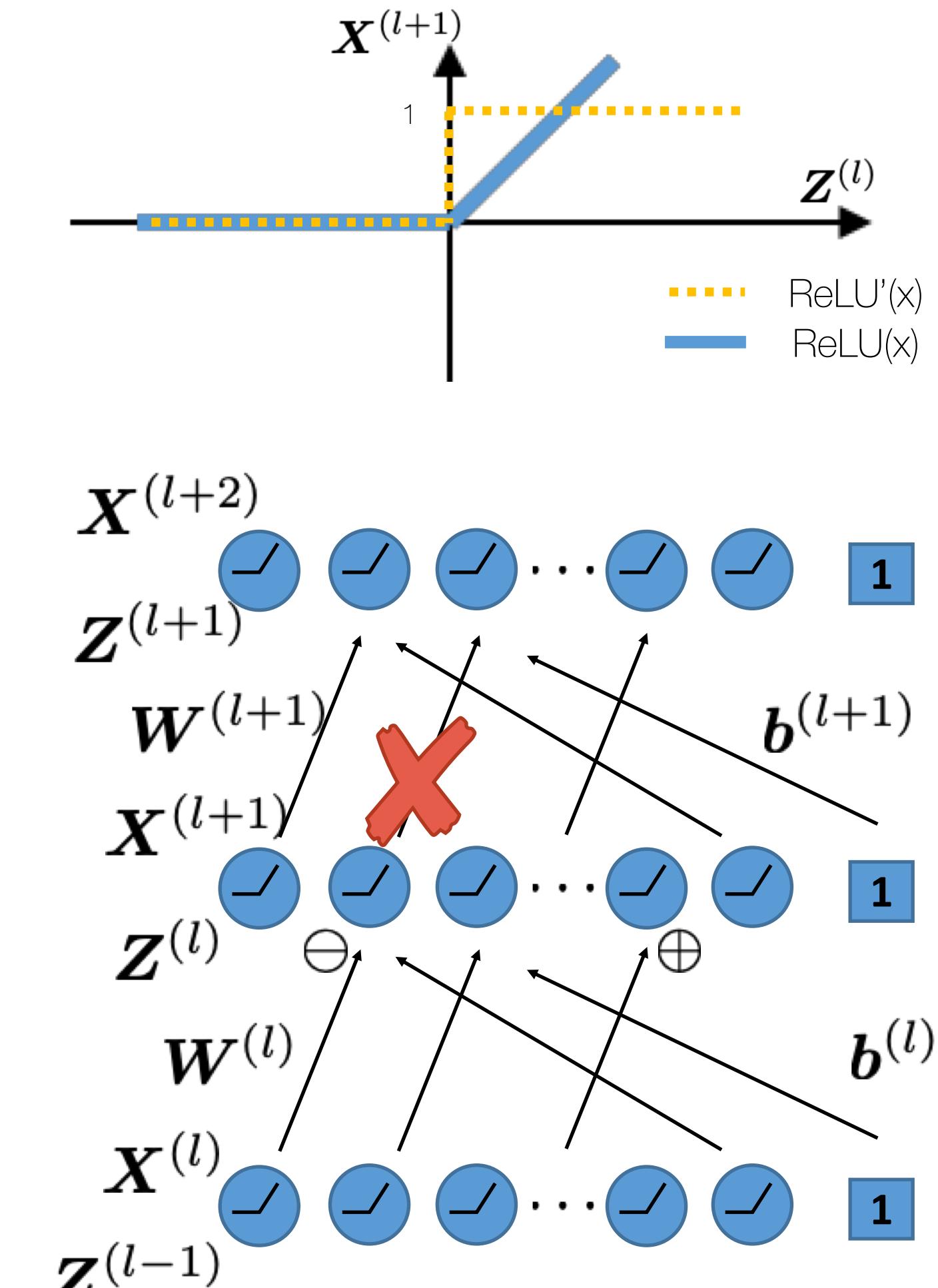
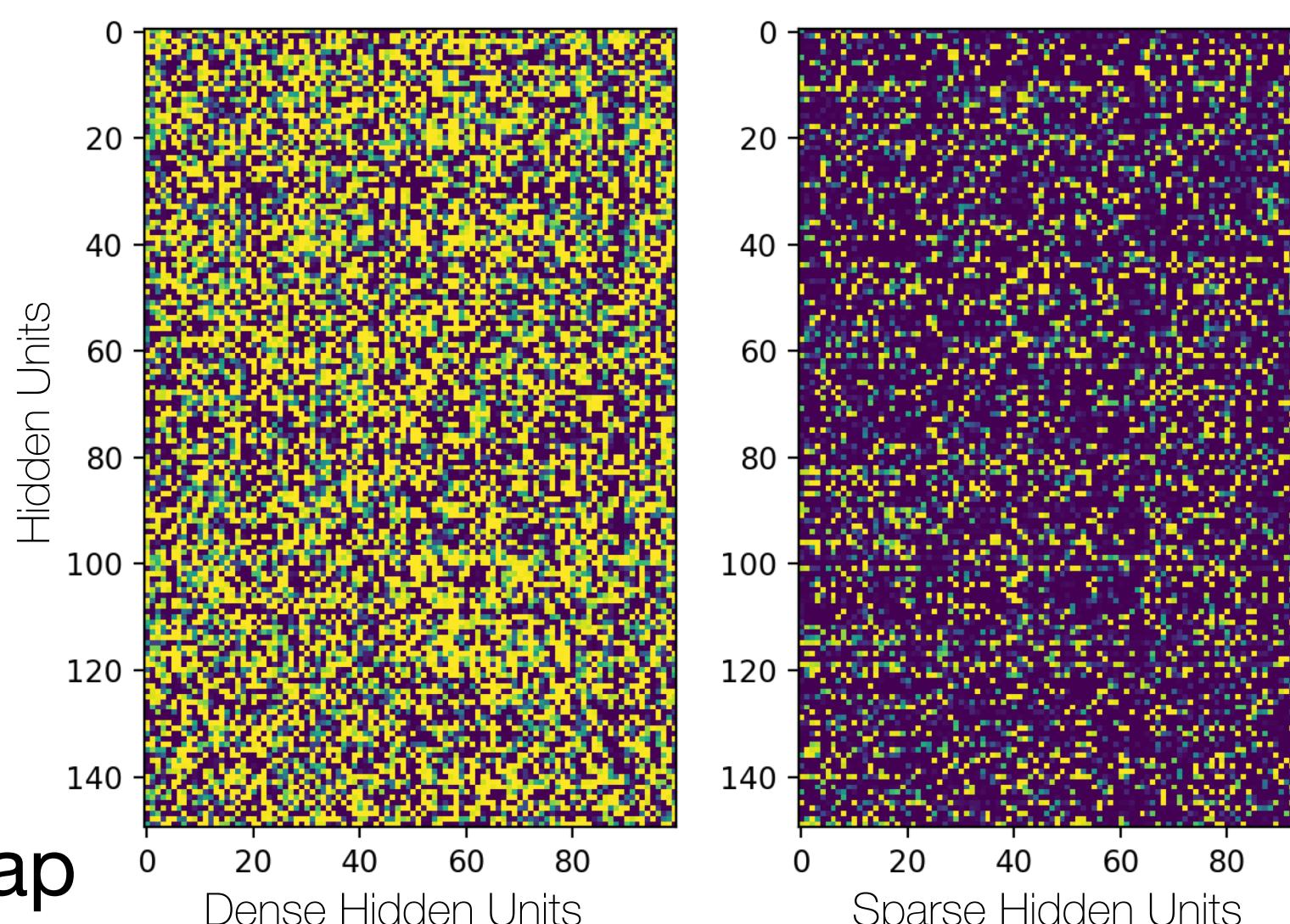
- Blocking BP path is NOT good

- But we do have half the pathways survived

$$\frac{\partial \mathbf{X}_{i,t}^{(l+1)}}{\partial \mathbf{Z}_{i,t}^{(l)}} = \begin{cases} 1 & \mathbf{Z}_{i,t}^{(l)} \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

- Another thing to note is that ReLU is cheap to calculate both for FP and BP

- Let's see what happens



INDIANA UNIVERSITY

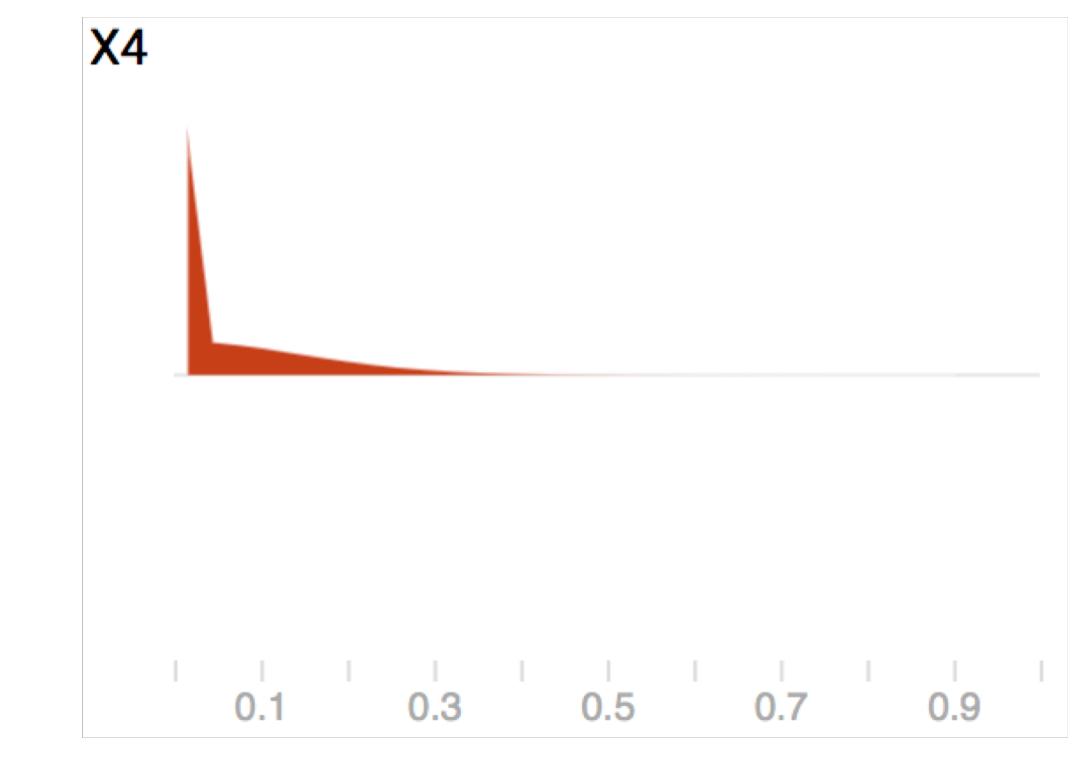
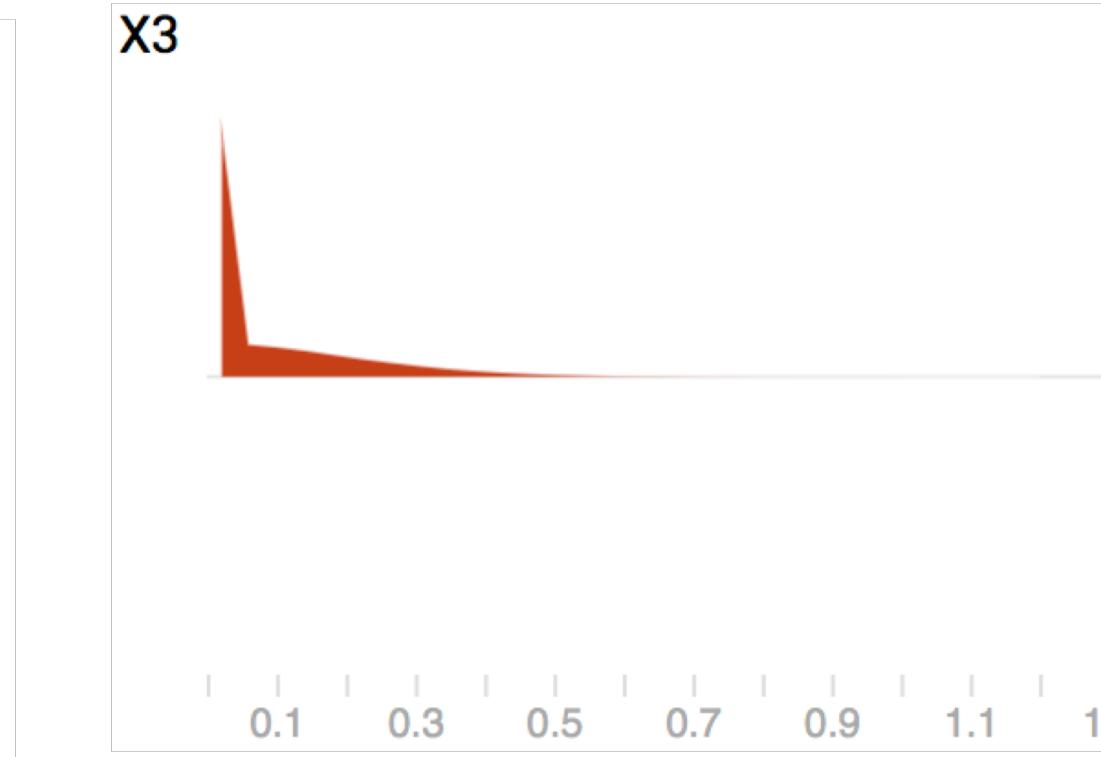
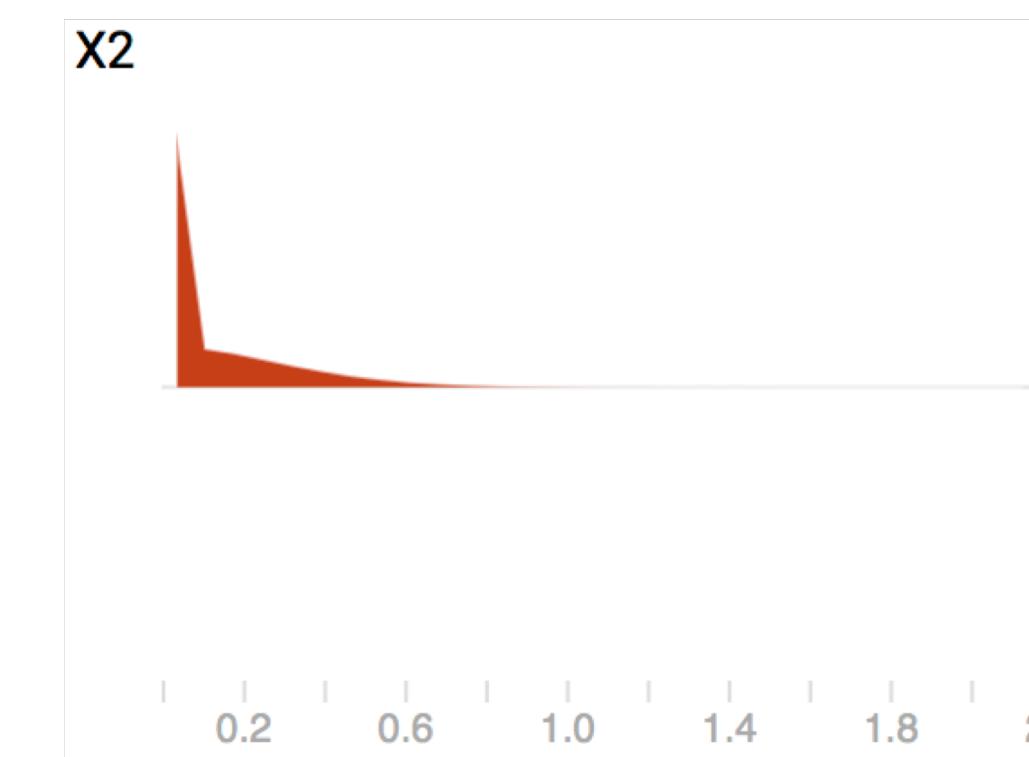
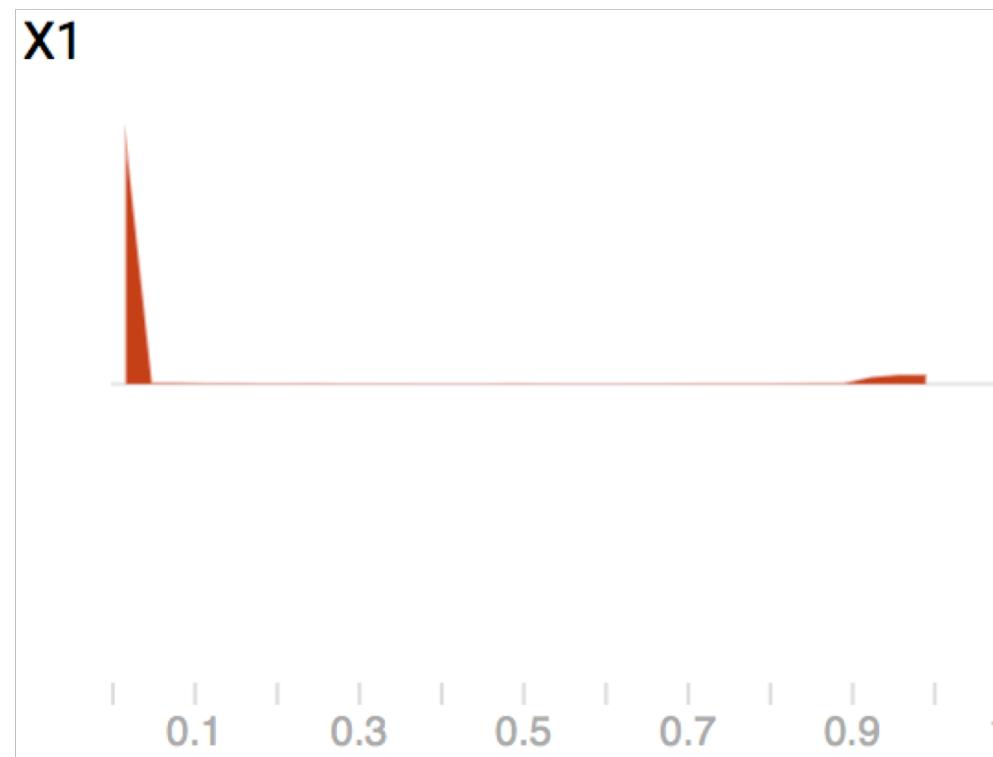
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Xavier Glorot, Antoine Bordes, and Yoshua Bengio. "Deep sparse rectifier neural networks." Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics. 2011.

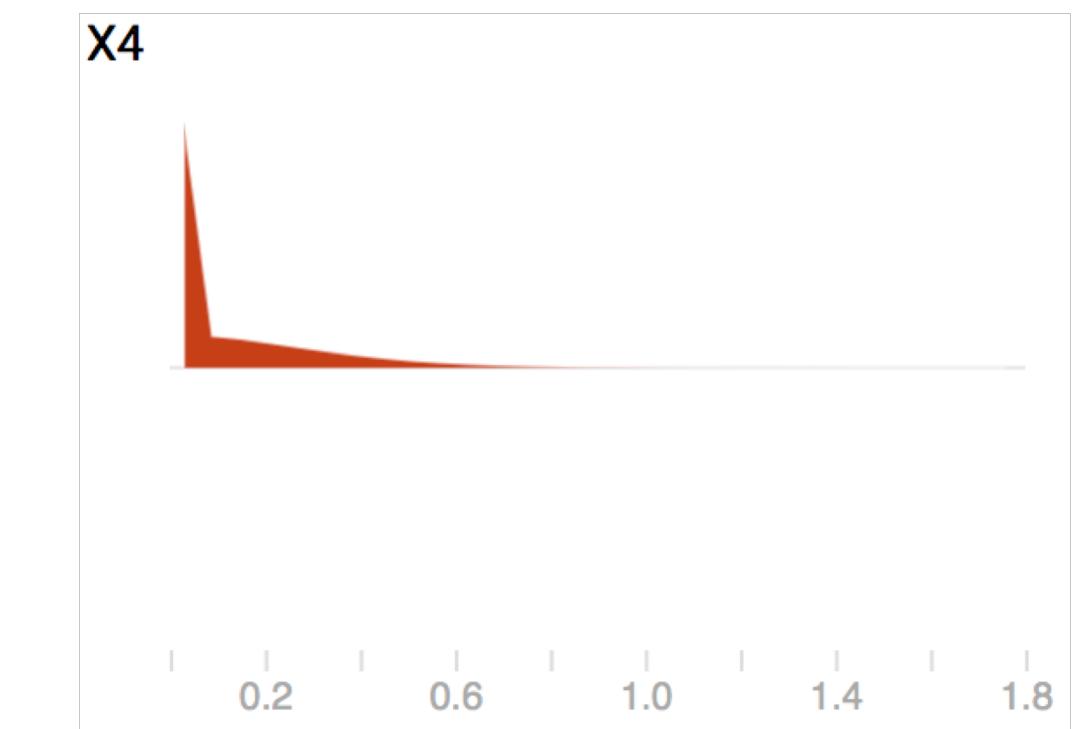
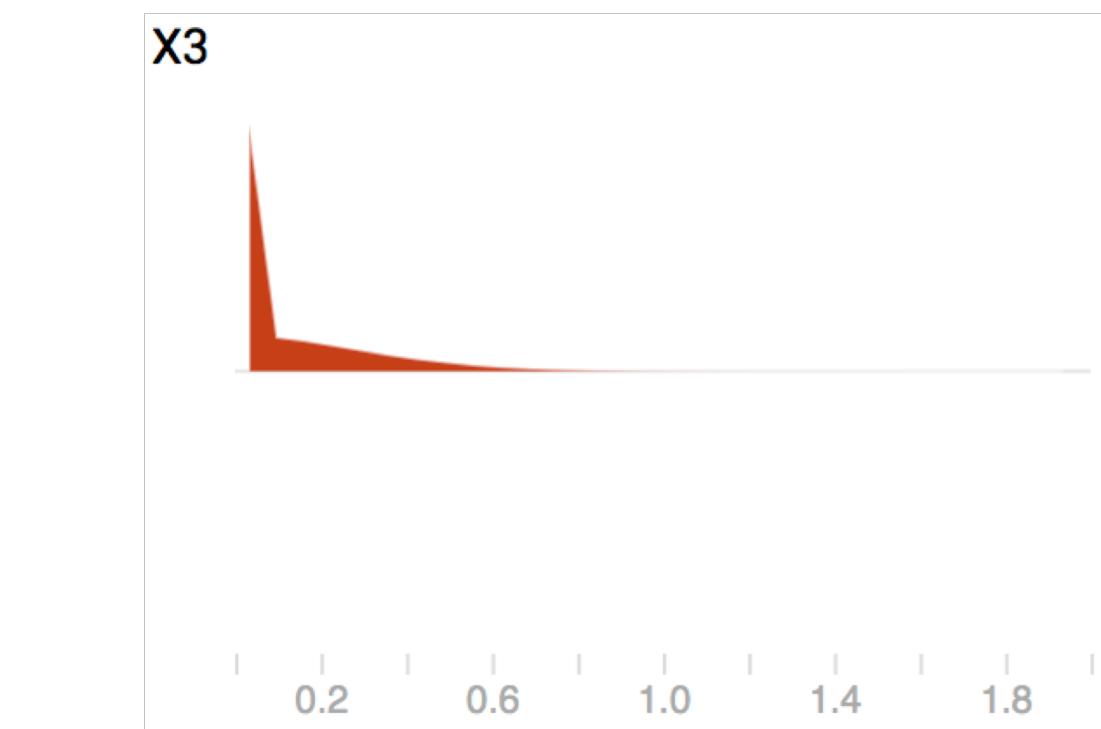
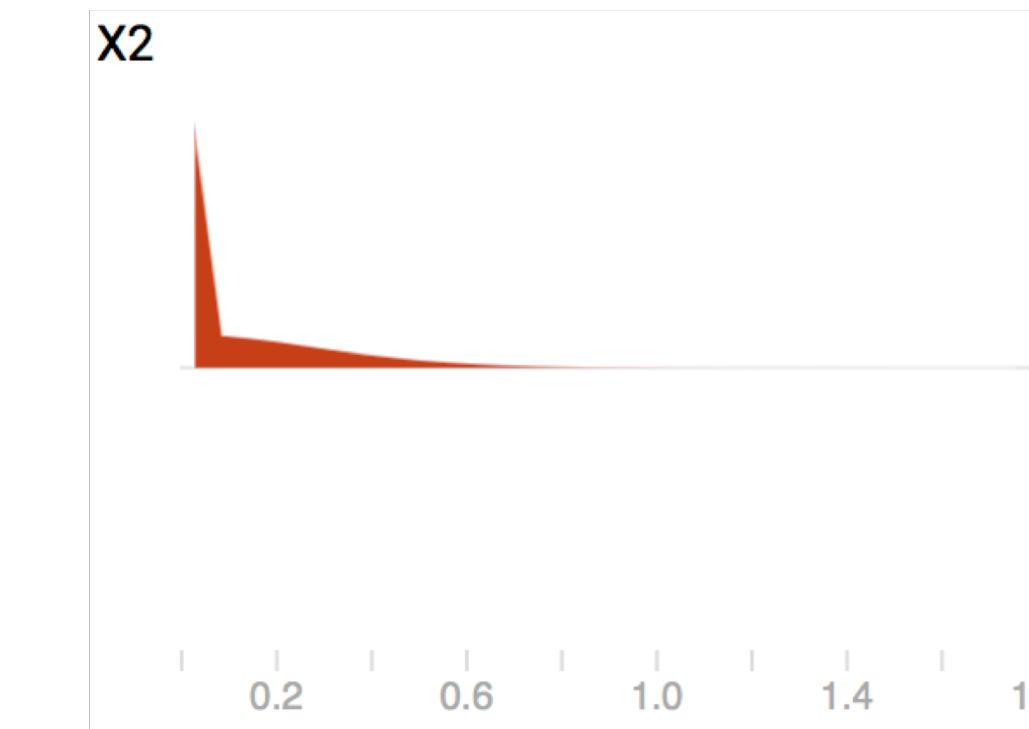
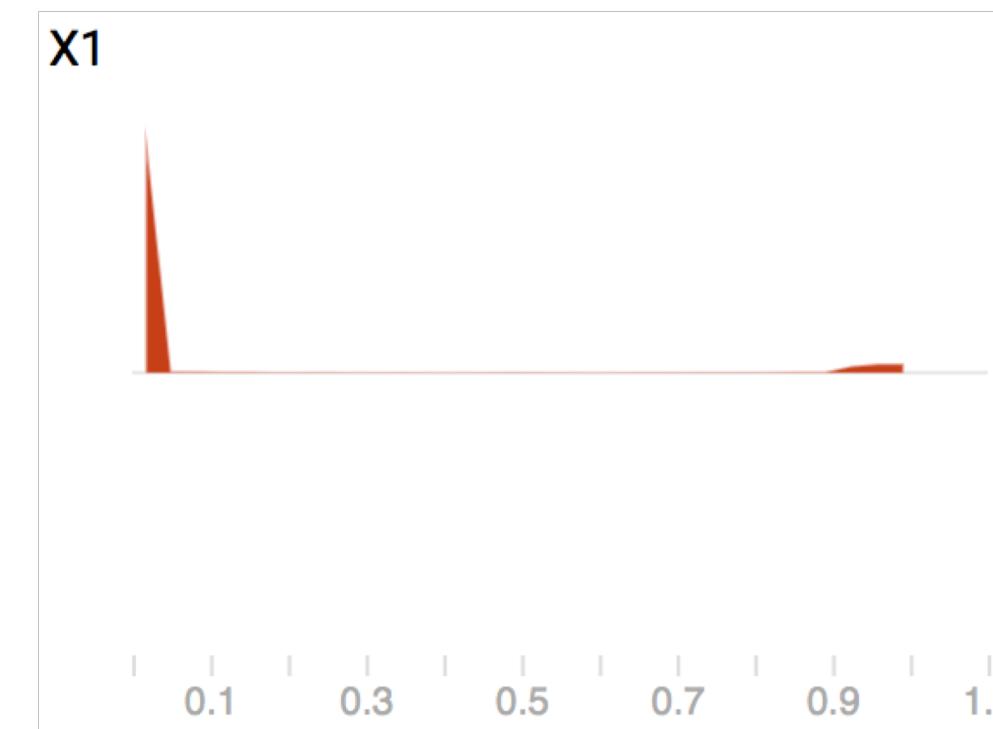
Activation Functions

-Rectified Linear Units (ReLU)

- ReLU with Xavier initialization (notice the different xtick ranges)



- ReLU with He initialization



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

He, Kaiming, et al. "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification." Proceedings of the IEEE international conference on computer vision. 2015.

Activation Functions

-He initialization for ReLU

- Variance of the linear combination $\mathbf{Z}_{i,t}^{(l+1)} = \sum_{j=1}^{N_j} \mathbf{W}_{i,j}^{(l+1)} \mathbf{X}_{j,t}^{(l+1)}$

$$\text{var}(\mathbf{Z}_{i,t}^{(l+1)}) = N_j \text{var}(\mathbf{W}_{i,j}^{(l+1)}) \text{var}(\mathbf{X}_{j,t}^{(l+1)}) + N_j \text{var}(\mathbf{W}_{i,j}^{(l+1)}) [E(\mathbf{X}_{j,t}^{(l+1)})]^2 + N_j \text{var}(\mathbf{X}_{j,t}^{(l+1)}) [E(\mathbf{W}_{i,j}^{(l+1)})]^2$$

This is correct only if your r.v.'s on the right hand side is centered

- We have no control over the units, but can still center the weights $\therefore [E(\mathbf{W}_{i,j}^{(l+1)})]^2 = 0$

$$\text{var}(\mathbf{Z}_{i,t}^{(l+1)}) = N_j \text{var}(\mathbf{W}_{i,j}^{(l+1)}) \left(\text{var}(\mathbf{X}_{j,t}^{(l+1)}) + [E(\mathbf{X}_{j,t}^{(l+1)})]^2 \right) = N_j \text{var}(\mathbf{W}_{i,j}^{(l+1)}) [E(\mathbf{X}_{j,t}^{(l+1)})^2] \stackrel{?}{=} \text{var}(\mathbf{Z}_{j,t}^{(l)})$$

- We look for $\text{var}(\mathbf{W}_{i,j}^{(l+1)})$ that meets the final equation

- Let's see what happens in the lower layer

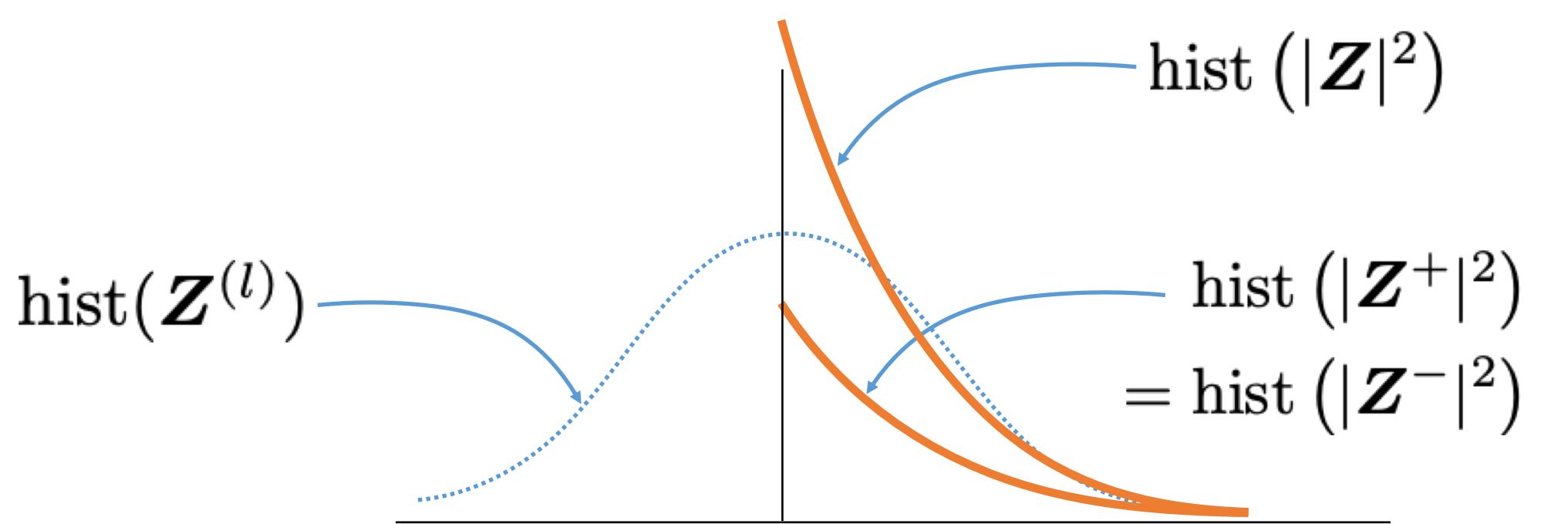
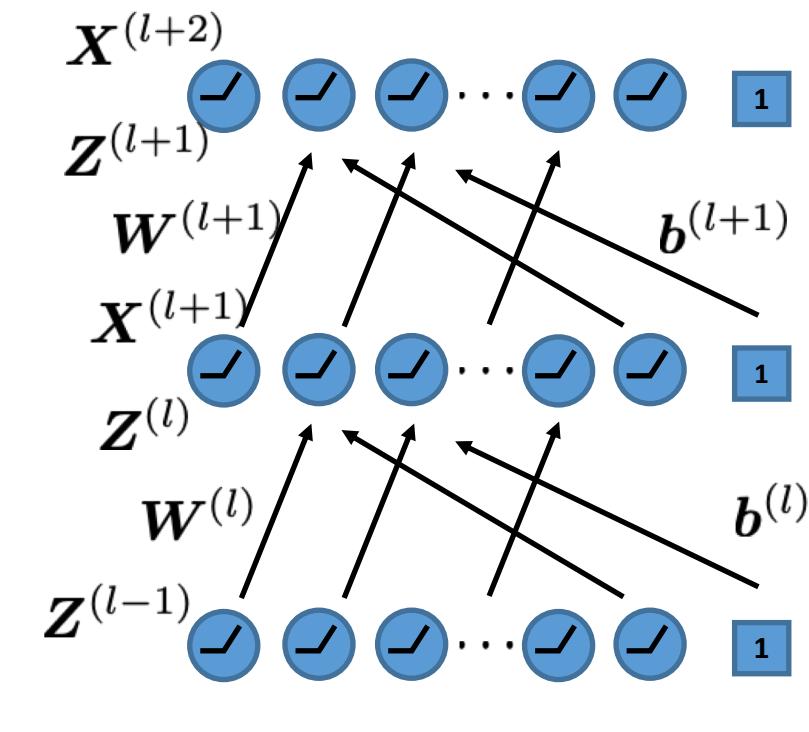
- Since we assume a centered symmetric distribution for $\mathbf{W}_{i,j}^{(l)}$, should $P(\mathbf{Z}_{i,j}^{(l)})$ be, too

$$\begin{aligned} E(\mathbf{X}_{j,t}^{(l+1)^2}) &= E(\text{ReLU}(\mathbf{Z}_{j,t}^{-^{(l)}})^2) / 2 + E(\text{ReLU}(\mathbf{Z}_{j,t}^{+^{(l)}})^2) / 2 \\ &= E((\mathbf{Z}_{j,t}^{+^{(l)}})^2) / 2 = E(\mathbf{Z}_{j,t}^{(l)^2}) / 2 = \text{var}(\mathbf{Z}_{j,t}^{(l)}) / 2 \end{aligned}$$

- Therefore,

$$\text{var}(\mathbf{Z}_{i,t}^{(l+1)}) = \frac{N_j}{2} \text{var}(\mathbf{W}_{i,j}^{(l+1)}) \text{var}(\mathbf{Z}_{j,t}^{(l)}) \stackrel{?}{=} \text{var}(\mathbf{Z}_{j,t}^{(l)})$$

$$\text{var}(\mathbf{W}_{i,j}^{(l+1)}) = \frac{2}{N_j}$$



INDIANA UNIVERSITY

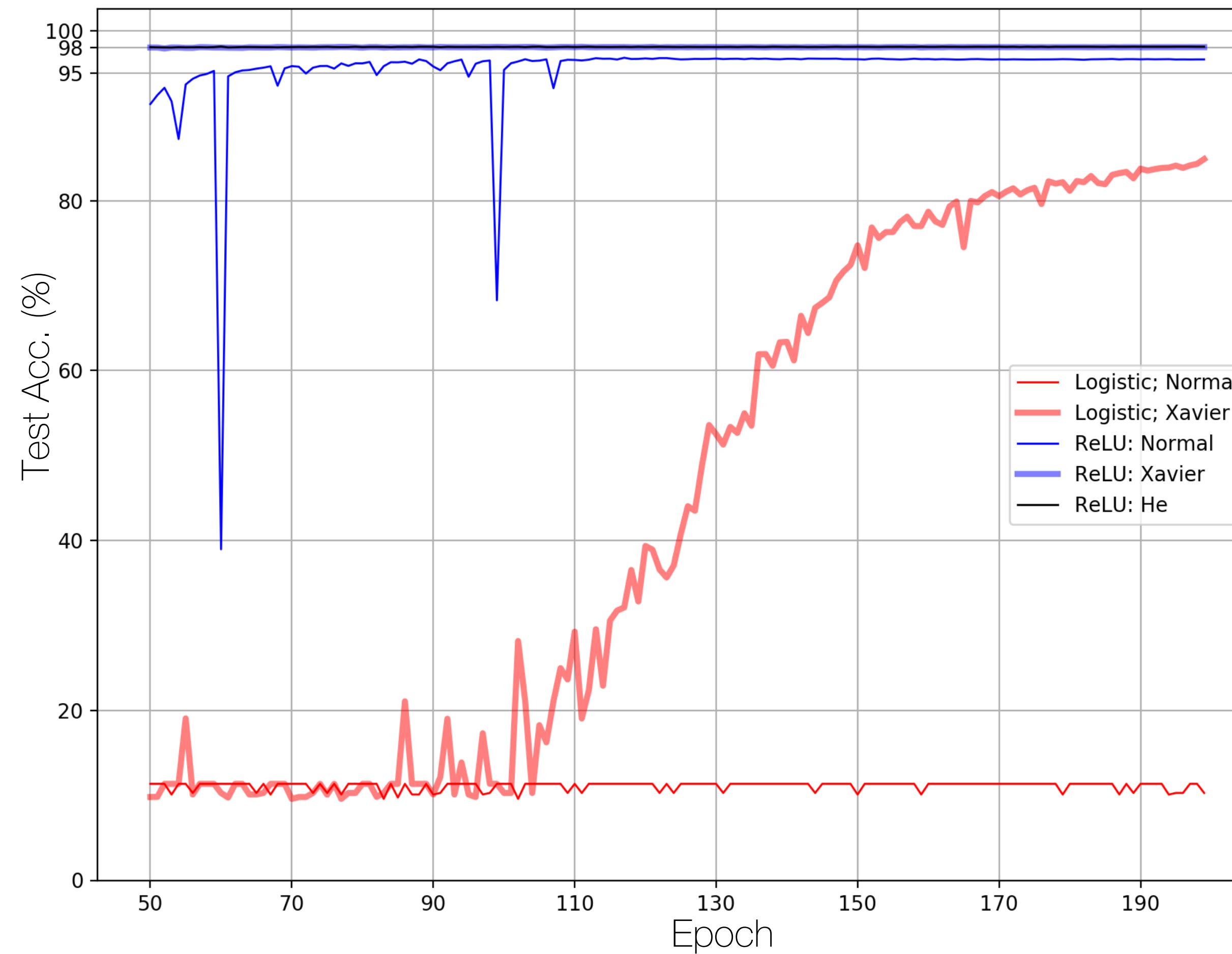
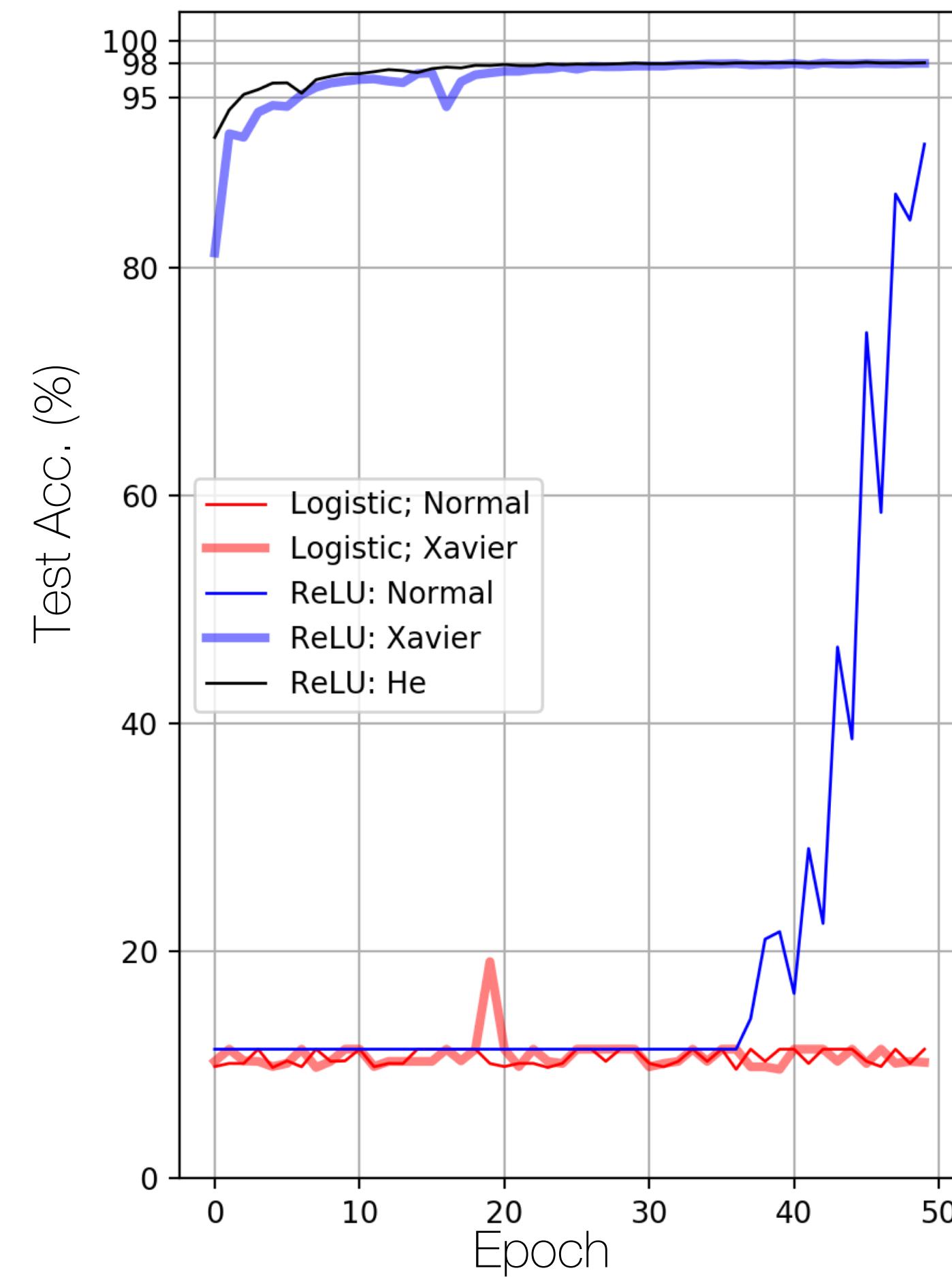
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

He, Kaiming, et al. "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification." Proceedings of the IEEE international conference on computer vision. 2015.

Activation Functions

-He initialization for ReLU

- For a 512X5 network for MNIST



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Activation Functions

-The other activation functions

Leaky ReLU

$$f(x) = \max(x, 0) + a \min(0, x)$$

Parametric ReLU

$$f(x) = \max(x, 0) + a \min(0, x)$$

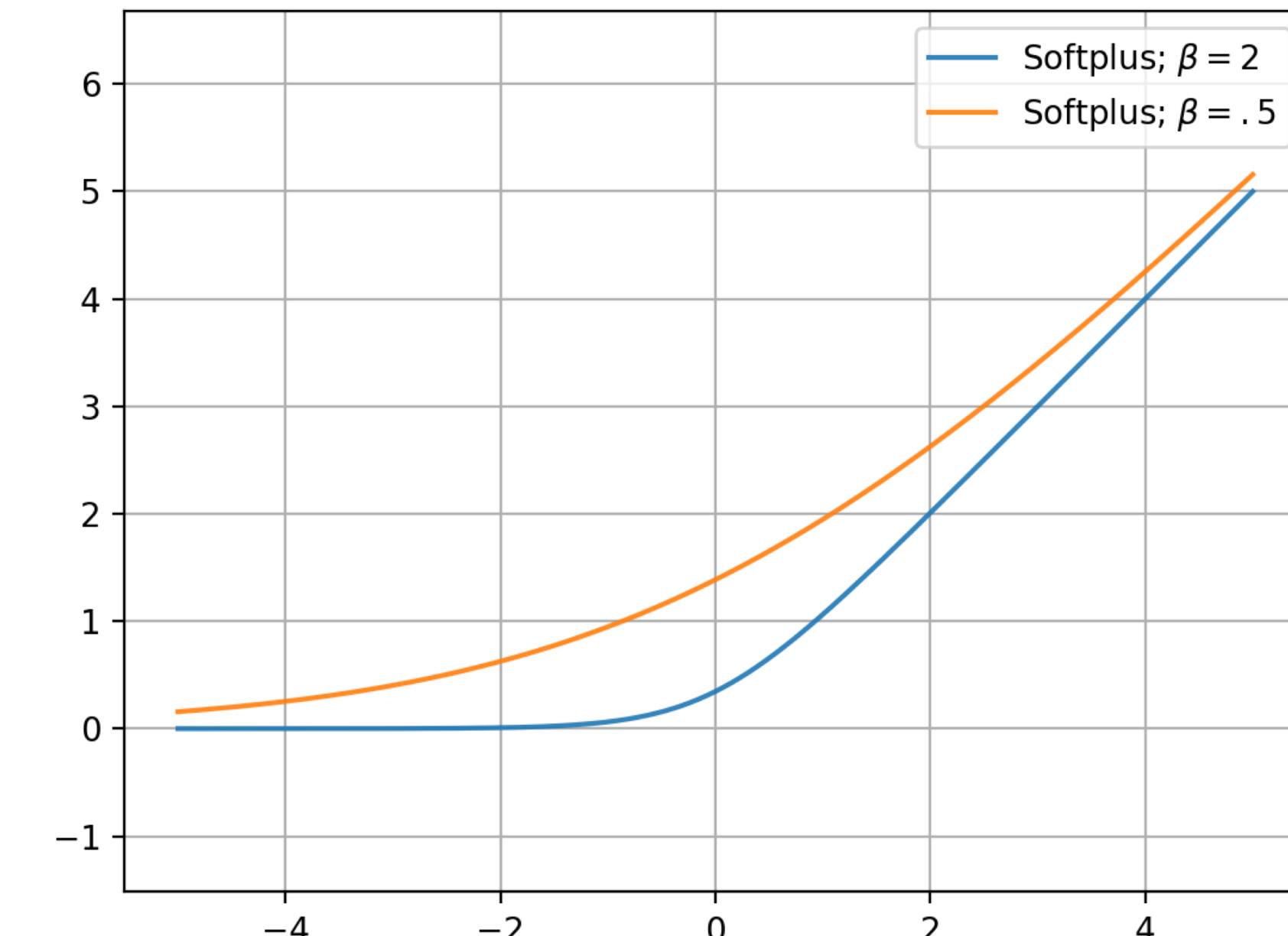
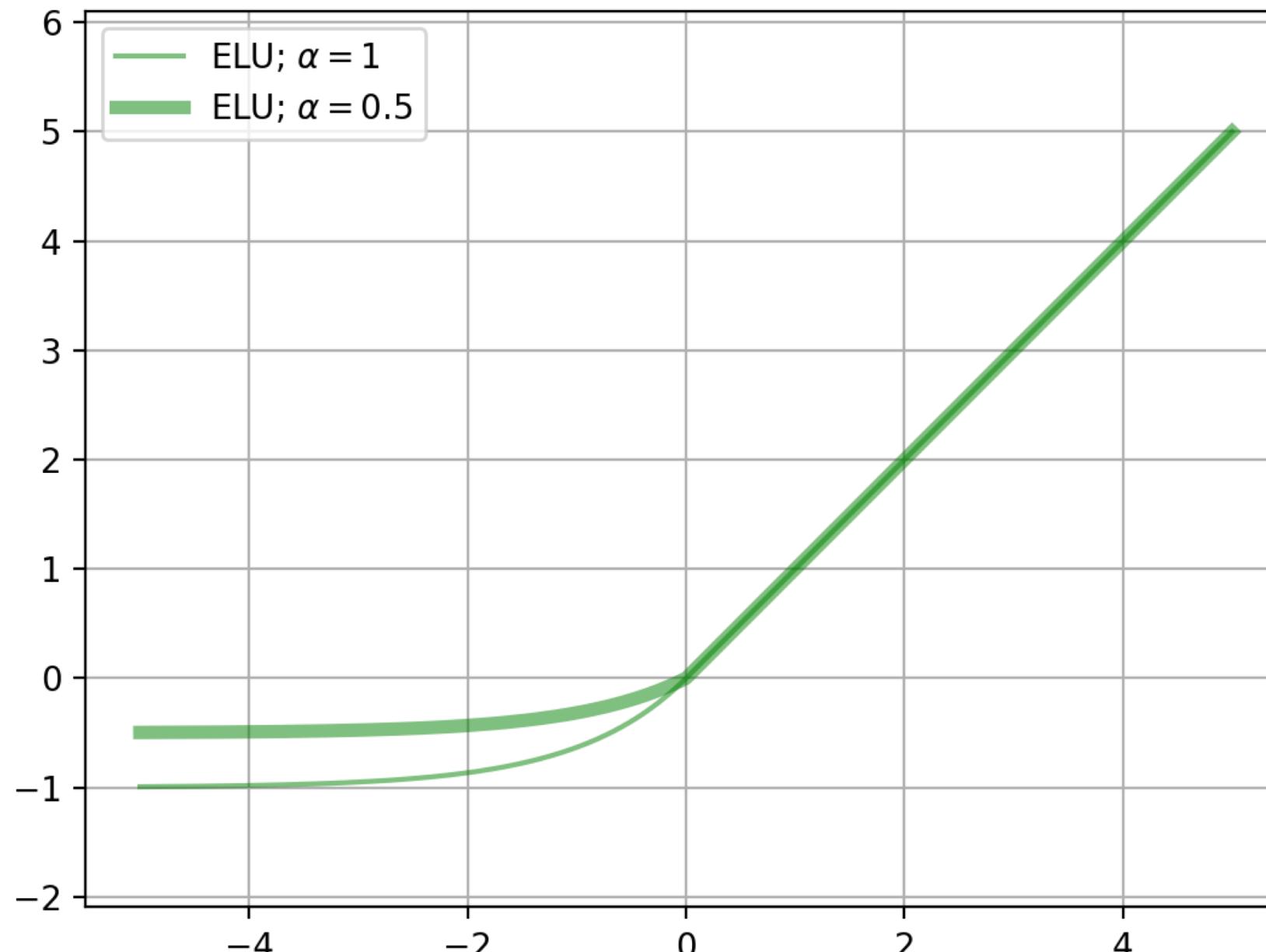
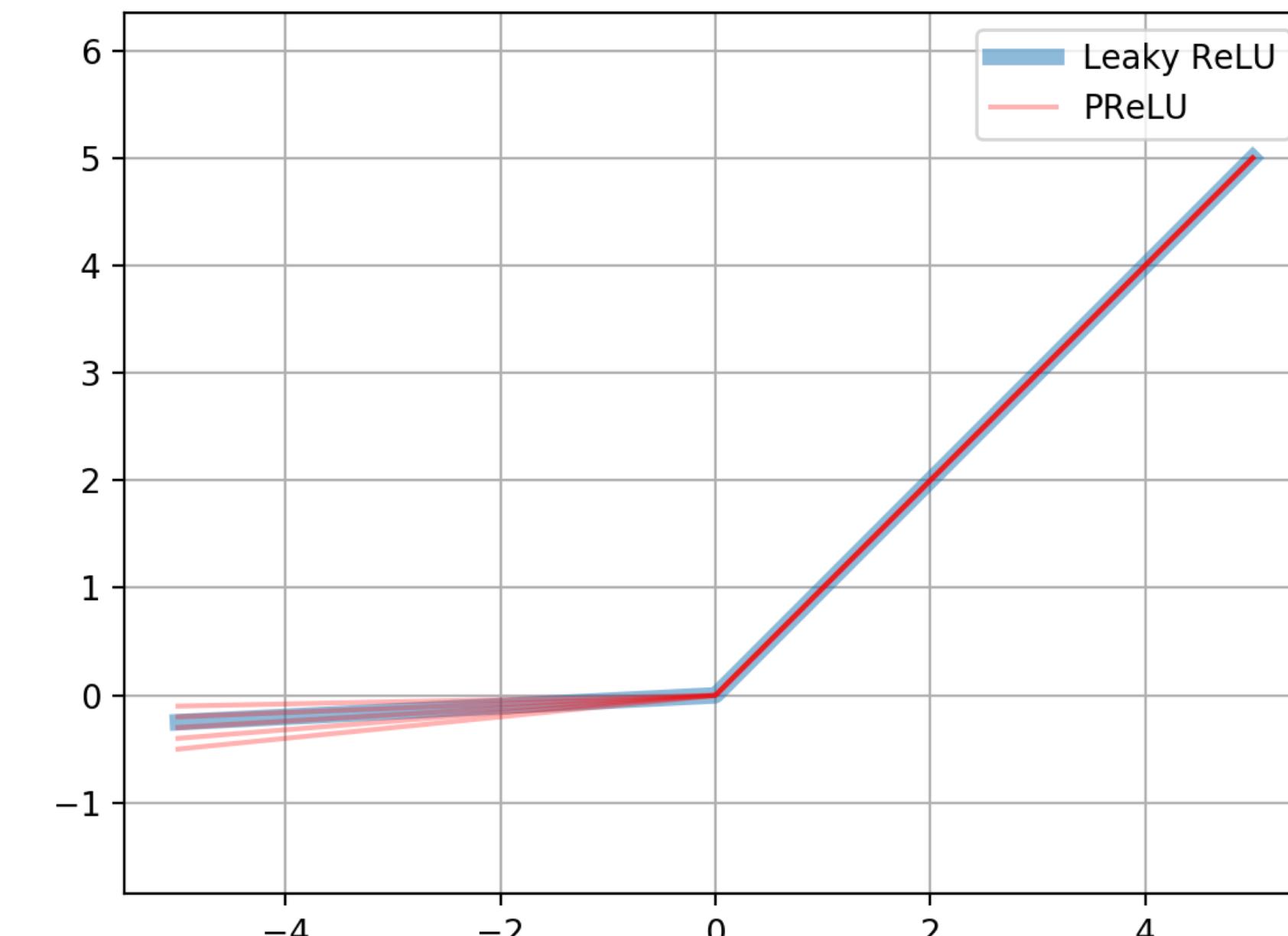
Softplus

$$f(x) = \frac{1}{\beta} \log(1 + \exp(\beta x))$$

Exponential LU

$$f(x) = \max(x, 0) + \min(0, \alpha(\exp(x) - 1))$$

You name it!



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING