

# **ENGR-E 533 “Deep Learning Systems”**

## **Lecture 05: Vanishing Gradients and Pretraining**

**Minje Kim**

Department of Intelligent Systems Engineering

Email: [minje@indiana.edu](mailto: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**

# What We've Learned So Far

- Feature extraction is important
  - But difficult
  - Usually involves unsupervised learning
- Supervised learning
  - Is easy if you start from a good set of features
  - Otherwise you'll need some nonlinearity
- A shallow neural network takes care of both tasks at the same time
  - The first layer learns features
  - The last layer learns the regressor
- Universal Approximation Theorem
  - If there are many many hidden units, a shallow neural network can approximate any (nonlinear) function
  - e.g. Sinusoidal functions
  - But how many?



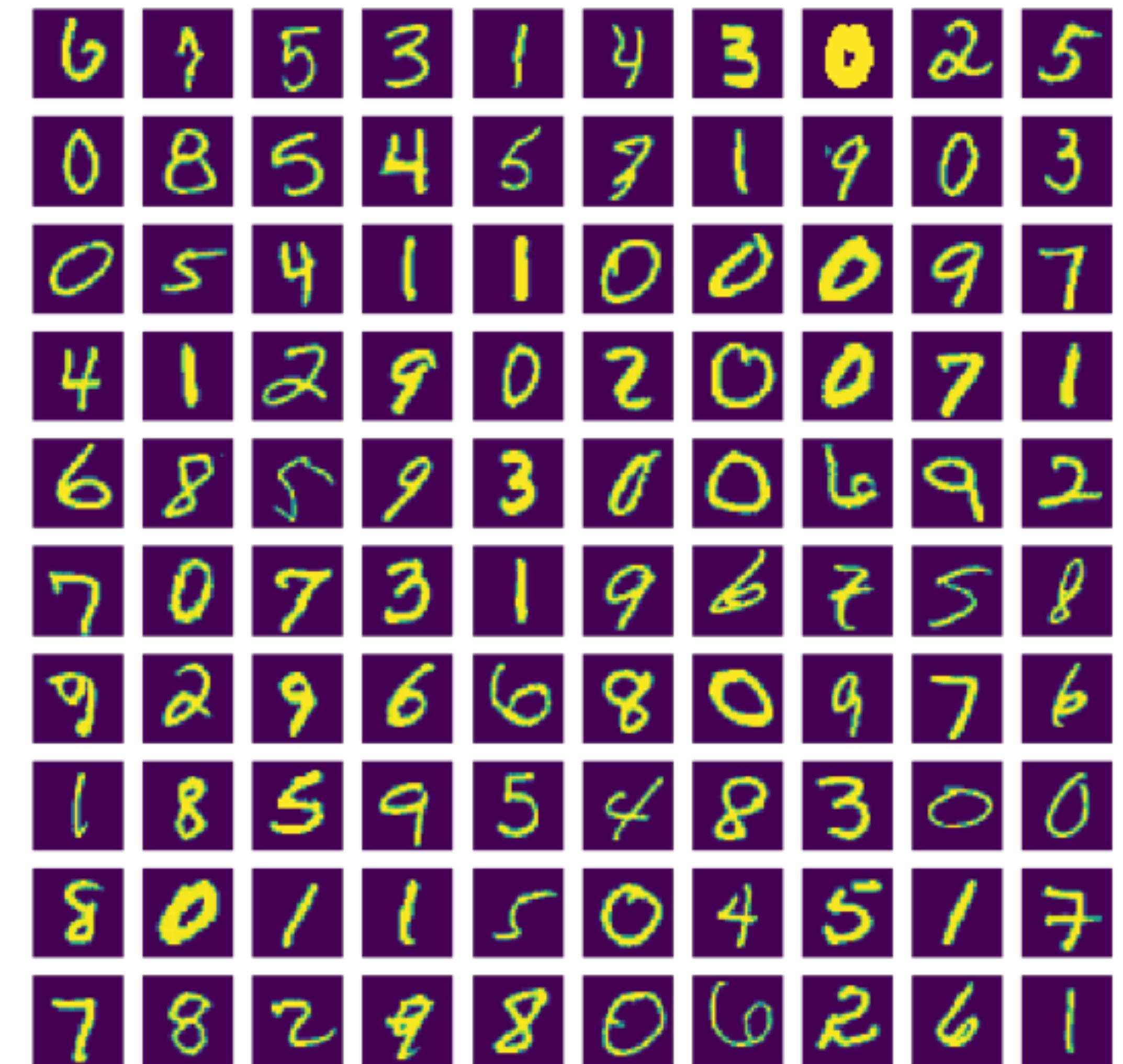
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# A Little More Difficult Problem

## -MNIST handwritten digit recognition

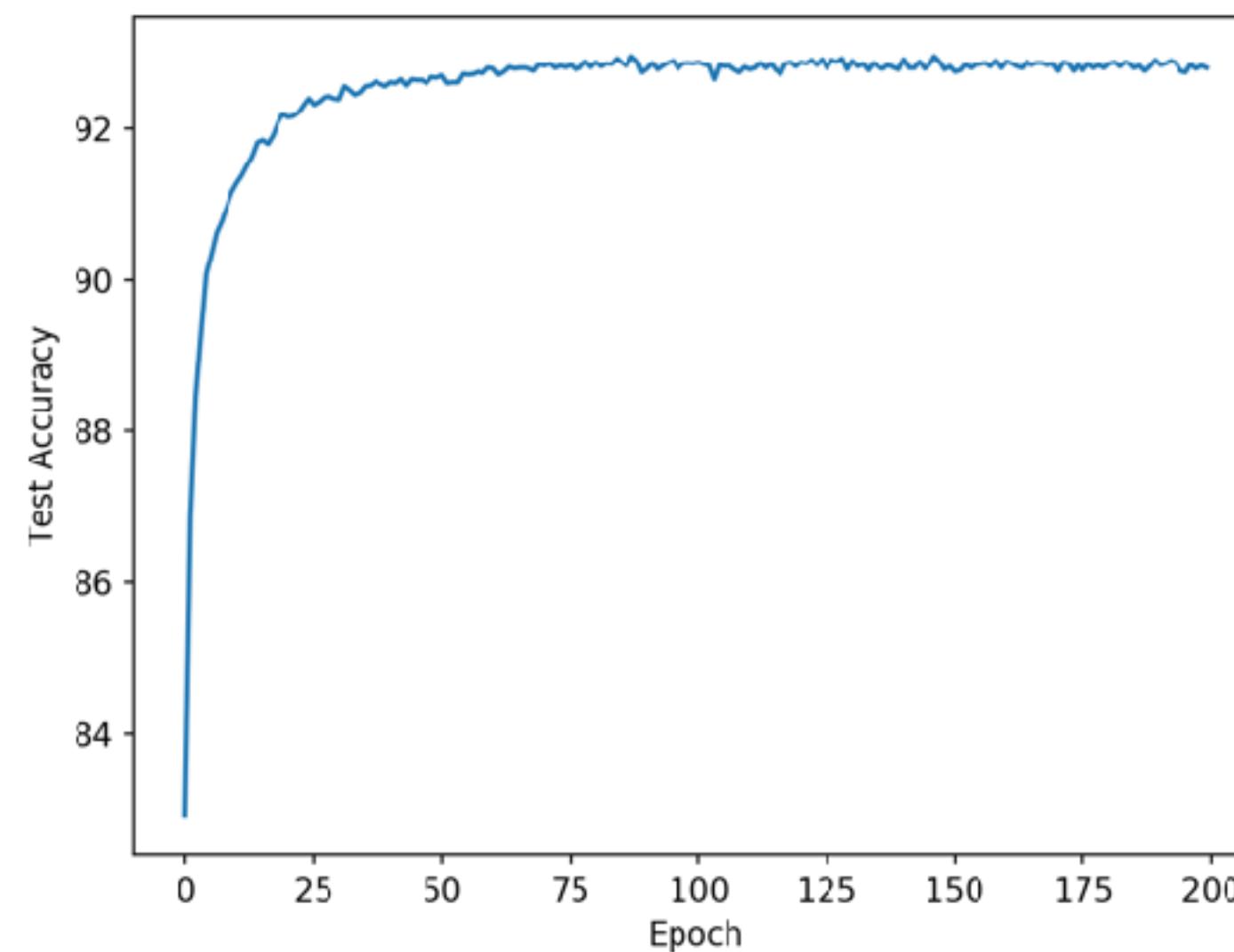
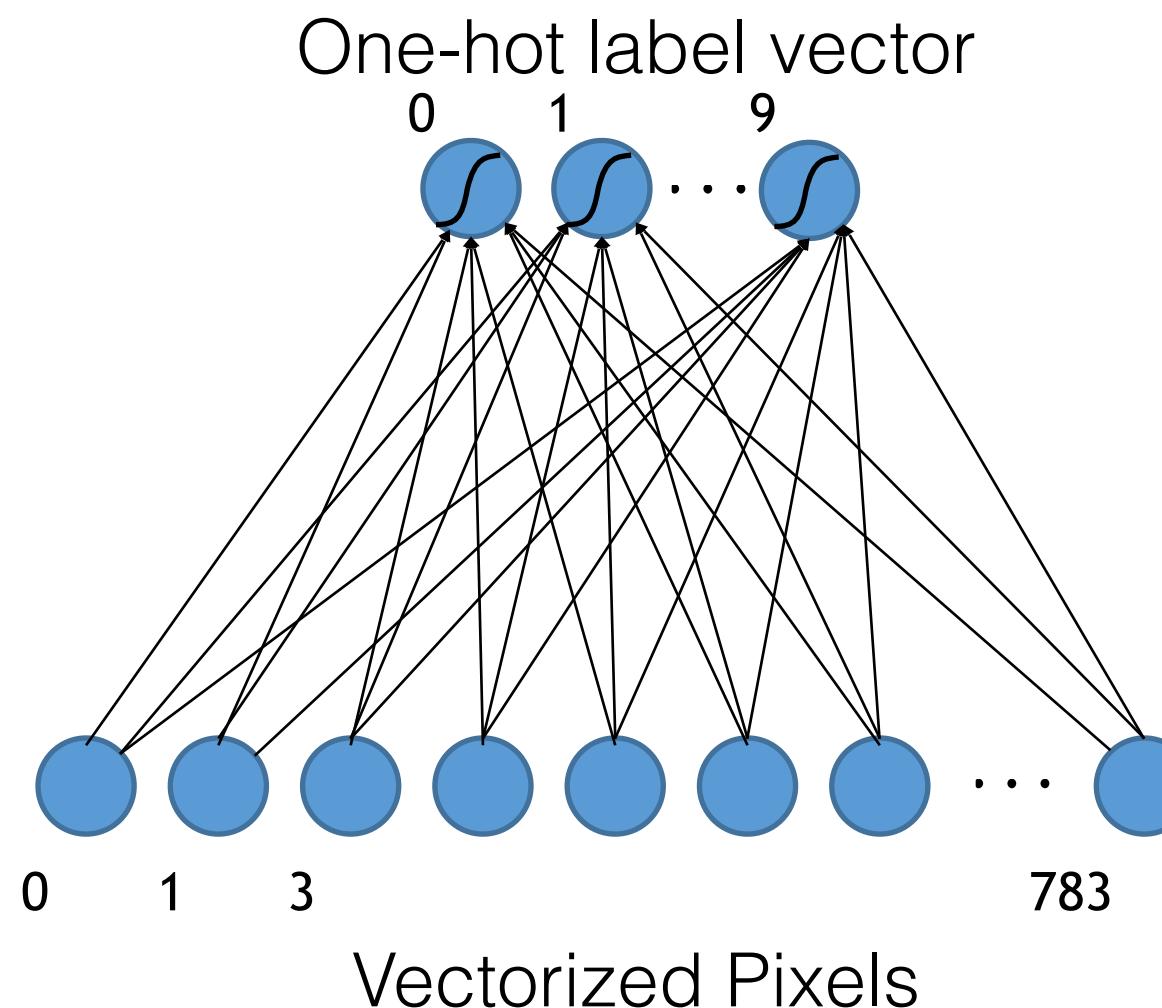
- Considered as the “Hello, World!” problem in deep learning
- 28X28=784 pixels per image
- 55,000 training images
- 10,000 testing images



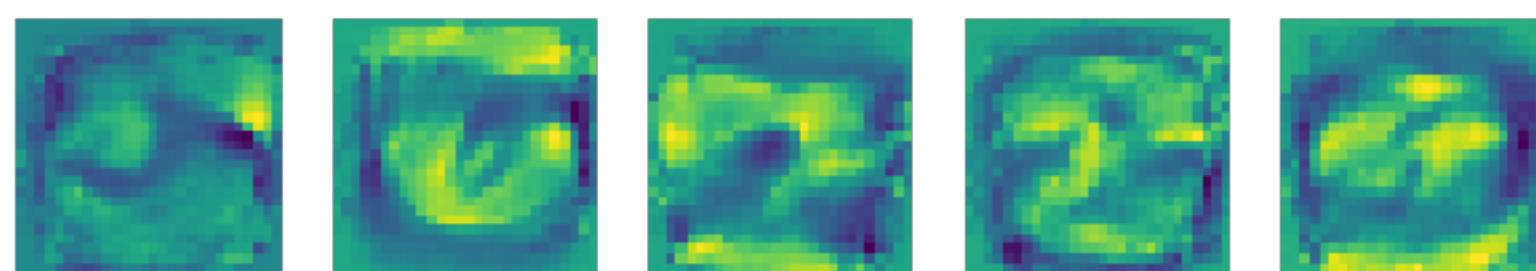
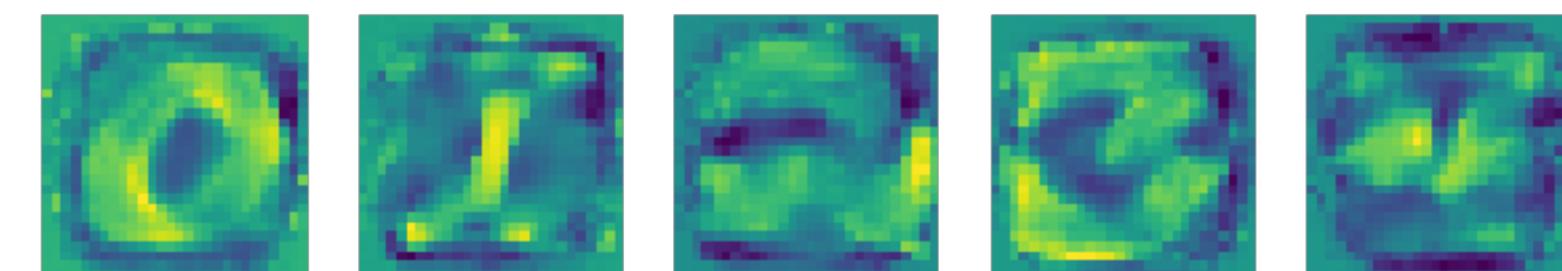
# A Little More Difficult Problem

-MNIST handwritten digit recognition

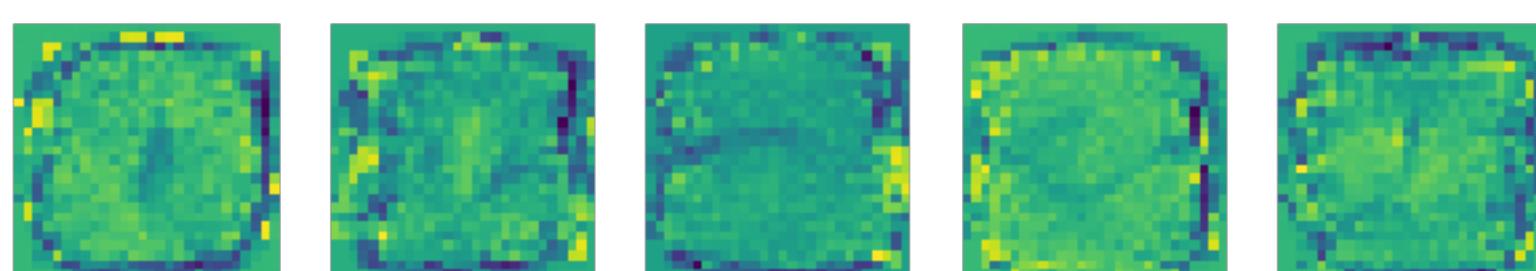
- A softmax classification with no feature extraction
  - i.e. Raw pixel intensities are the features



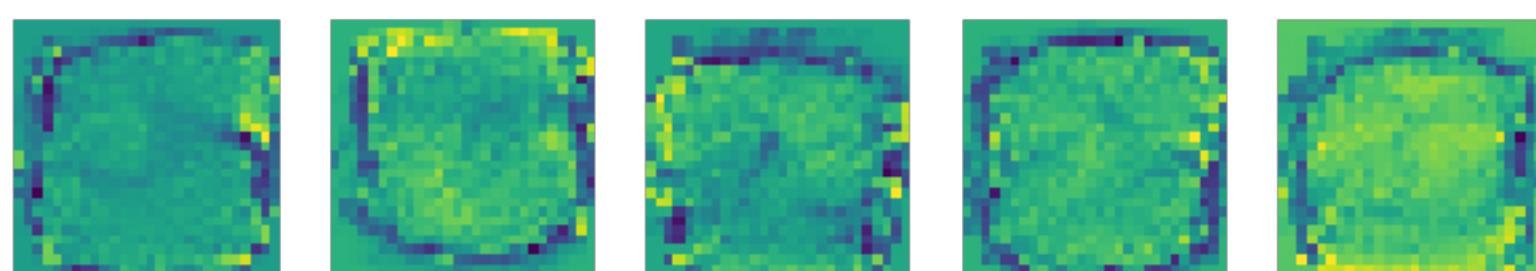
- It's not bad, but how about some feature extraction?
  - YALT last time



Weights at 20th epoch



Weights at 200th epoch



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Feature Extraction

## -Sparse AE

- Let's extract some features, and then do the softmax

- Two cases: naïve AE versus sparse AE

- The AE objective function  $\arg \min_{\mathbf{W}, \mathbf{H}} \mathcal{D}(\mathbf{X} || \mathbf{W}^\dagger \mathbf{H}) + \lambda g(\mathbf{H})$ , where  $\mathbf{H} = \sigma(\mathbf{W} \mathbf{X})$

- With the sparsity constraint

$$\hat{\alpha} = \frac{1}{T} \sum_{j=0}^{T-1} \mathbf{H}_{i,j} \quad \text{Average activation (logistic)}$$

$$g(\mathbf{H}_{i,:}) = KL(\alpha || \hat{\alpha}) = -\alpha \log \left( \frac{\alpha}{\hat{\alpha}} \right) - (1 - \alpha) \log \left( \frac{1 - \alpha}{1 - \hat{\alpha}} \right)$$

$$g(\mathbf{H}) = \sum_i g(\mathbf{H}_{i,:}) \quad \text{Desired probability of activations}$$

- For example,  $x = [1, 0, 1, 0, 0, 0, 0, 0, 0, 0]$      $x = [0.8, 0, 0.9, 0, 0.2, 0, 0.1, 0, 0, 0]$

$$\hat{\alpha} = 0.2$$

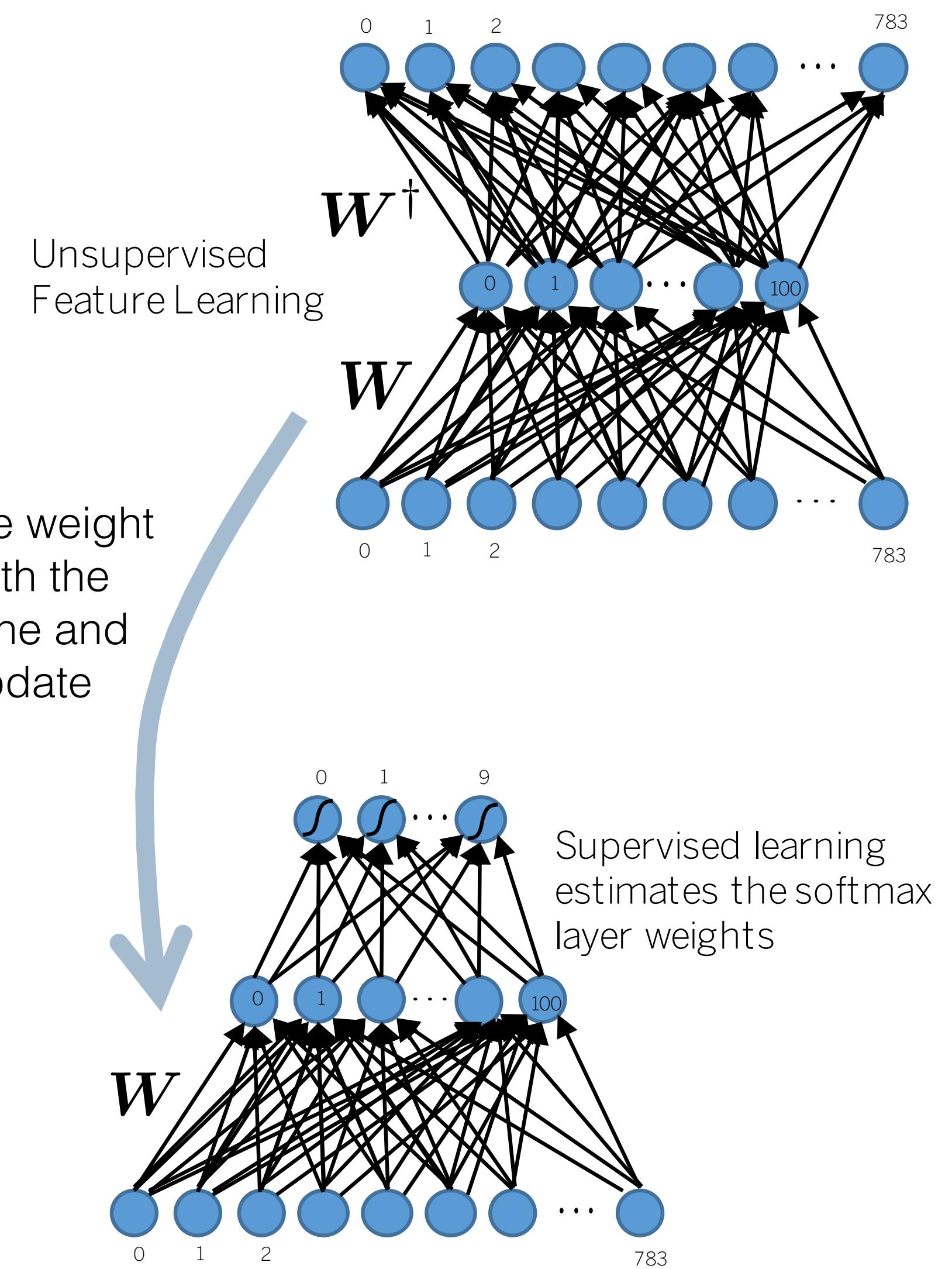
$$\hat{\alpha} = 0.2$$

$$\text{If } \alpha = 0.2 \quad g(x) = 0$$

- A less sparse one

$$x = [0.8, 1, 0.9, 1, 0.2, 0, 0.1, 0, 0, 0] \quad g(x) = 0.092$$

$$\hat{\alpha} = 0.4$$



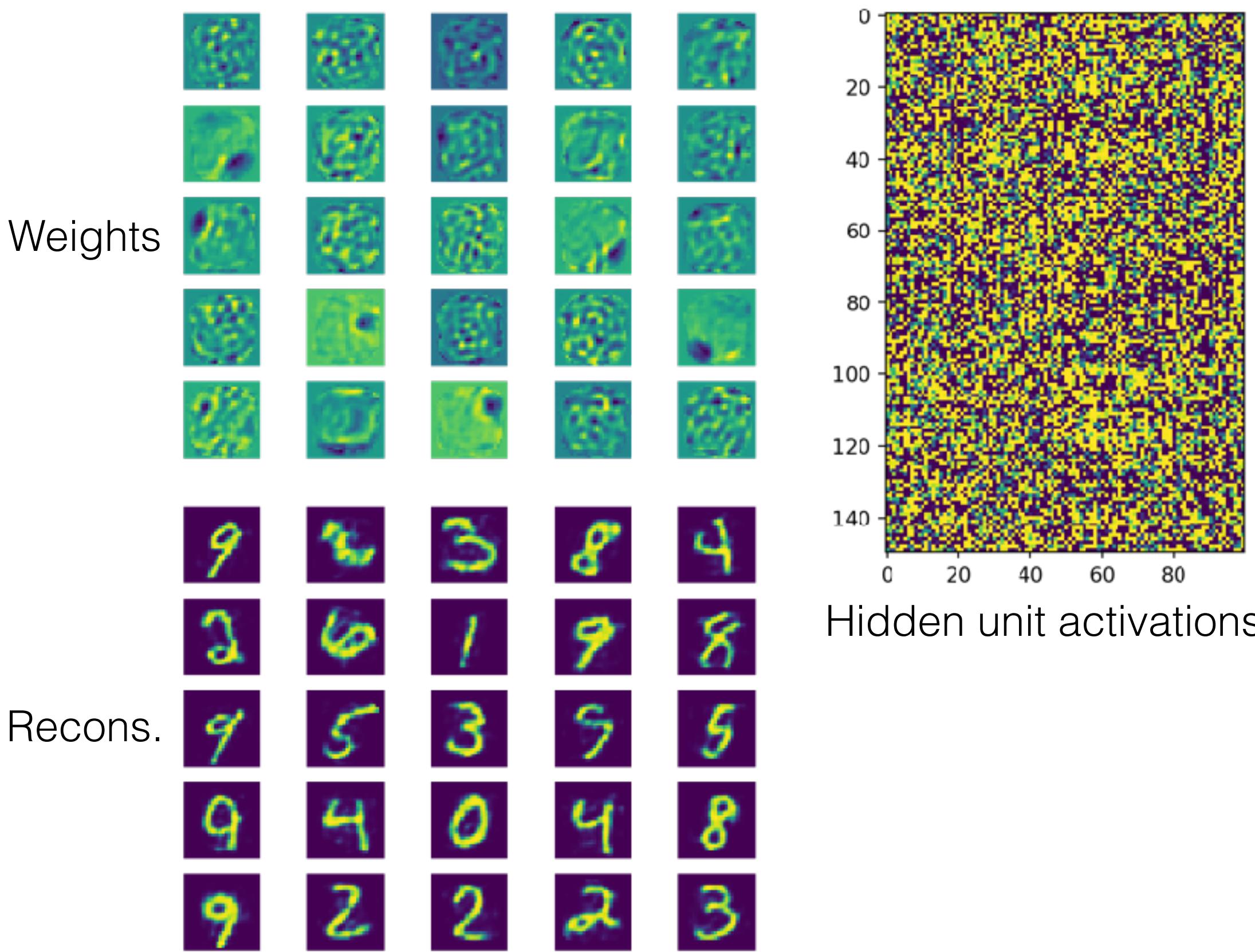
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

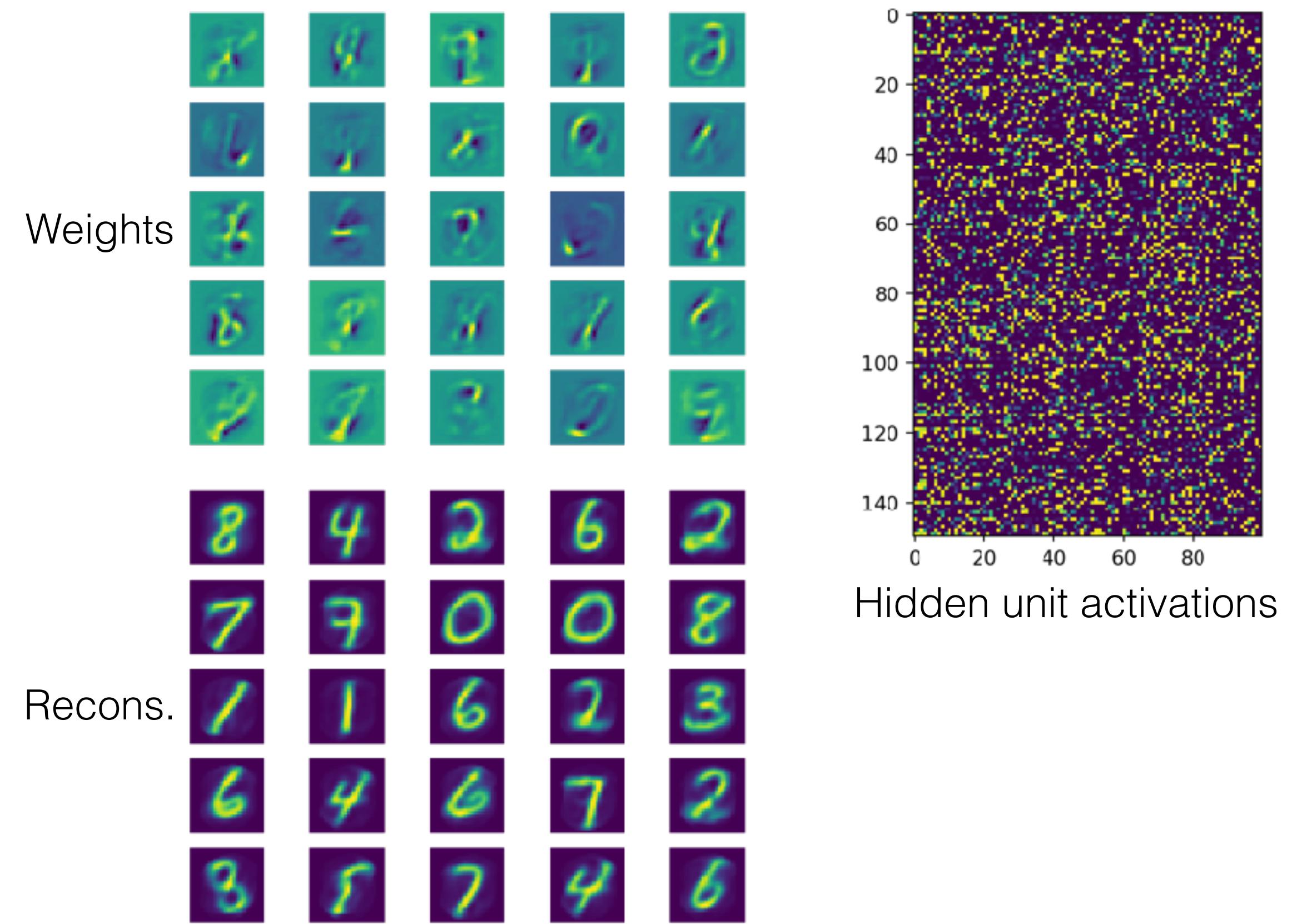
# Feature Extraction

## -Sparse AE

- Case A: Naïve AE
  - No sparsity constraint



- Case B: Sparse AE
  - $\lambda=0.1$ ,  $\alpha=0.2$



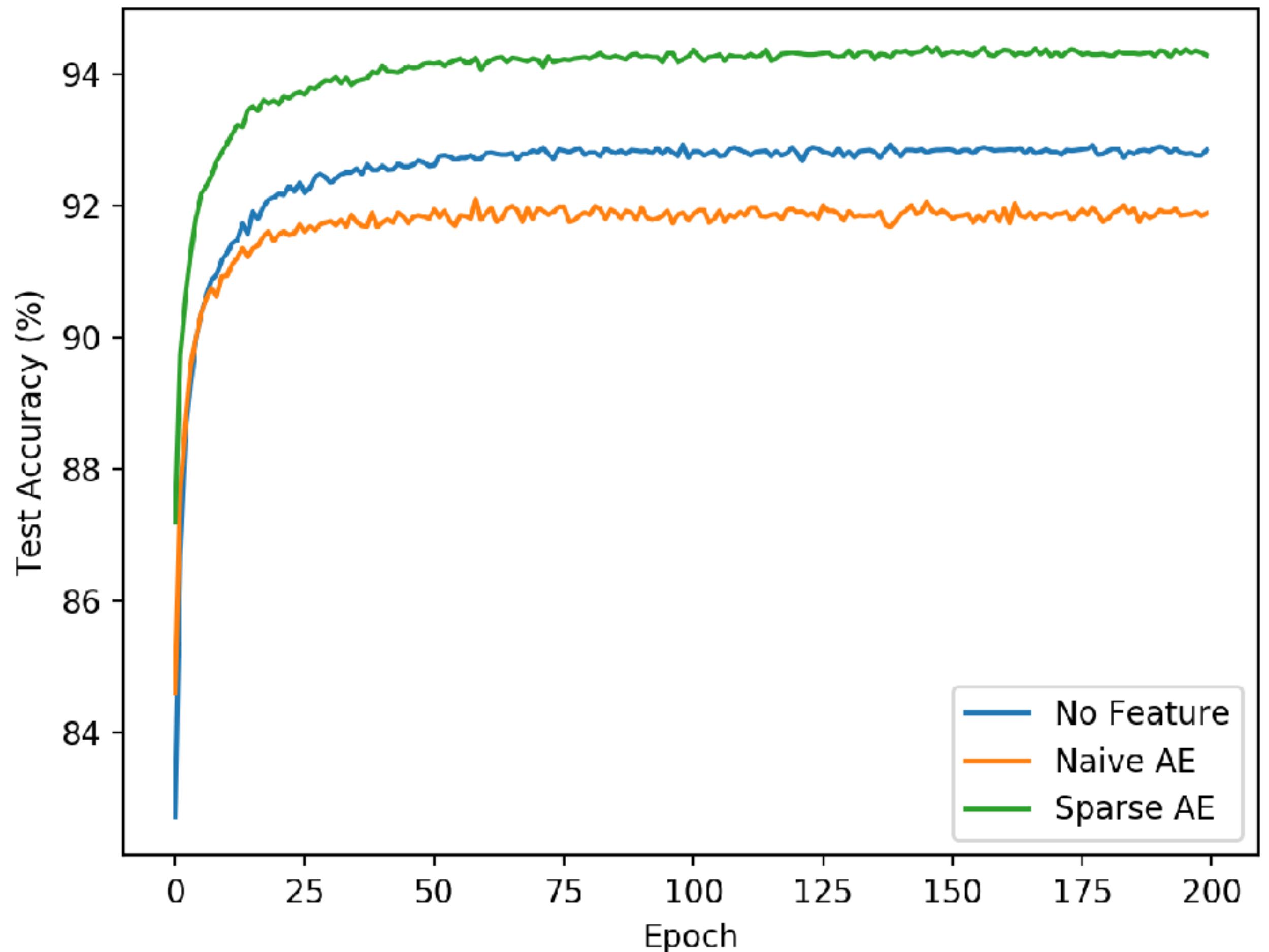
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Feature Extraction + Softmax

- Any better than the raw pixels?

- A direct softmax classification on the raw pixel intensities (blue line)
  - Works pretty well
  - No features
- If your features are NOT good enough (naïve AE)
  - They degrade the performance
- If your features are good enough (sparse AE)
  - They improve the performance

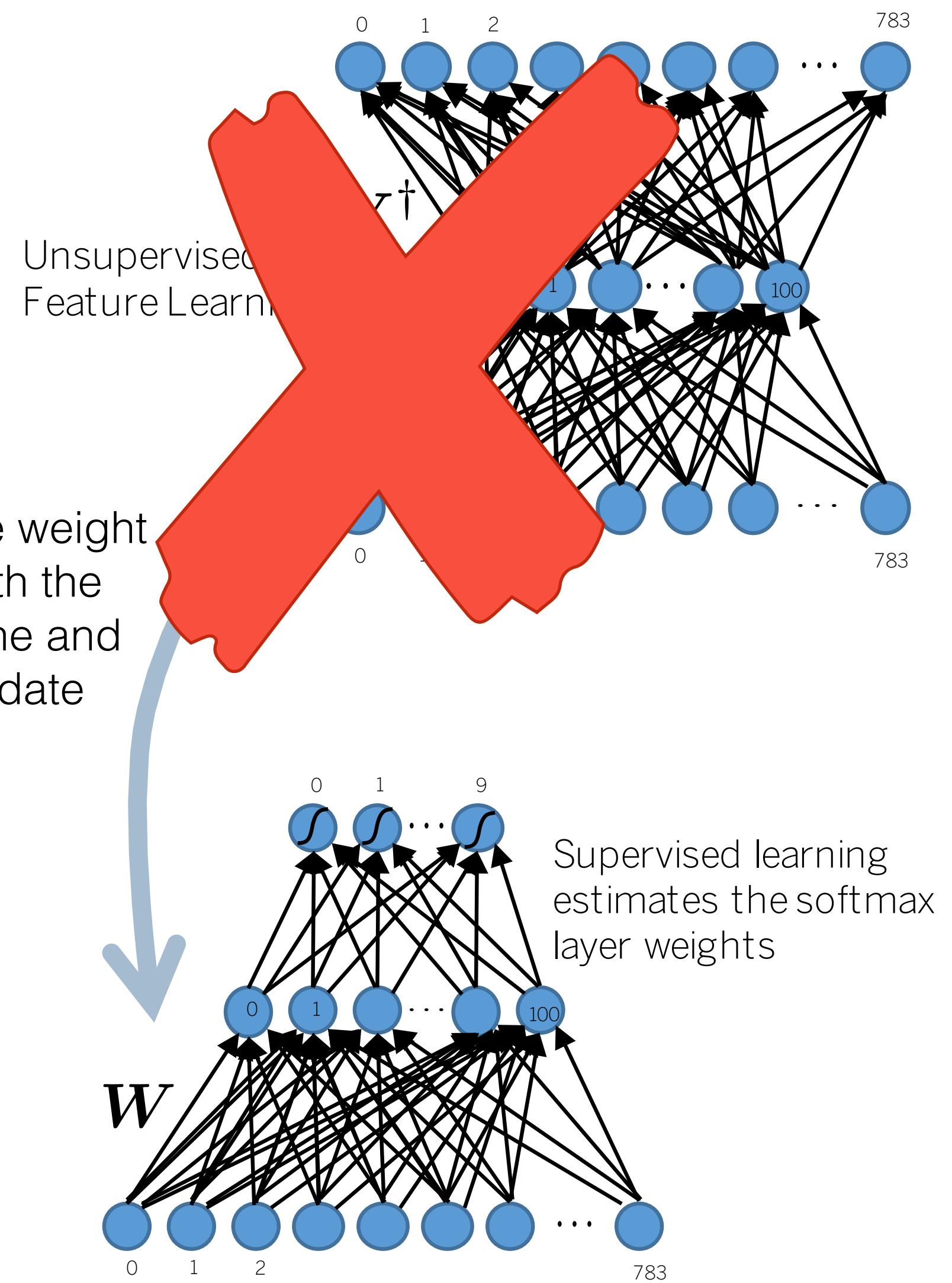


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

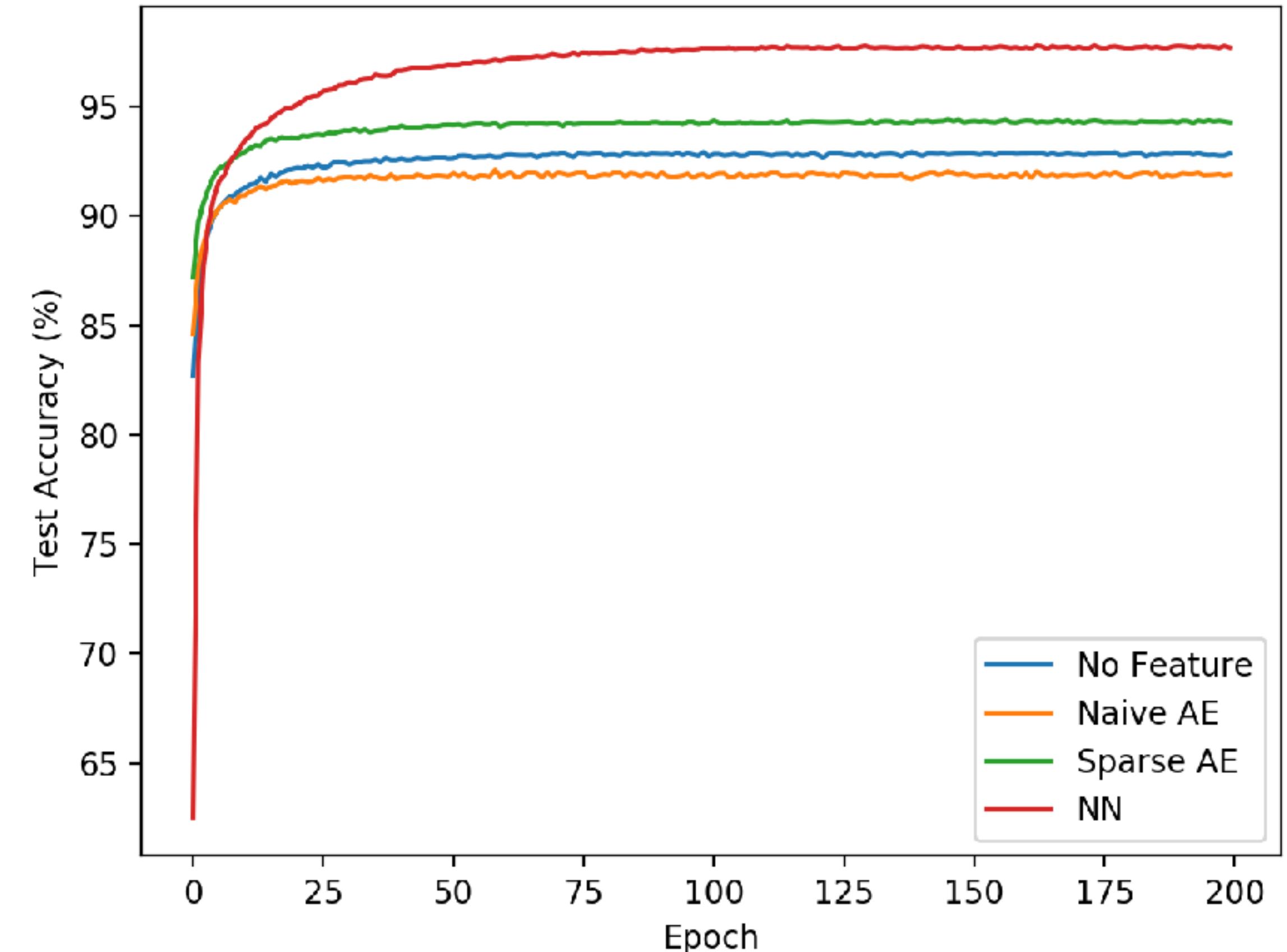
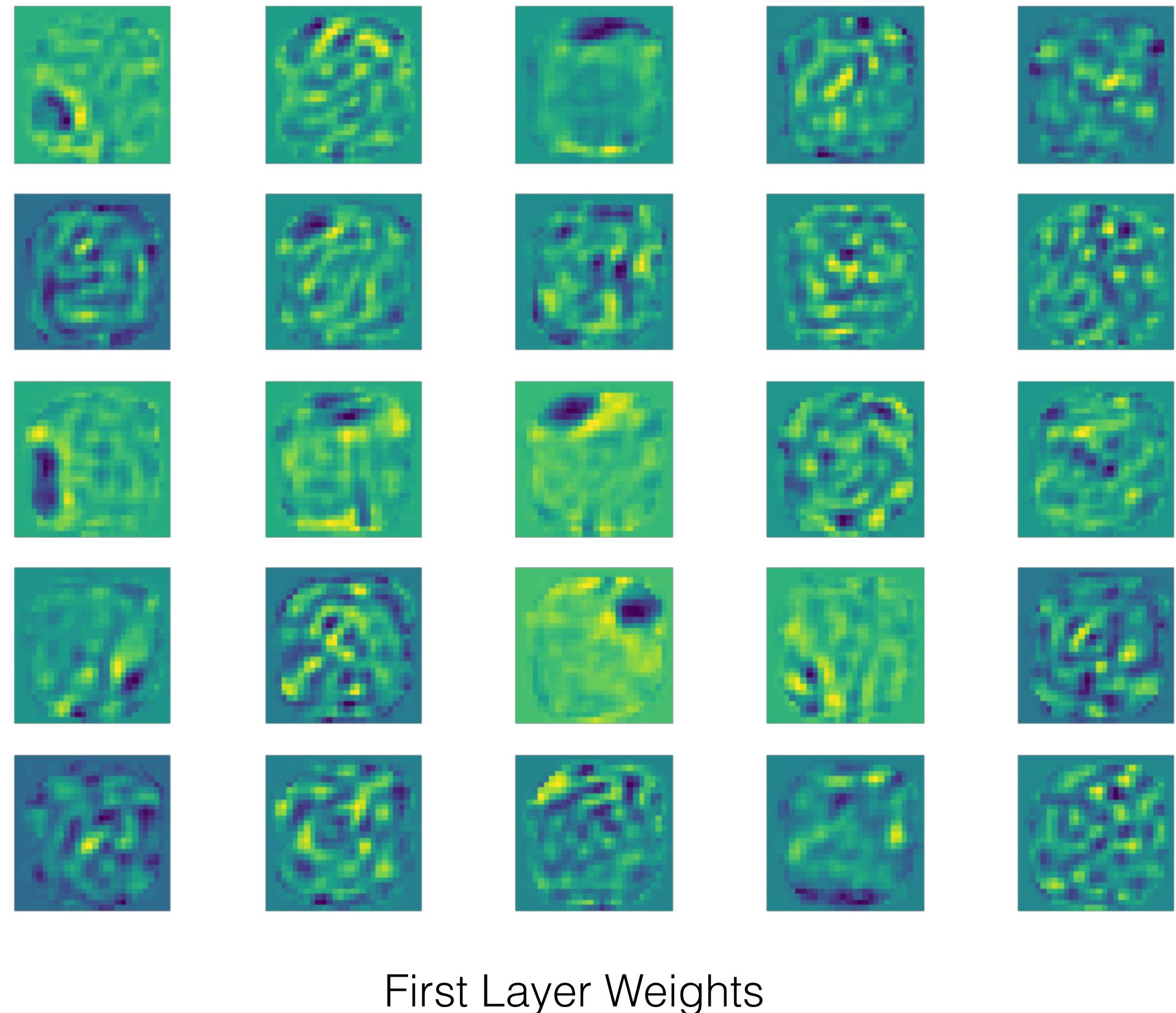
# Feature Extraction + Softmax?

- Big question: why bother learning the features?
  - Feature learning suffers from too many different decisions to make
    - e.g. regularization, constraints, ranks, error functions, activation functions, etc.
    - Could be unpredictable until you go ahead and test them out
  - Supervised learning has a lot of options, too
    - e.g. choice of the kernel functions, generative models, hyperparameters, etc.
- We can learn a neural network that has one hidden layer
  - Randomly initialize both weight matrices in the two layers
  - Estimate them using backpropagation
  - Does it work better than the previous methods with feature engineering?
    - See the graph in the next page



# Shallow Neural Networks

- No feature learning, but a hidden layer
  - Feature learning is absorbed in the NN learning process



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Backpropagation (Revisited)

-Now that we know the first layer weights should be updated, too

○  $i=0$

- Initialize parameters with small random numbers
- Calculate the prediction using all training samples (i.e. input and target pairs)

$$\mathbf{Z}^{(1)} \leftarrow \mathbf{W}^{(1)} \mathbf{X}^{(1)} + \mathbf{b}^{(1)}$$

$$\mathbf{X}^{(2)} \leftarrow \sigma(\mathbf{Z}^{(1)})$$

$$\mathbf{Z}^{(2)} \leftarrow \mathbf{W}^{(2)} \mathbf{X}^{(2)} + \mathbf{b}^{(2)}$$

$$\hat{\mathbf{Y}} \leftarrow g(\mathbf{Z}^{(2)})$$

- Calculate the error (cost)

$$\mathcal{E} = \sum_t \mathcal{D}(\mathbf{Y}_{:,t} \| \hat{\mathbf{Y}}_{:,t})$$

- Update the parameters

$$\mathbf{W}^{(2)} \leftarrow \mathbf{W}^{(2)} - \rho \frac{\partial \mathcal{E}}{\partial \mathbf{W}^{(2)}}$$

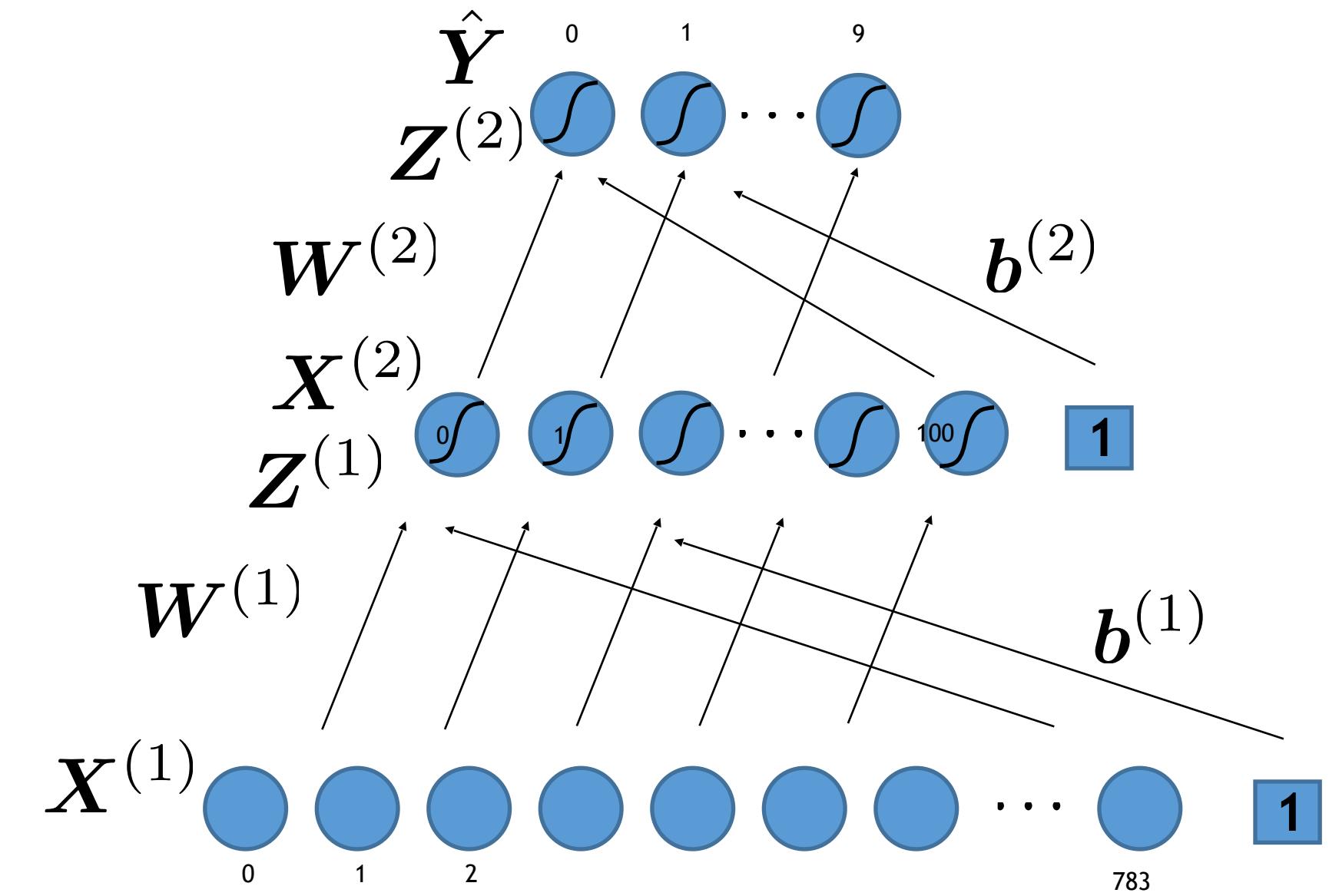
$$\mathbf{W}^{(1)} \leftarrow \mathbf{W}^{(1)} - \rho \frac{\partial \mathcal{E}}{\partial \mathbf{W}^{(1)}}$$

$$\mathbf{b}^{(2)} \leftarrow \mathbf{b}^{(2)} - \rho \frac{\partial \mathcal{E}}{\partial \mathbf{b}^{(2)}}$$

$$\mathbf{b}^{(1)} \leftarrow \mathbf{b}^{(1)} - \rho \frac{\partial \mathcal{E}}{\partial \mathbf{b}^{(1)}}$$

○  $i>0$

- Repeat 2-4



# Backpropagation (Revisited)

- Now that we know the first layer weights should be updated, too

- For the last layer

$$\frac{\partial \mathcal{E}}{\partial \mathbf{W}^{(2)}} = \frac{\partial \mathcal{E}}{\partial \hat{\mathbf{Y}}} \frac{\partial \hat{\mathbf{Y}}}{\partial \mathbf{Z}^{(2)}} \frac{\partial \mathbf{Z}^{(2)}}{\partial \mathbf{W}^{(2)}} = \underbrace{(\hat{\mathbf{Y}} - \mathbf{Y})}_{\text{BP Error } \Delta^{(2)}} \mathbf{X}^{(2)\top} \quad \frac{\partial \mathcal{E}}{\partial \mathbf{b}^{(2)}} = (\hat{\mathbf{Y}} - \mathbf{Y}) \mathbf{1}$$

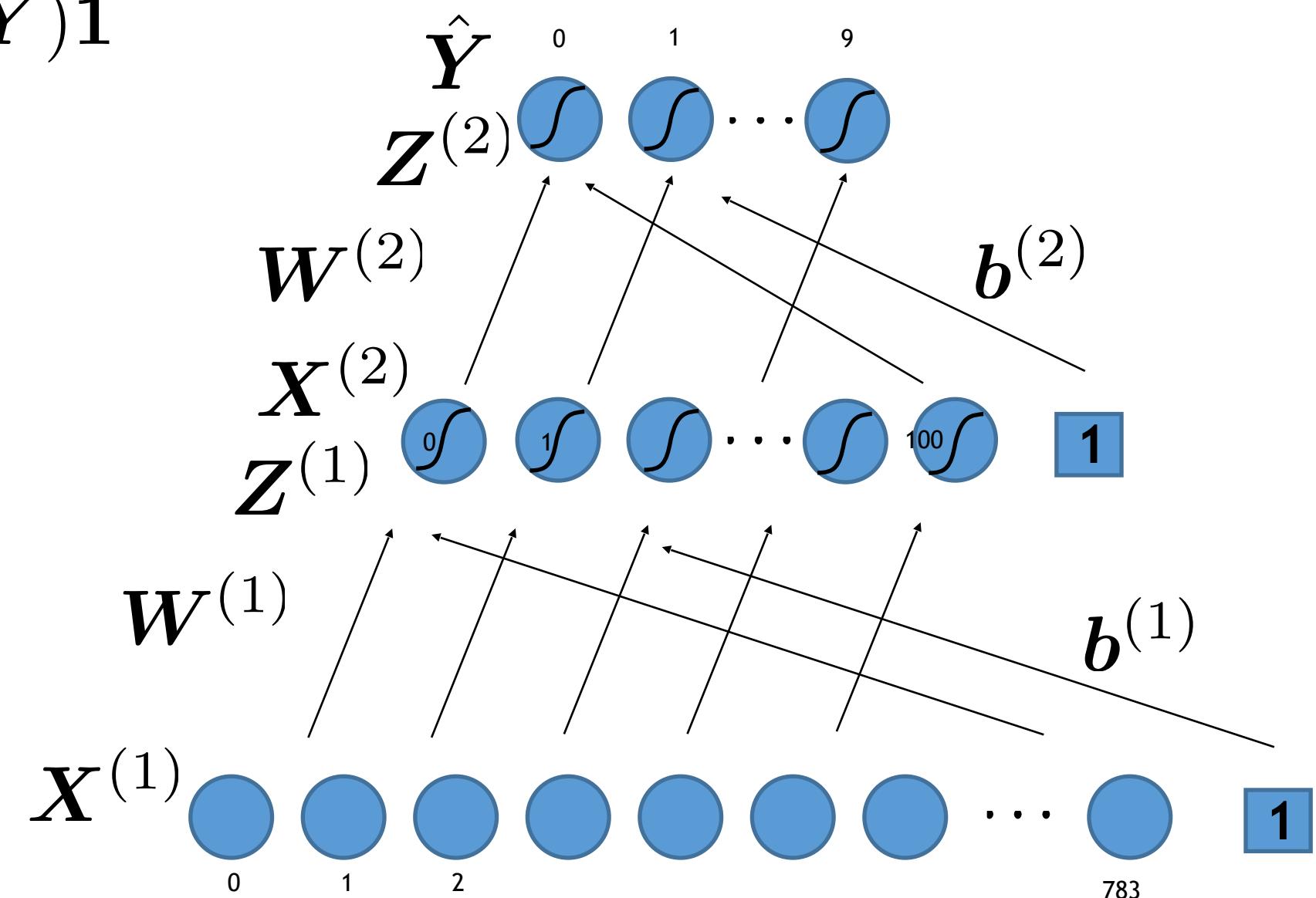
- Backpropagation Error  $\Delta^{(l)}$ :  
the derivative of the cost function w.r.t. the input to the ( $l$ )-th hidden layer unit

- For the first layer

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \mathbf{W}^{(1)}} &= \frac{\partial \mathcal{E}}{\partial \hat{\mathbf{Y}}} \frac{\partial \hat{\mathbf{Y}}}{\partial \mathbf{Z}^{(2)}} \frac{\partial \mathbf{Z}^{(2)}}{\partial \mathbf{X}^{(2)}} \frac{\partial \mathbf{X}^{(2)}}{\partial \mathbf{Z}^{(1)}} \frac{\partial \mathbf{Z}^{(1)}}{\partial \mathbf{W}^{(1)}} \\ &= \underbrace{(\mathbf{W}^{(2)\top} (\hat{\mathbf{Y}} - \mathbf{Y}))}_{\text{BP Error } \Delta^{(2)}} \odot \sigma'(\mathbf{Z}^{(1)}) \mathbf{X}^{(1)\top} \quad \frac{\partial \mathcal{E}}{\partial \mathbf{b}^{(1)}} = \Delta^{(1)} \mathbf{1} \end{aligned}$$

- The gradient for  $\mathbf{W}^{(l)}$  is [BP Error at ( $l$ )] X [Input from the previous layer]
- $\Delta^{(l)}$  : [BP Error at ( $l+1$ )] X [Weights at ( $l+1$ )] X [Derivative of the activation at ( $l$ )]

$$\begin{aligned} \mathbf{Z}^{(1)} &\leftarrow \mathbf{W}^{(1)} \mathbf{X}^{(1)} + \mathbf{b}^{(1)} \\ \mathbf{X}^{(2)} &\leftarrow \sigma(\mathbf{Z}^{(1)}) \\ \mathbf{Z}^{(2)} &\leftarrow \mathbf{W}^{(2)} \mathbf{X}^{(2)} + \mathbf{b}^{(2)} \\ \hat{\mathbf{Y}} &\leftarrow g(\mathbf{Z}^{(2)}) \end{aligned}$$

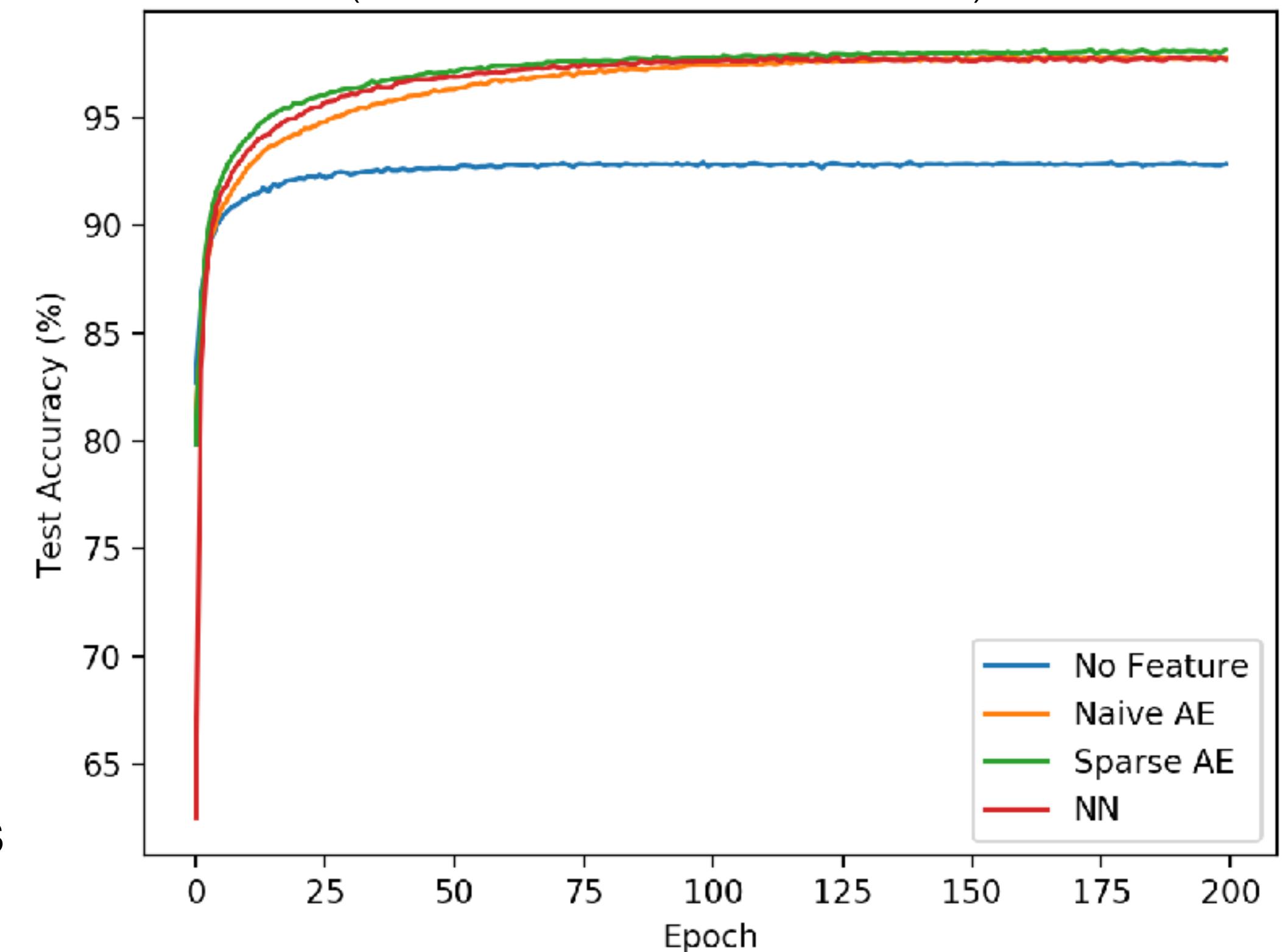


# A Step Towards Deep Neural Networks?

## -BP on DNNs?

- We noticed that a separate feature learning procedure is NOT super helpful for softmax classification
  - A direct BP down to the first layer is doable and works better
    - Why is it better?
    - The feature learning part (i.e. the first layer weights) is aware of the discrimination objective
  - It might NOT be true if the last layer is more complicated (e.g. SVM)
    - Maybe the features don't really have to be that great
- Another option
  - What if we initialize the first layer with the trained AE weights, and then update weights of both layers?
    - See the graph
  - A better initialization doesn't seem to contribute to the convergence (why?)
    - BP just gets down to the first layer
- Conclusion (so far): we don't need pretraining for shallow NNs
- What if we make the network deeper?

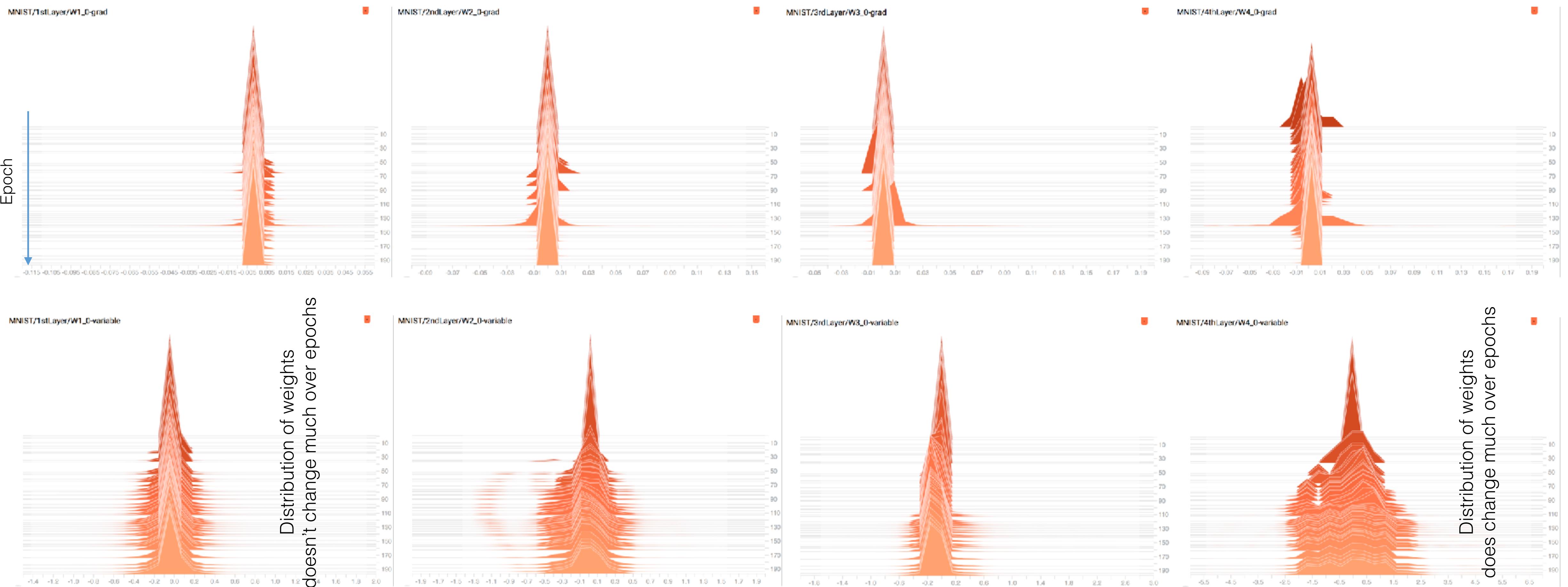
Full BP training of shallow net with one hidden layer  
(different initialization schemes)



# Vanishing Gradients

-BP doesn't get down to the first layer

- The vanishing gradient problem



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Vanishing Gradients

## -Why does it happen?

- It's due to the sigmoid function. What's wrong with them?
  - Their gradients are small when the function is saturated
- Suppose a DNN with L hidden layers

$$\mathbf{Z}^{(1)} \leftarrow \mathbf{W}^{(1)} \mathbf{X}^{(1)} + \mathbf{b}^{(1)}$$

$$\mathbf{Z}^{(2)} \leftarrow \mathbf{W}^{(2)} \mathbf{X}^{(2)} + \mathbf{b}^{(2)}$$

$$\cdots$$

$$\mathbf{Z}^{(L)} \leftarrow \mathbf{W}^{(L)} \mathbf{X}^{(L)} + \mathbf{b}^{(L)}$$

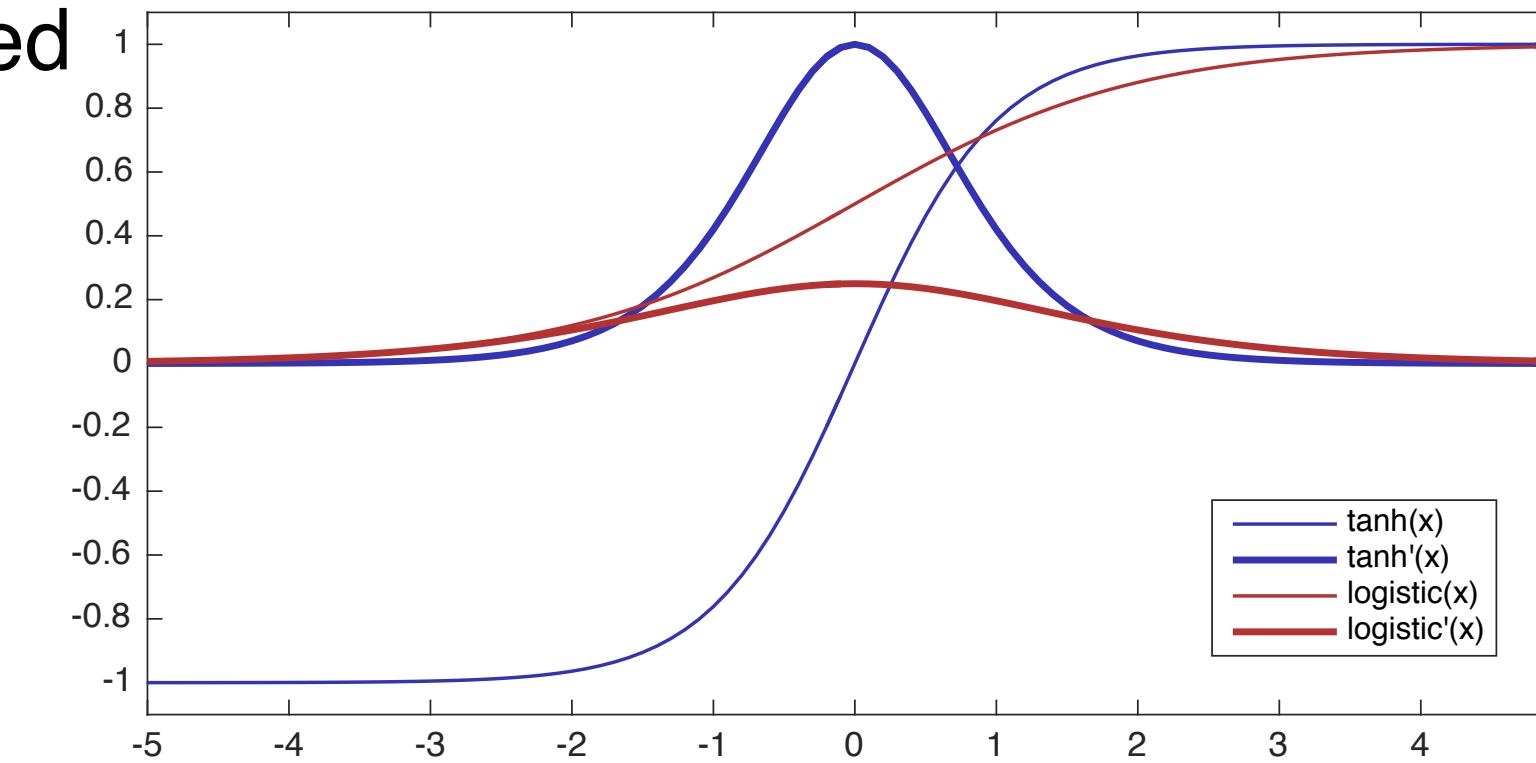
$$\mathbf{Z}^{(L+1)} \leftarrow \mathbf{W}^{(L+1)} \mathbf{X}^{(L+1)} + \mathbf{b}^{(L+1)}$$

$$\mathbf{X}^{(2)} \leftarrow \sigma(\mathbf{Z}^{(1)})$$

$$\mathbf{X}^{(3)} \leftarrow \sigma(\mathbf{Z}^{(2)})$$

$$\mathbf{X}^{(L+1)} \leftarrow \sigma(\mathbf{Z}^{(L)})$$

$$\hat{\mathbf{Y}} \leftarrow \text{softmax}(\mathbf{Z}^{(L+1)})$$



## Gradients per layer

$$\frac{\partial \mathcal{E}}{\partial \mathbf{W}^{(L+1)}} = \frac{\partial \mathcal{E}}{\partial \hat{\mathbf{Y}}} \frac{\partial \hat{\mathbf{Y}}}{\partial \mathbf{Z}^{(L+1)}} \frac{\partial \mathbf{Z}^{(L+1)}}{\partial \mathbf{W}^{(L+1)}} = (\hat{\mathbf{Y}} - \mathbf{Y}) \mathbf{X}^{(L+1)\top}$$

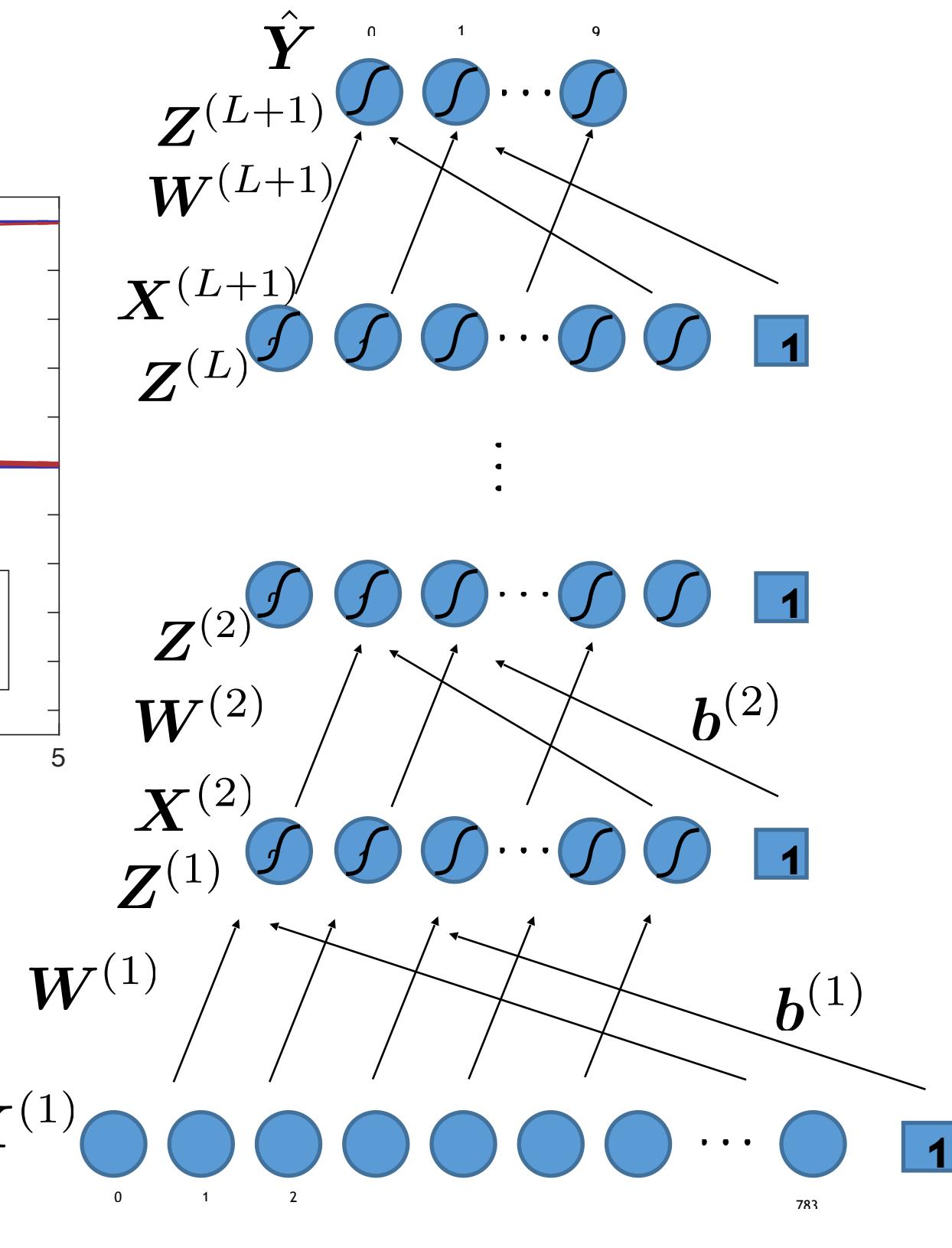
BP error  $\Delta^{(L+1)}$

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

$\Delta^{(L)}$

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

$\Delta^{(1)}$   
 $L$  speed bumps



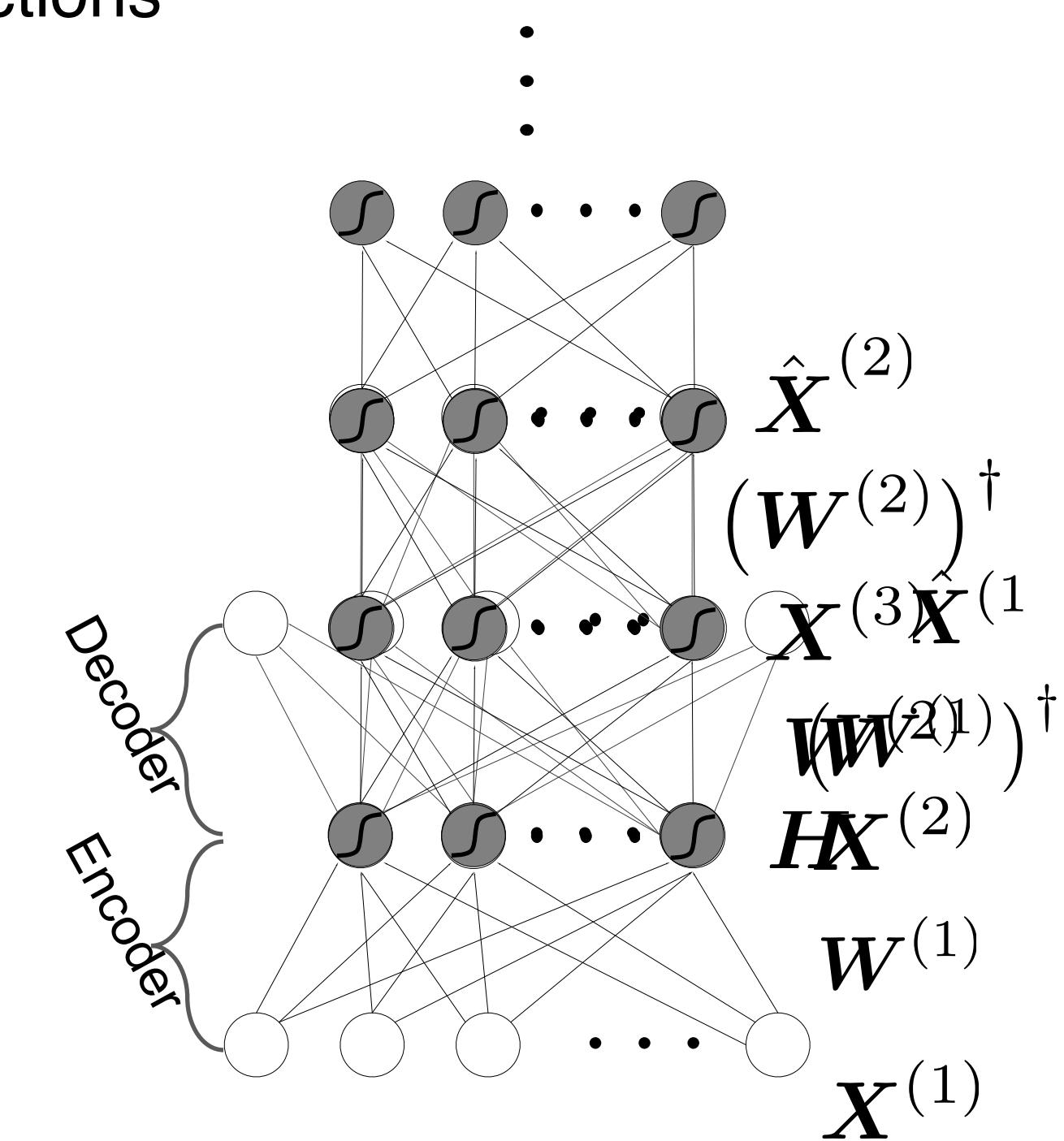
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Pretraining

## -Greedy layer-wise unsupervised learning

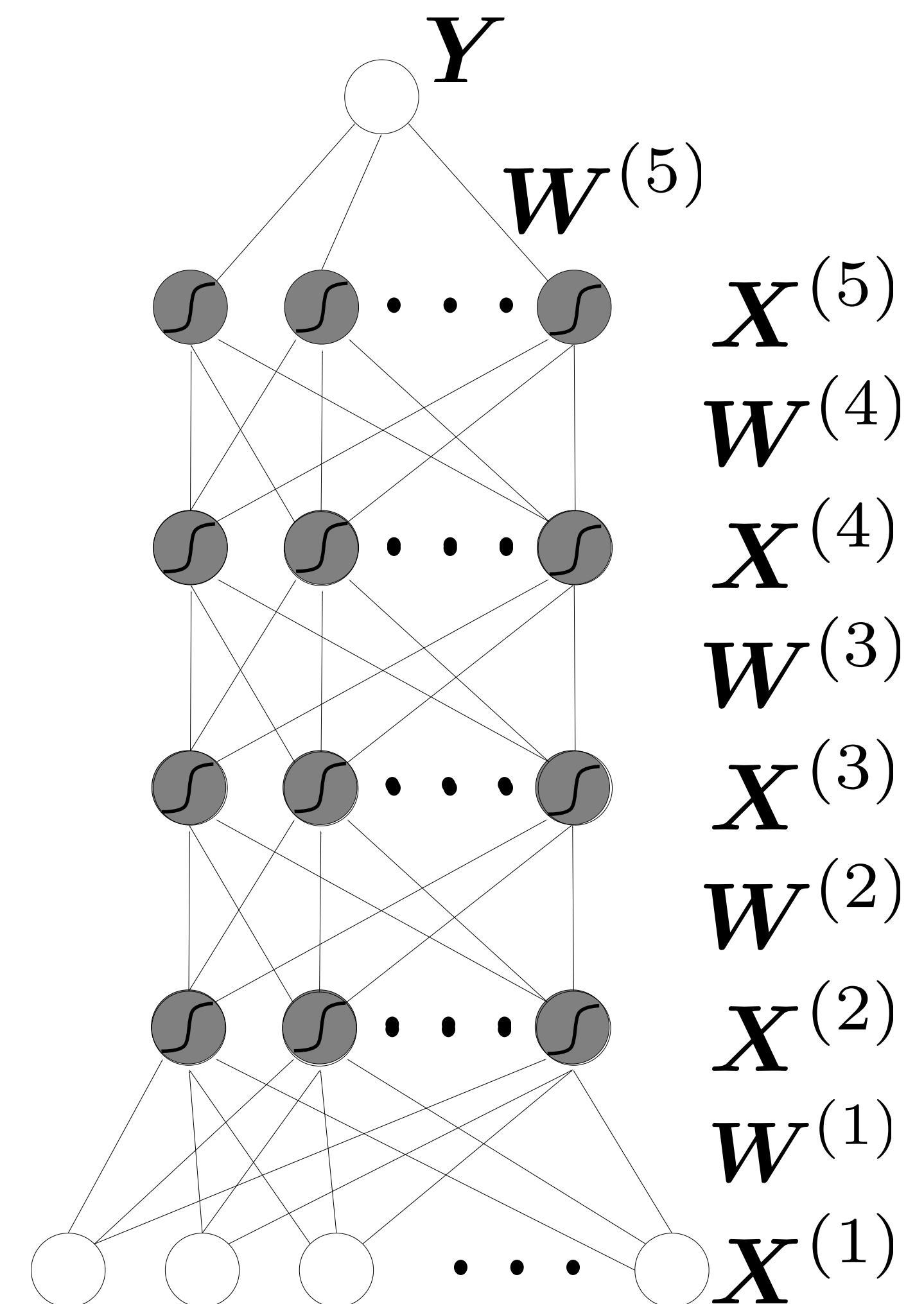
- A primitive method to train a DNN
  - The BP error vanishes in the lower layer due to the use of sigmoid functions
  - Pretraining learns the lower layer weights in a different way
- Stacked AE
  - For the 1st layer
$$\hat{\mathbf{X}}^{(1)} \approx \mathbf{W}^{(1)\dagger} \sigma(\mathbf{W}^{(1)} \mathbf{X}^{(1)})$$
  - For the 2nd layer
$$\mathbf{X}^{(2)} \leftarrow \sigma(\mathbf{W}^{(1)} \mathbf{X}^{(1)})$$
$$\hat{\mathbf{X}}^{(2)} \approx \mathbf{W}^{(2)\dagger} \sigma(\mathbf{W}^{(2)} \mathbf{X}^{(2)})$$
  - For the ( $l$ )-th layer
$$\mathbf{X}^{(l)} \leftarrow \sigma(\mathbf{W}^{(l-1)} \mathbf{X}^{(l-1)})$$
$$\hat{\mathbf{X}}^{(l)} \approx \mathbf{W}^{(l)\dagger} \sigma(\mathbf{W}^{(l)} \mathbf{X}^{(l)})$$
- What kinds of autoencoders are there?
  - Hinton's first pretraining method used Restricted Boltzman Machine
  - Andrew Ng's group showed sparse AEs work well
  - Bengio's group suggested denoising AEs



# Pretraining

## -Greedy layer-wise unsupervised learning

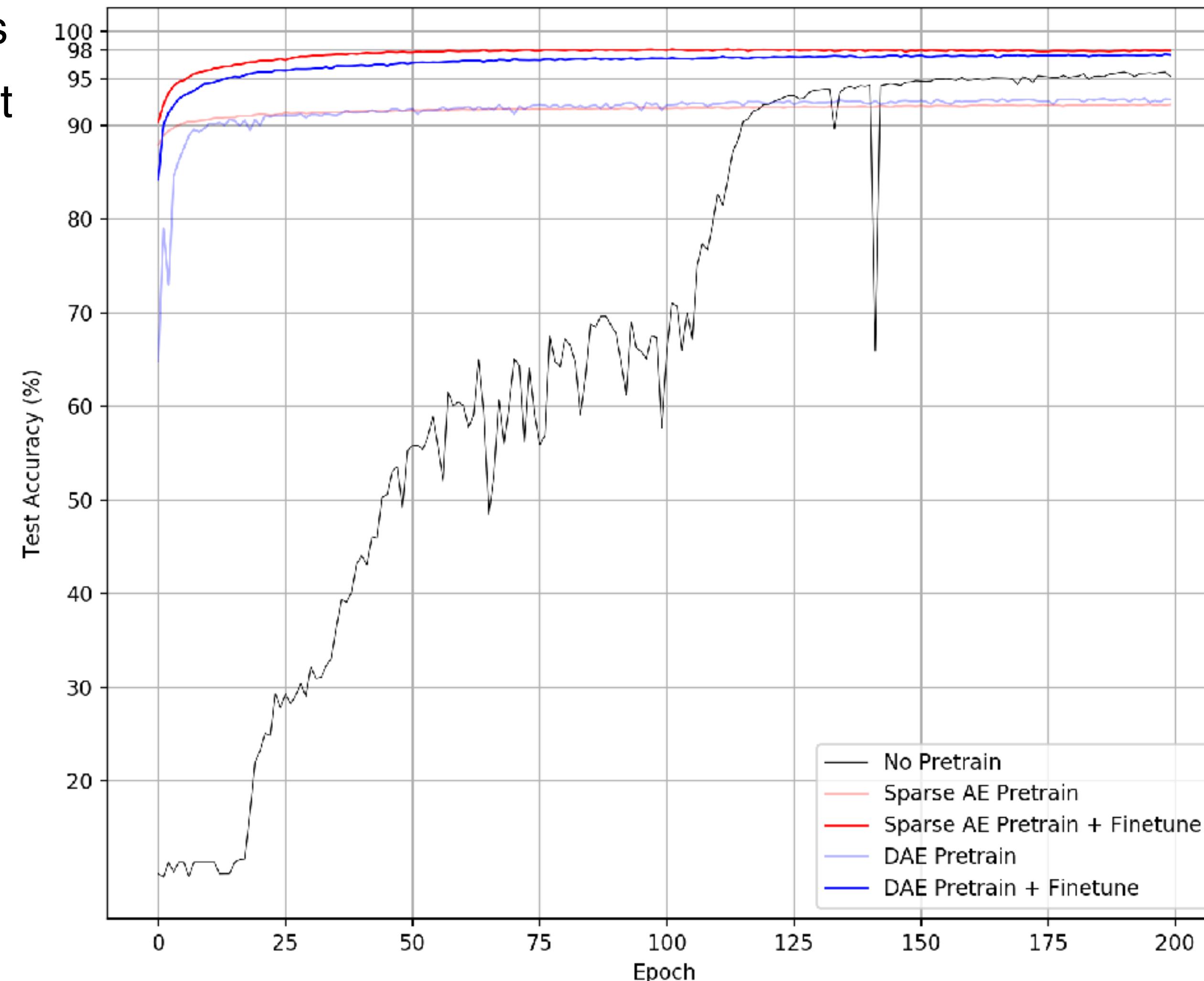
- Now that you know some better weights, the DNN can be initialized with the pretrained weights
  - Any problems with this?
- Each layer-wise AE creates approximation error
  - which propagates up to the last layer
  - Last layer weights in the stacked AE are not in good quality
  - What can we do?
- Finetuning
  - This procedure is not different from a regular BP
  - It might not be very effective to update the lower layers
  - But can fix the last layer weights
  - A good solution for the stacked AE
- The pretraining-and-finetuning pipeline was the first deep learning method
- But pretraining is known to be difficult to train



# Pretraining and Finetuning

- Does this work well?

- 100X3 Logistic Hidden Units
- Stochastic Gradient Descent



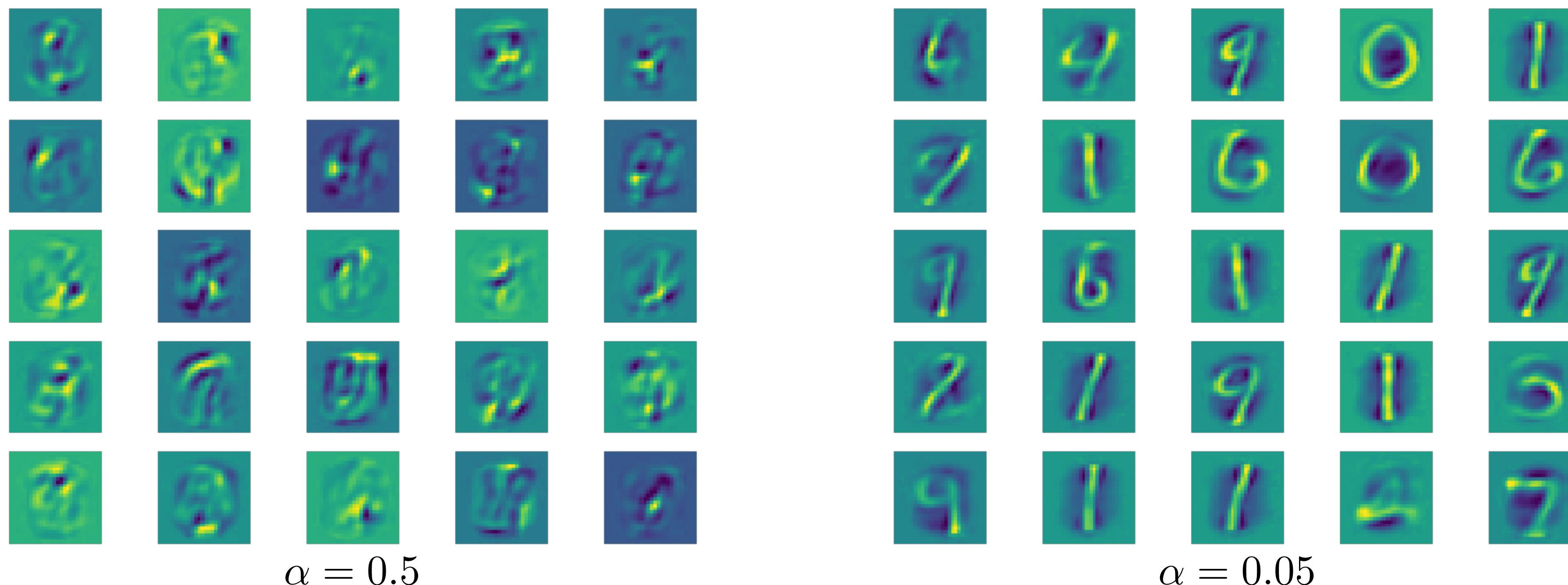
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Pretraining and Finetuning

- Does this work well?

- The thing is, training such autoencoders is tricky in a lot of senses
  - There are some criteria known to work (e.g. sparsity), but it's difficult to know until we try them out
  - It's unsupervised learning, so there's a chance of overfitting
  - Some hyperparameters to tune as usual
    - e.g. What would be the right amount of sparsity?



# Pretraining and Finetuning

- Does this work well?

- More importantly, the higher layers are hard to visualize to figure out the best way to train
- It turned out I'm not the only one who doesn't like this approach
- What would be the other ways to train DNNs?
  - Choosing the activation function wisely
  - Securing the pathway of BP error
  - Better optimization methods help, too



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Reading

- Hinton, Geoffrey E., Simon Osindero, and Yee-Whye Teh. "A fast learning algorithm for deep belief nets." *Neural computation* 18.7 (2006): 1527-1554.
- Lee, Honglak, Chaitanya Ekanadham, and Andrew Y. Ng. "Sparse deep belief net model for visual area V2." *Advances in neural information processing systems*. 2008.
- Vincent, Pascal, et al. "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion." *Journal of Machine Learning Research* 11.Dec (2010): 3371-3408.



**Thank You!**