

# CS5841/EE5841 Machine Learning

## Lecture 11: Neural Networks, extra details

Evan Lucas



Michigan Tech

# Putting it all together

- Neural networks are multiple layers of linear classifiers
- We train them through optimization with gradient descent methods
- Implementation steps and finer details

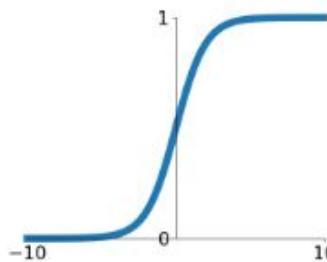


Michigan Tech

# Activation functions

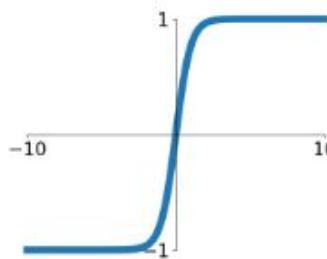
**Sigmoid**

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



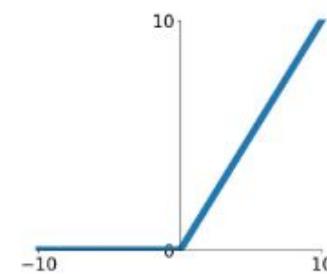
**tanh**

$$\tanh(x)$$



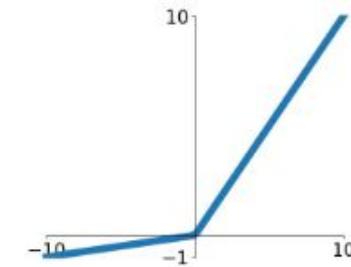
**ReLU**

$$\max(0, x)$$



**Leaky ReLU**

$$\max(0.1x, x)$$

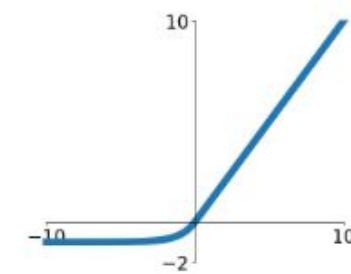


**Maxout**

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

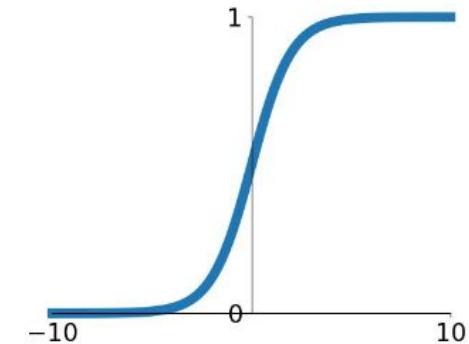
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Michigan Tech

# Sigmoid

- Sigmoid is historically popular because of Logistic Regression lineage AND easily interpreted because in range  $[0,1]$ 
  - What issues could this cause?

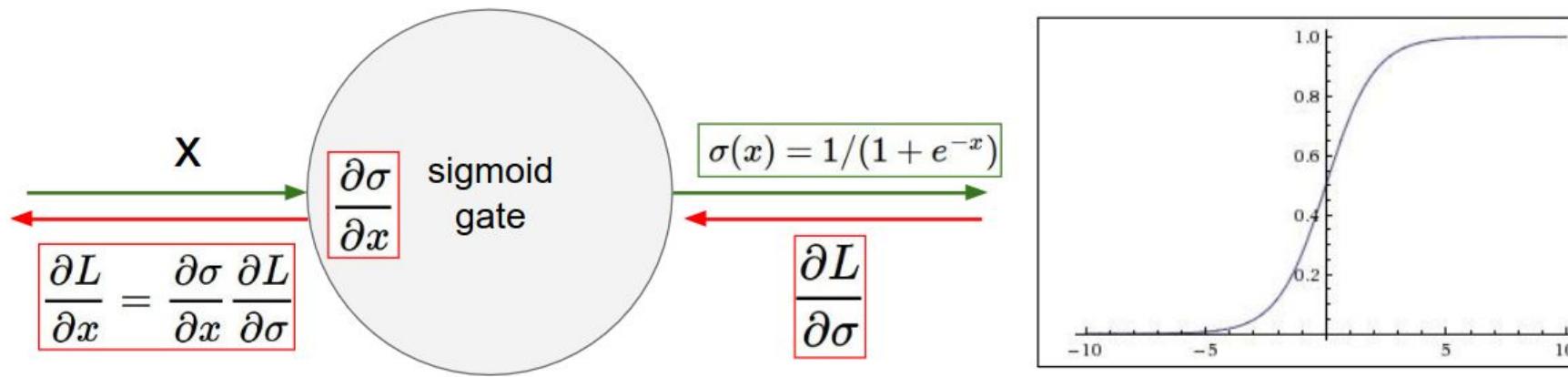


**Sigmoid**



**Michigan Tech**

# Issue 1: Saturation



What happens when  $x = -10$ ?  
What happens when  $x = 0$ ?  
What happens when  $x = 10$ ?



Michigan Tech

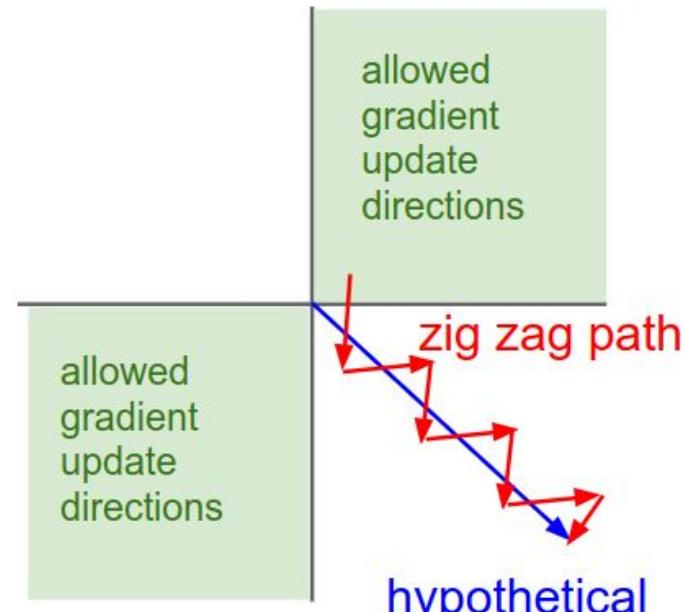
# Issue 2: Not zero centered

Consider what happens when the input to a neuron is always positive...

This is somewhat of an edge case and using smaller batch sizes helps prevent it

$$f \left( \sum_i w_i x_i + b \right)$$

What can we say about the gradients on  $w$ ?  
Always all positive or all negative :(

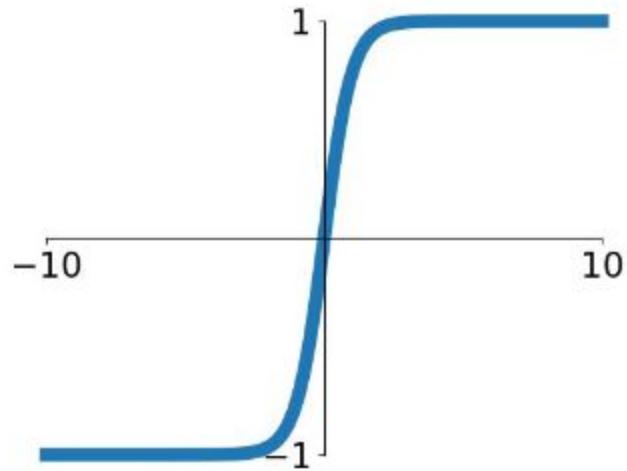


# Issue 3: `exp()` is computationally expensive



Michigan Tech

# **tanh()**



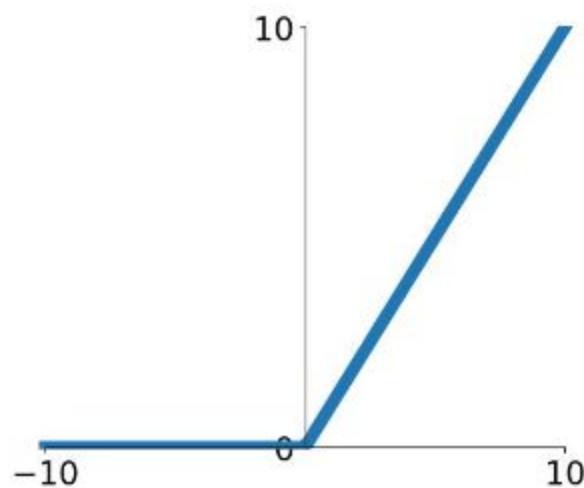
**tanh(x)**

- Squashes numbers to range [-1,1]
- zero centered (nice)
- still kills gradients when saturated :(



**Michigan Tech**

# ReLU (ie. $\max(0,x)$ )



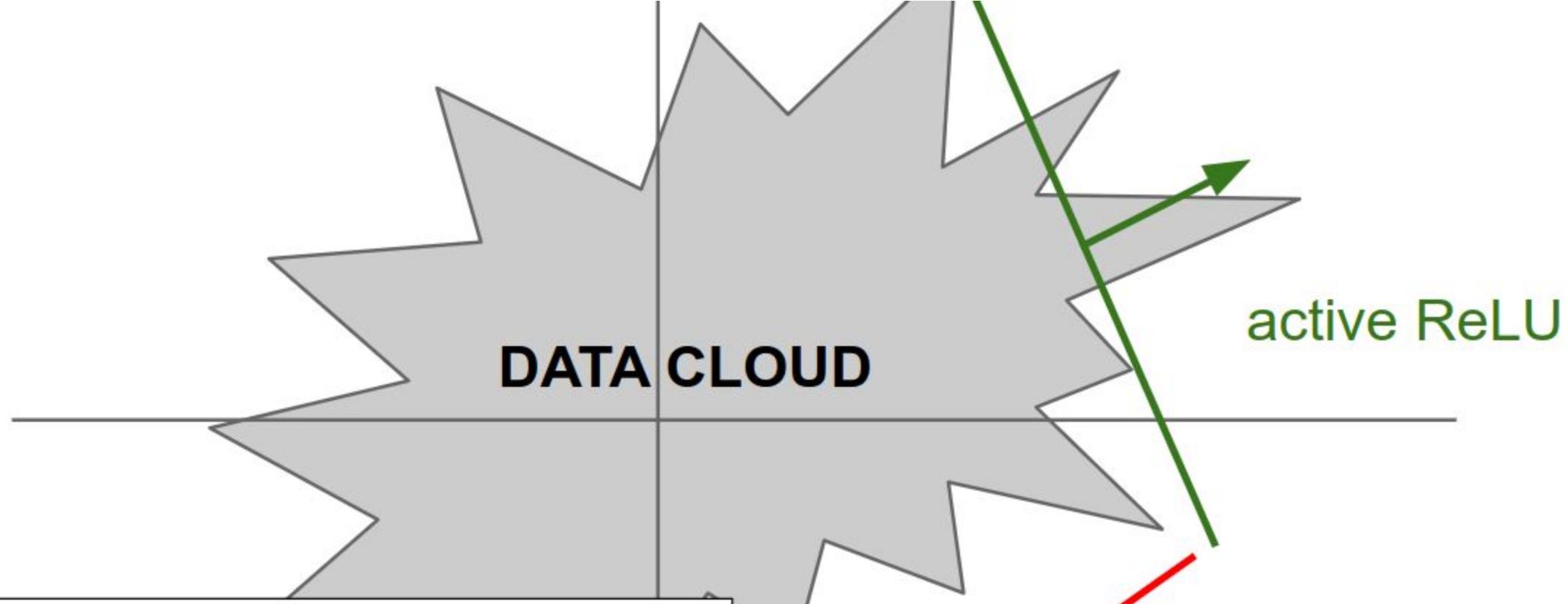
- Does not saturate (in +region)
  - Very computationally efficient
  - Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- 
- Not zero-centered output
  - An annoyance:

**ReLU**  
(Rectified Linear Unit)

hint: what is the gradient when  $x < 0$ ?

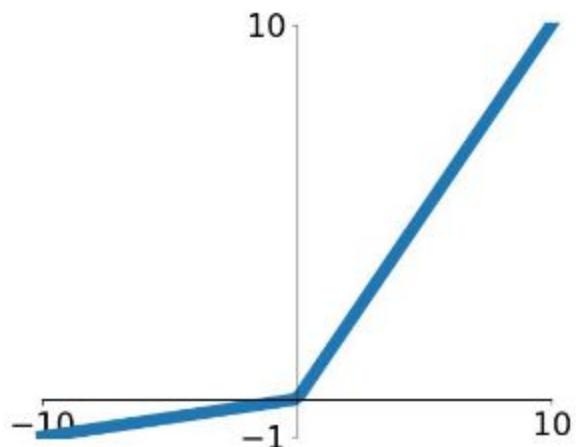


Michigan Tech



# Activation Functions

[Mass et al., 2013]  
[He et al., 2015]



- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

## Leaky ReLU

$$f(x) = \max(0.01x, x)$$



Michigan Tech

# Others

<b>ReLU</b> $\max(0, x)$	<b>GELU</b> $\frac{x}{2} \left( 1 + \tanh \left( \sqrt{\frac{2}{\pi}} (x + ax^3) \right) \right)$	<b>PRelu</b> $\max(0, x)$
<b>ELU</b> $\begin{cases} x & \text{if } x > 0 \\ \alpha(x \exp x - 1) & \text{if } x < 0 \end{cases}$	<b>Swish</b> $\frac{x}{1 + \exp -x}$	<b>SELU</b> $\alpha(\max(0, x) + \min(0, \beta(\exp x - 1)))$
<b>SoftPlus</b> $\frac{1}{\beta} \log(1 + \exp(\beta x))$	<b>Mish</b> $x \tanh \left( \frac{1}{\beta} \log(1 + \exp(\beta x)) \right)$	<b>RReLU</b> $\begin{cases} x & \text{if } x \geq 0 \\ ax & \text{if } x < 0 \text{ with } a \sim \mathcal{R}(l, u) \end{cases}$



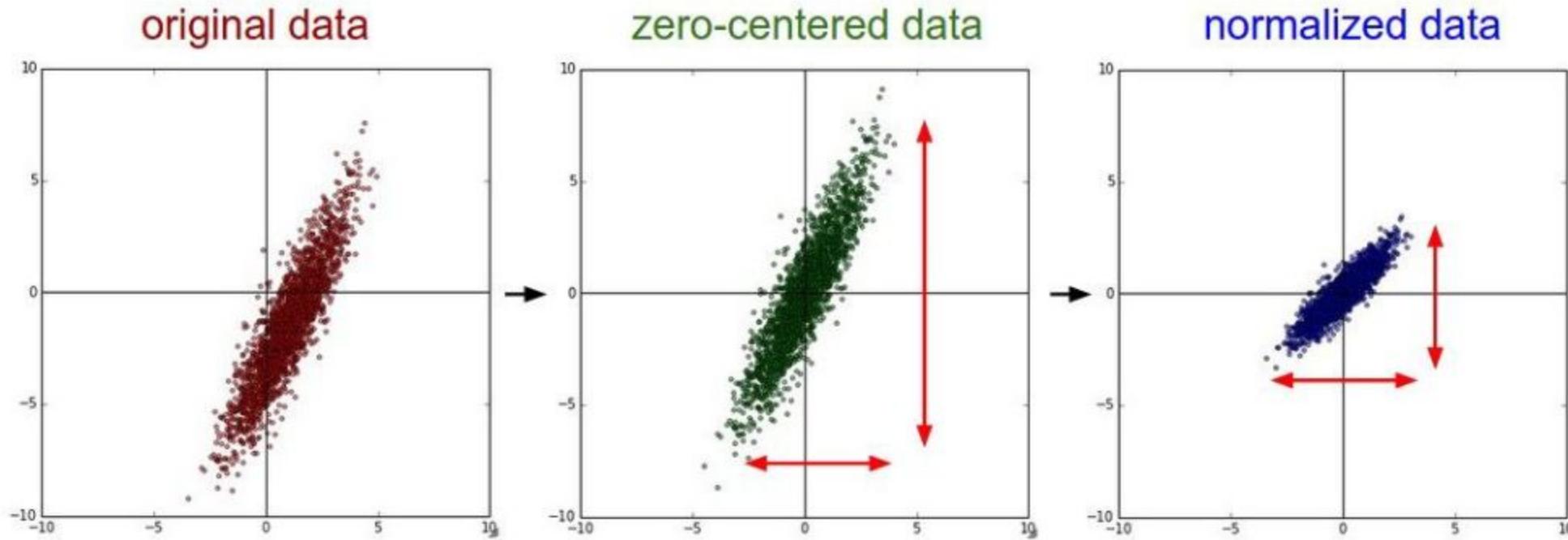
# General best practices

- ReLU usually works
- Consider Leaky ReLU/Maxout/ELU
- Play with other ones as desired
- Don't use sigmoid



Michigan Tech

# Data preprocessing



```
X -= np.mean(X, axis = 0)
```

```
X /= np.std(X, axis = 0)
```

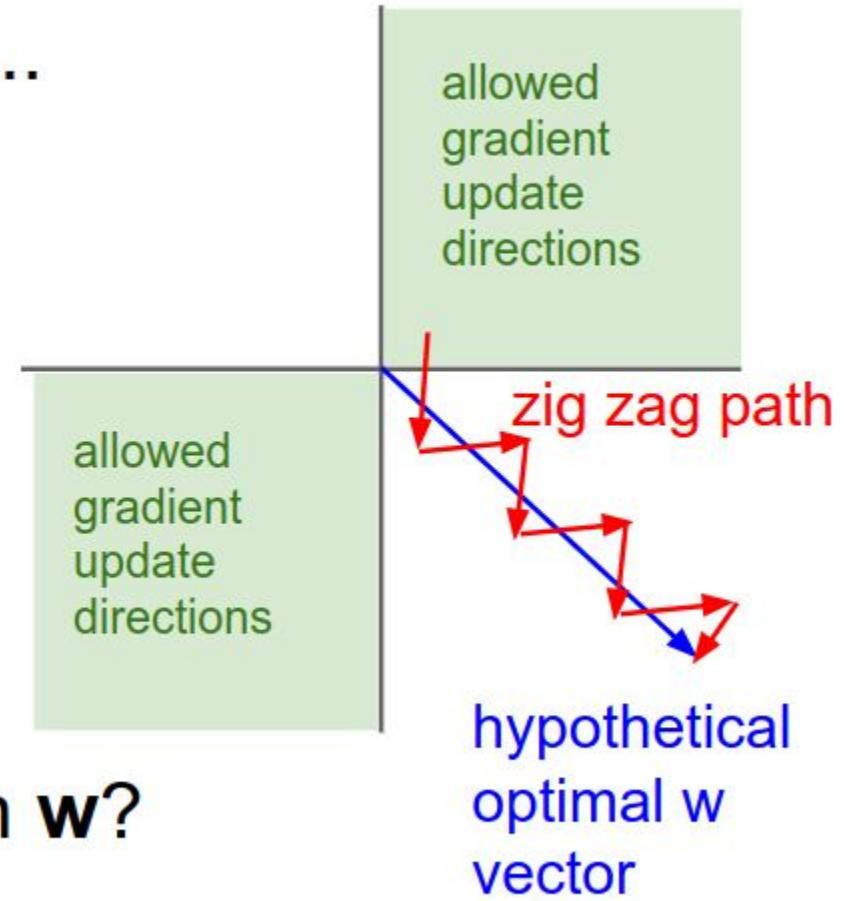
(Assume  $X$  [NxD] is data matrix,  
each example in a row)



Michigan Tech

Remember: Consider what happens when the input to a neuron is always positive...

$$f \left( \sum_i w_i x_i + b \right)$$

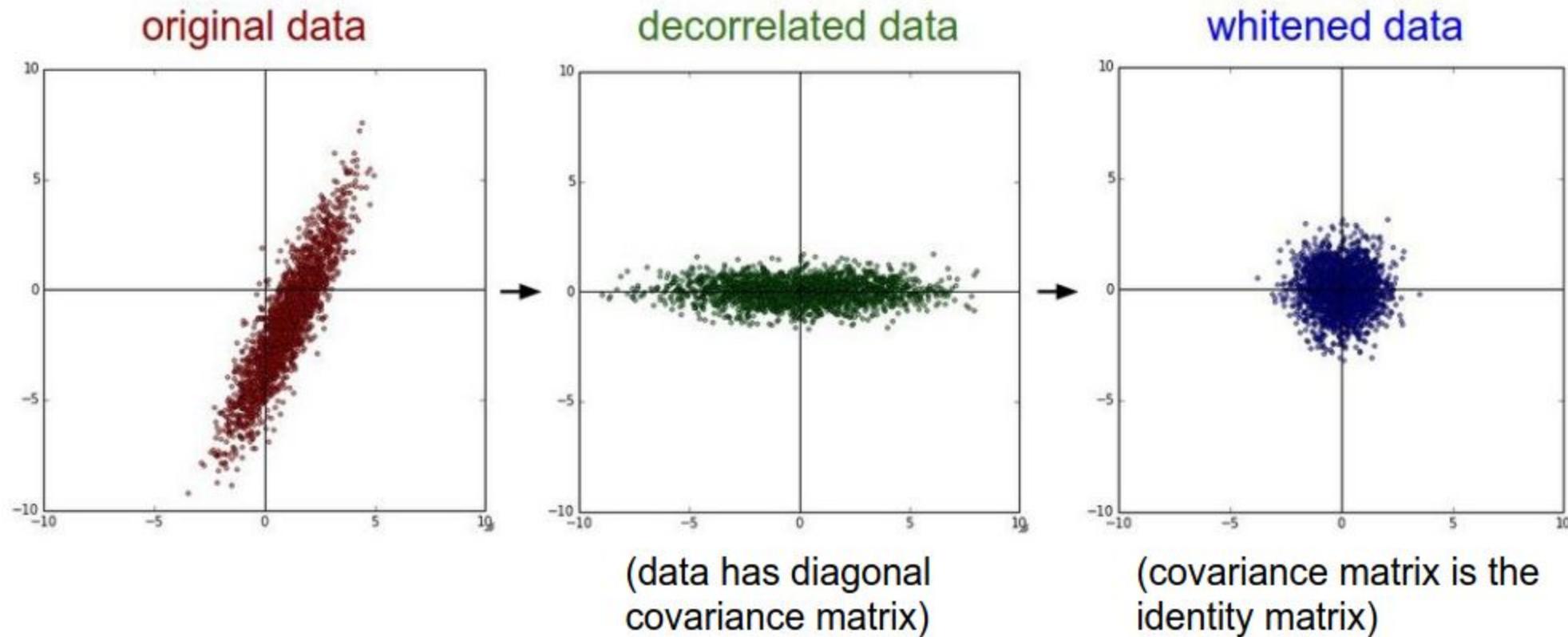


What can we say about the gradients on  $w$ ?  
Always all positive or all negative :(  
(this is also why you want zero-mean data!)



Michigan Tech

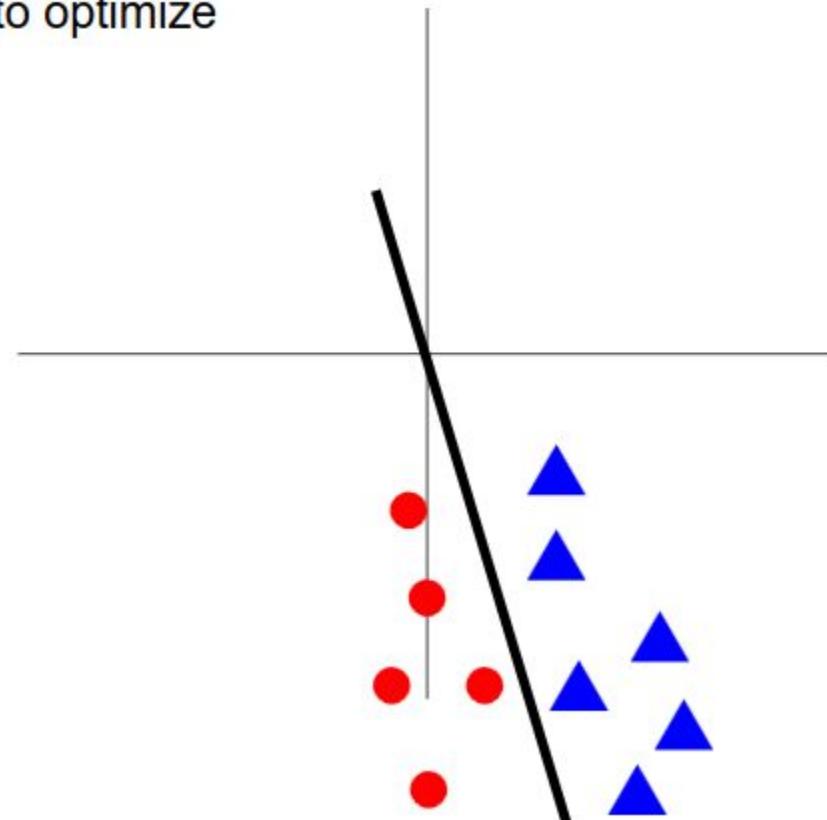
# Other practices



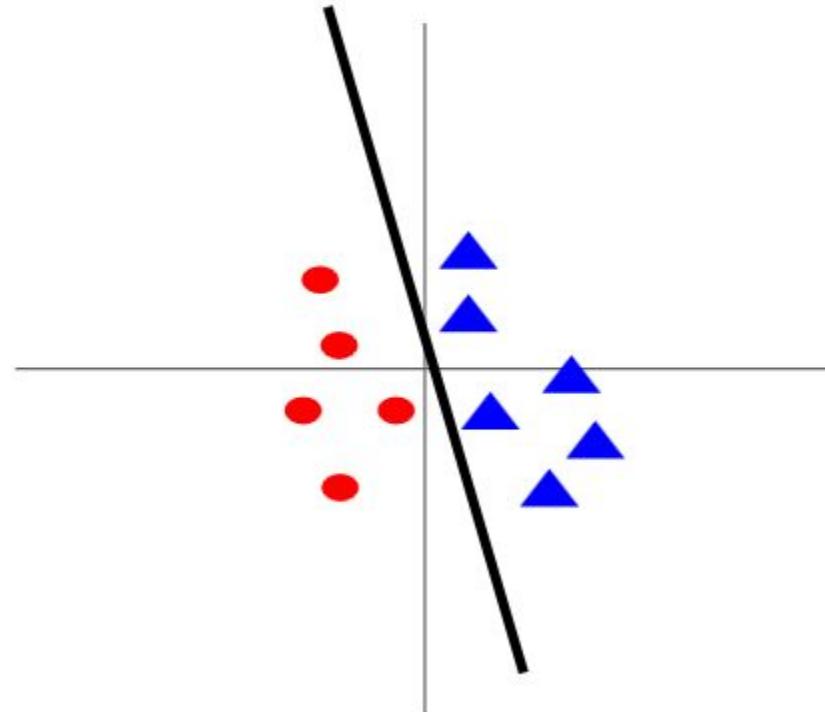
Michigan Tech

# Preprocessing effects

**Before normalization:** classification loss very sensitive to changes in weight matrix; hard to optimize



**After normalization:** less sensitive to small changes in weights; easier to optimize



Michigan Tech

# Other data specific pre-processing

For images:

- Subtract mean image
- Subtract per-channel mean
- Subtract per-channel mean and divide by per-channel std
- Typically don't do PCA or whitening



Michigan Tech

# Weight initialization

- Random
  - Works pretty well for small networks
  - Values must be small
  - zero mean,  $\text{e-}2$  std typical
  - Problematic with big networks



Michigan Tech

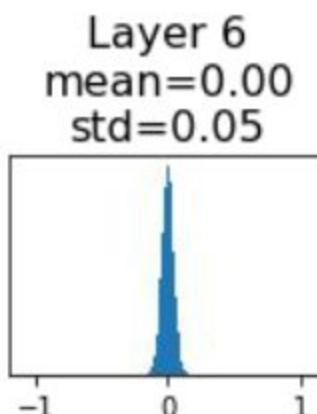
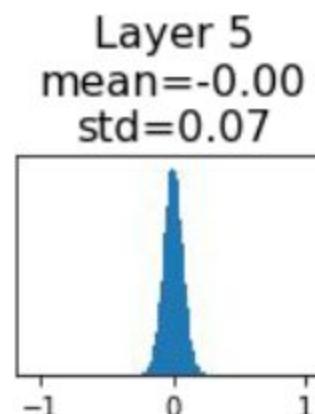
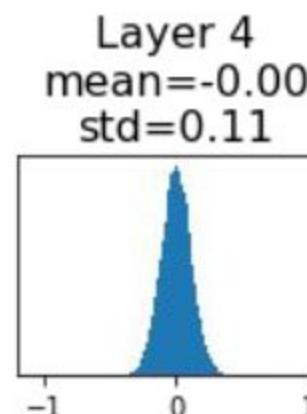
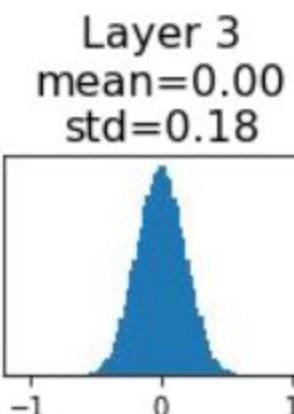
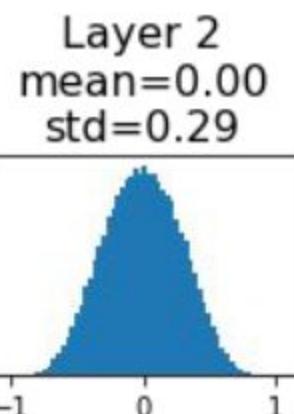
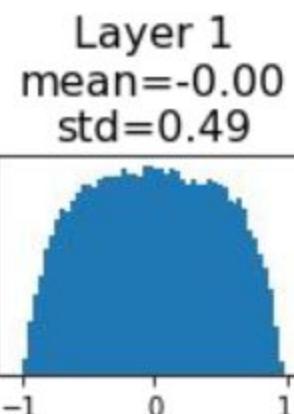
# Weight Initialization: Activation statistics

```
dims = [4096] * 7      Forward pass for a 6-layer  
hs = []                  net with hidden size 4096  
x = np.random.randn(16, dims[0])  
for Din, Dout in zip(dims[:-1], dims[1:]):  
    w = 0.01 * np.random.randn(Din, Dout)  
    x = np.tanh(x.dot(w))  
    hs.append(x)
```

All activations tend to zero for deeper network layers

**Q:** What do the gradients  $dL/dW$  look like?

**A:** All zero, no learning =(



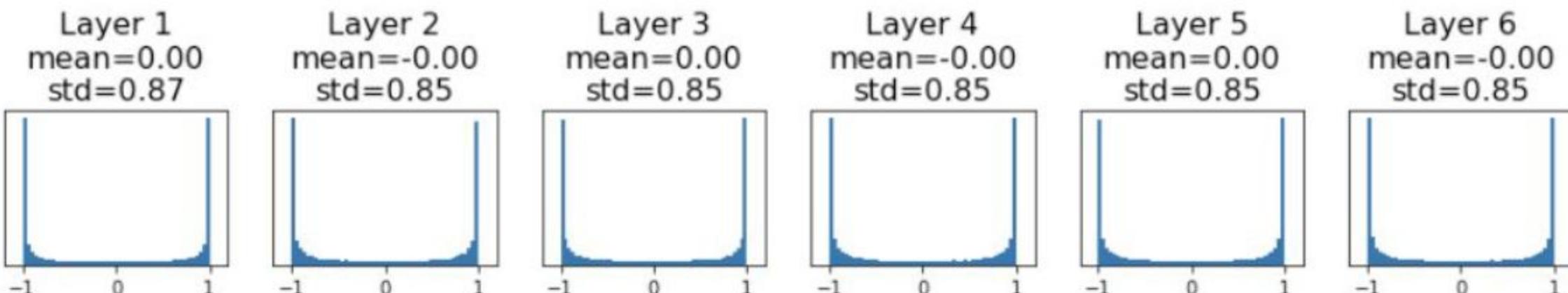
# Weight Initialization: Activation statistics

```
dims = [4096] * 7      Increase std of initial  
hs = []                  weights from 0.01 to 0.05  
x = np.random.randn(16, dims[0])  
for Din, Dout in zip(dims[:-1], dims[1:]):  
    W = 0.05 * np.random.randn(Din, Dout)  
    x = np.tanh(x.dot(W))  
    hs.append(x)
```

All activations saturate

**Q:** What do the gradients look like?

**A:** Local gradients all zero, no learning =(



Michigan Tech  
1885

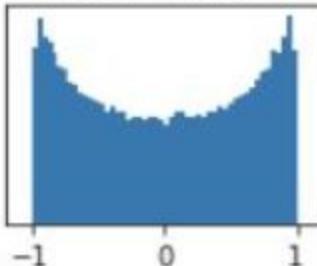
# Weight Initialization: “Xavier” Initialization

```
dims = [4096] * 7          "Xavier" initialization:  
hs = []                      std = 1/sqrt(Din)  
x = np.random.randn(16, dims[0])  
for Din, Dout in zip(dims[:-1], dims[1:]):  
    W = np.random.randn(Din, Dout) / np.sqrt(Din)  
    x = np.tanh(x.dot(W))  
    hs.append(x)
```

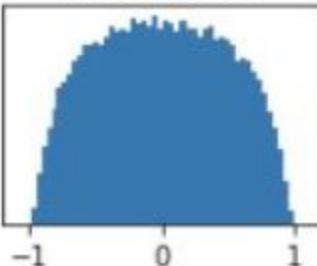
“Just right”: Activations are nicely scaled for all layers!

For conv layers,  $Din$  is  $\text{kernel\_size}^2 * \text{input\_channels}$

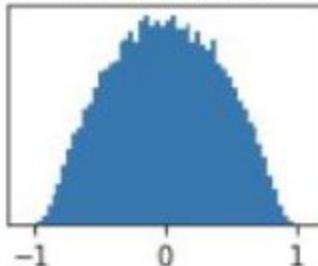
Layer 1  
mean=-0.00  
std=0.63



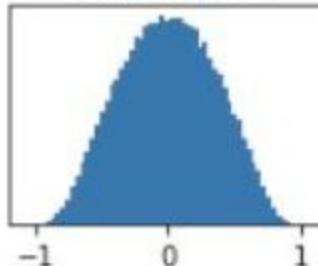
Layer 2  
mean=-0.00  
std=0.49



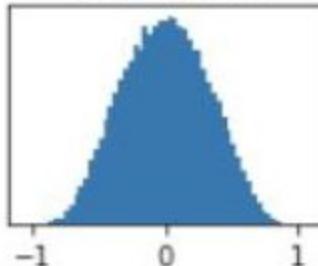
Layer 3  
mean=0.00  
std=0.41



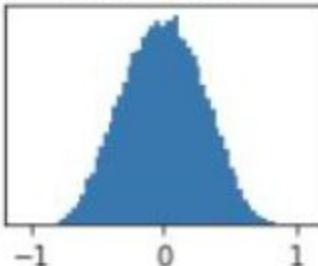
Layer 4  
mean=0.00  
std=0.36



Layer 5  
mean=0.00  
std=0.32



Layer 6  
mean=-0.00  
std=0.30



Glorot and Bengio, “Understanding the difficulty of training deep feedforward neural networks”, AISTAT 2010



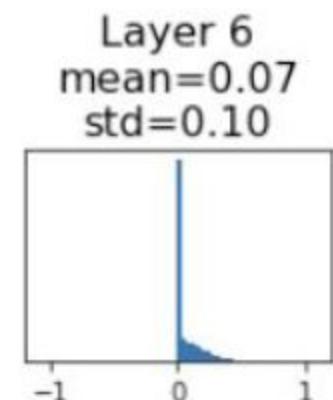
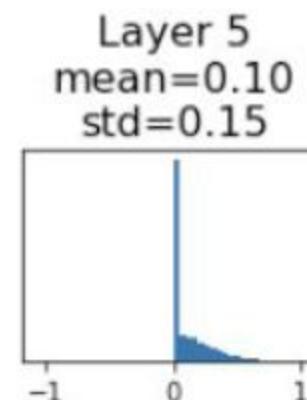
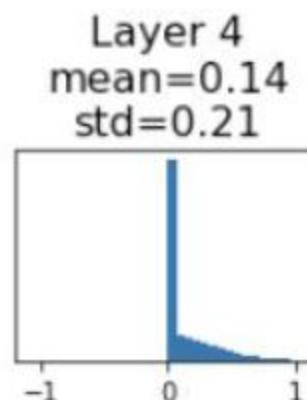
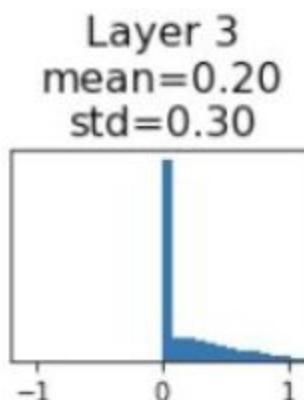
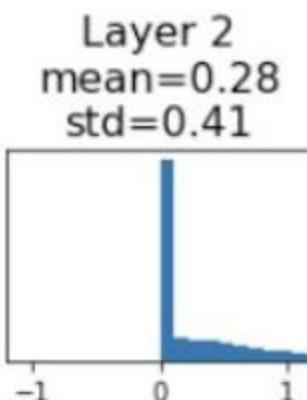
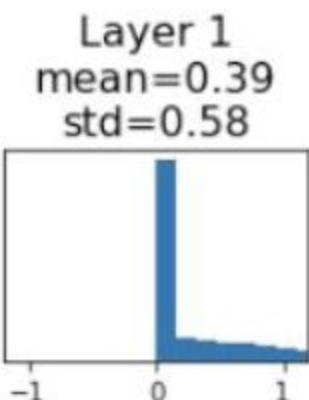
Michigan Tech

# Weight Initialization: What about ReLU?

```
dims = [4096] * 7      Change from tanh to ReLU
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.maximum(0, x.dot(W))
    hs.append(x)
```

Xavier assumes zero centered activation function

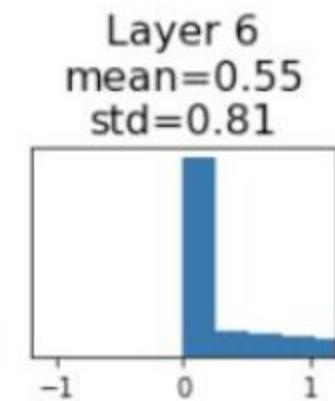
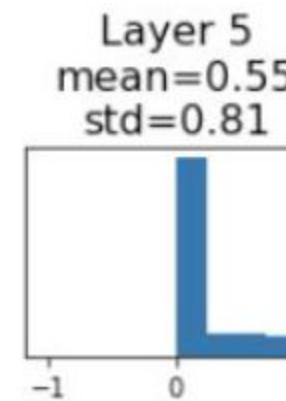
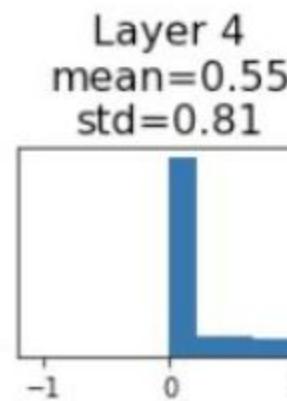
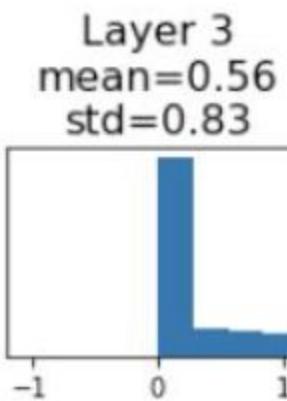
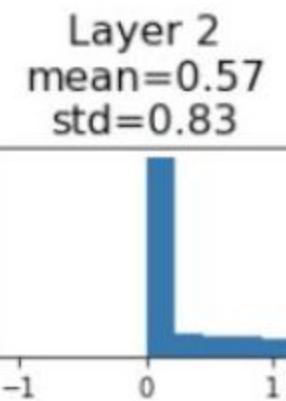
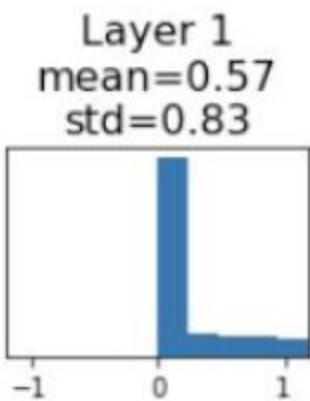
Activations collapse to zero again, no learning =(



# Weight Initialization: Kaiming / MSRA Initialization

```
dims = [4096] * 7  ReLU correction: std = sqrt(2 / Din)
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) * np.sqrt(2/Din)
    x = np.maximum(0, x.dot(W))
    hs.append(x)
```

“Just right”: Activations are nicely scaled for all layers!



He et al, “Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification”, ICCV 2015



Michigan Tech

# Proper initialization is an active area of research...

***Understanding the difficulty of training deep feedforward neural networks***

by Glorot and Bengio, 2010

***Exact solutions to the nonlinear dynamics of learning in deep linear neural networks*** by Saxe et al, 2013

***Random walk initialization for training very deep feedforward networks*** by Sussillo and Abbott, 2014

***Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification*** by He et al., 2015

***Data-dependent Initializations of Convolutional Neural Networks*** by Krähenbühl et al., 2015

***All you need is a good init***, Mishkin and Matas, 2015

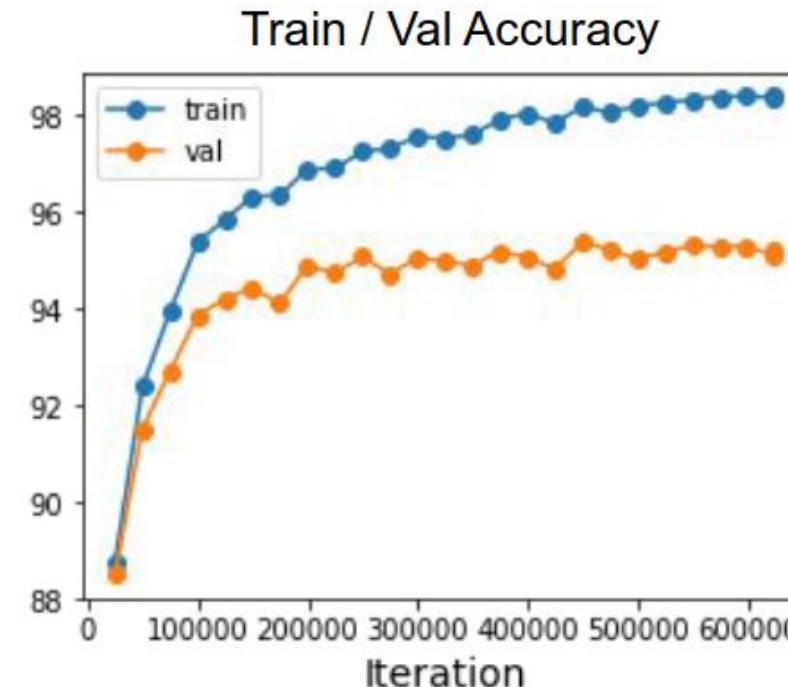
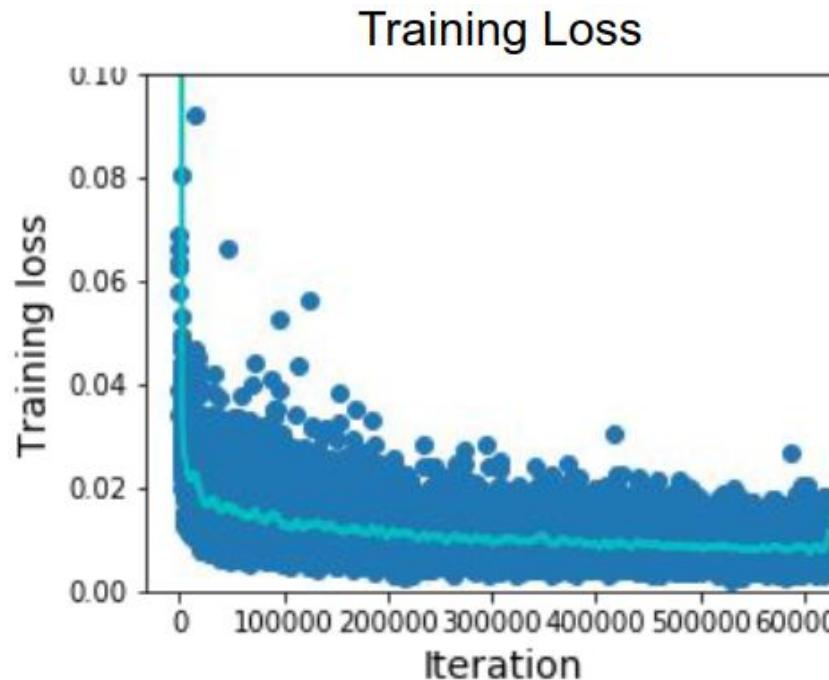
***Fixup Initialization: Residual Learning Without Normalization***, Zhang et al, 2019

***The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks***, Frankle and Carbin, 2019



Michigan Tech

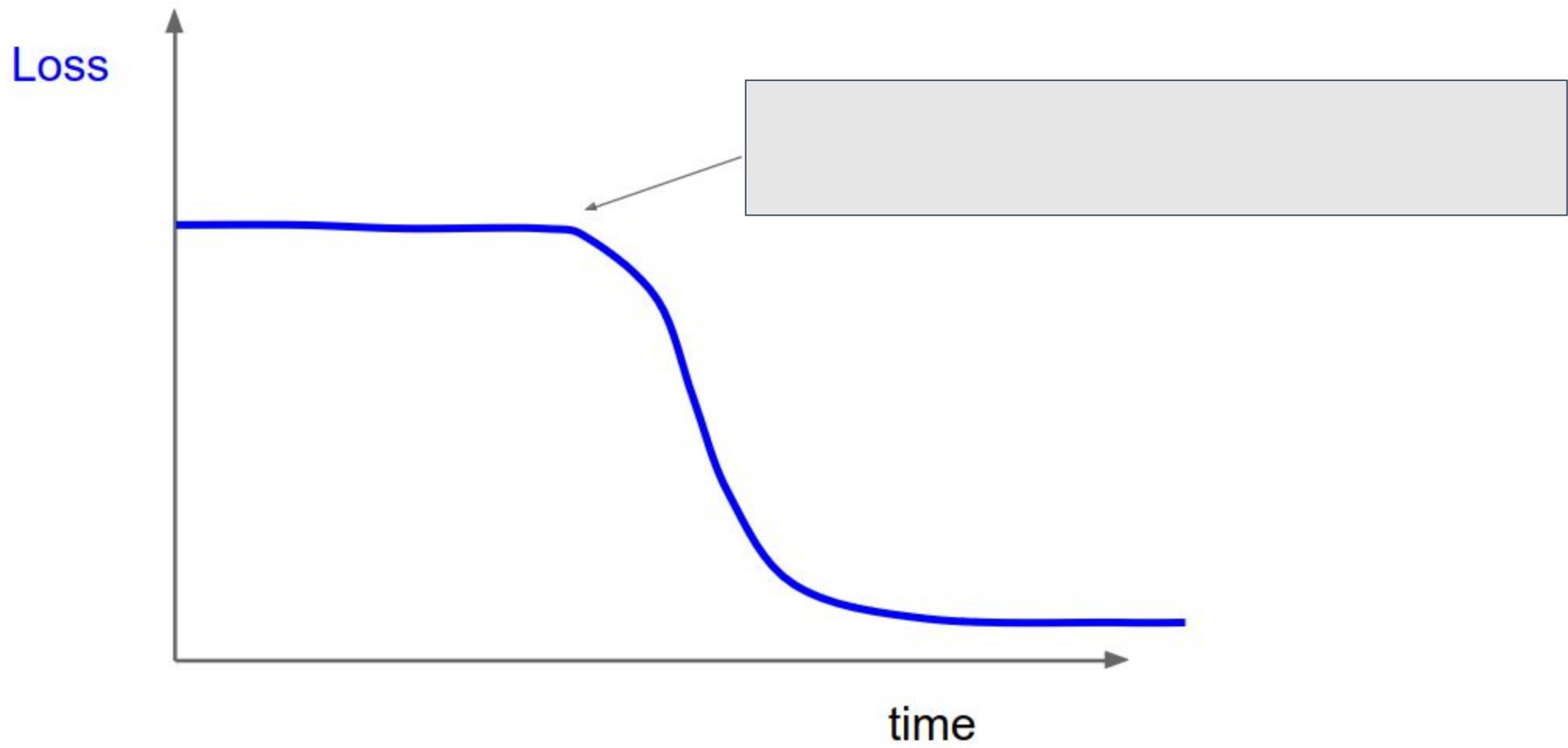
# Reading loss curves



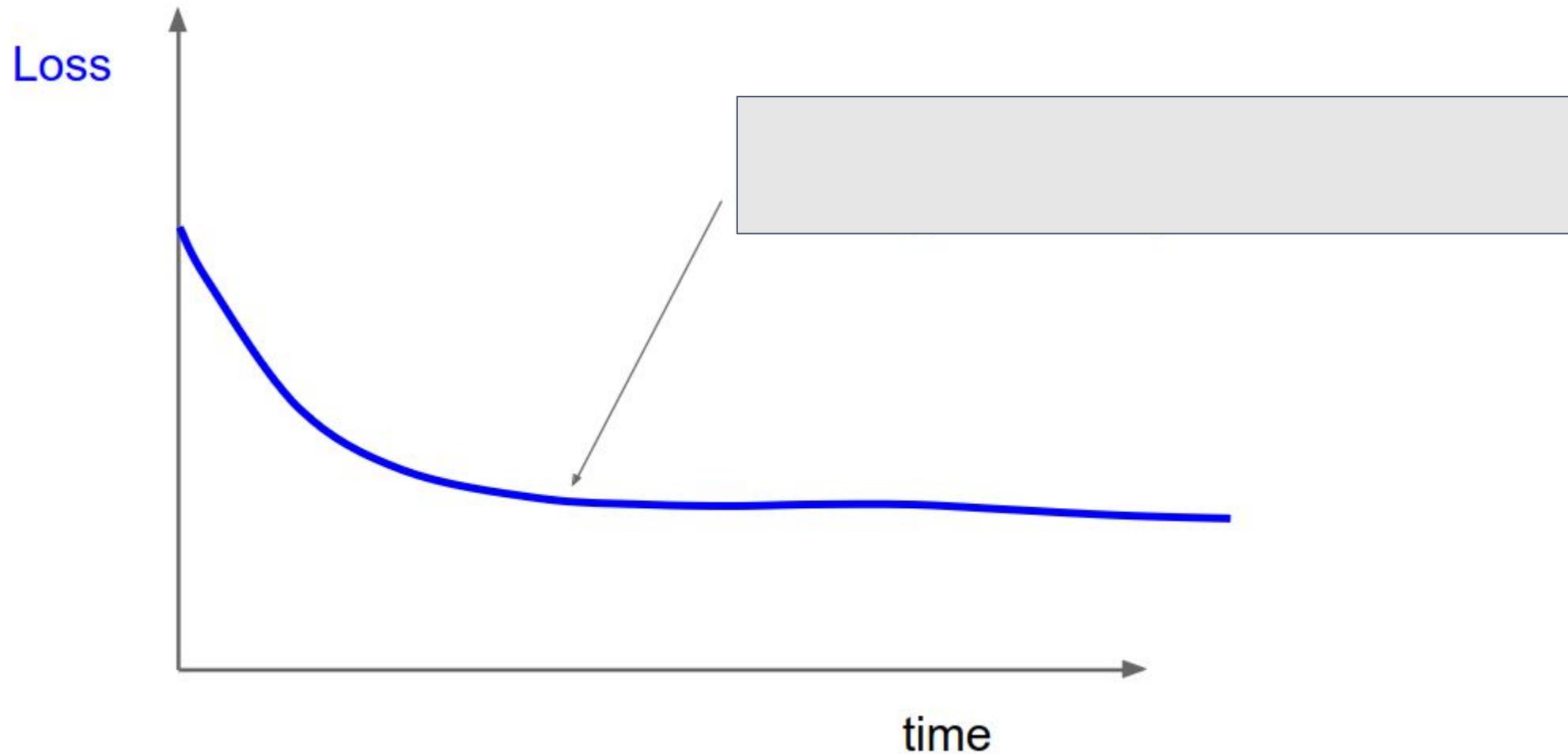
Losses may be noisy, use a  
scatter plot and also plot moving  
average to see trends better



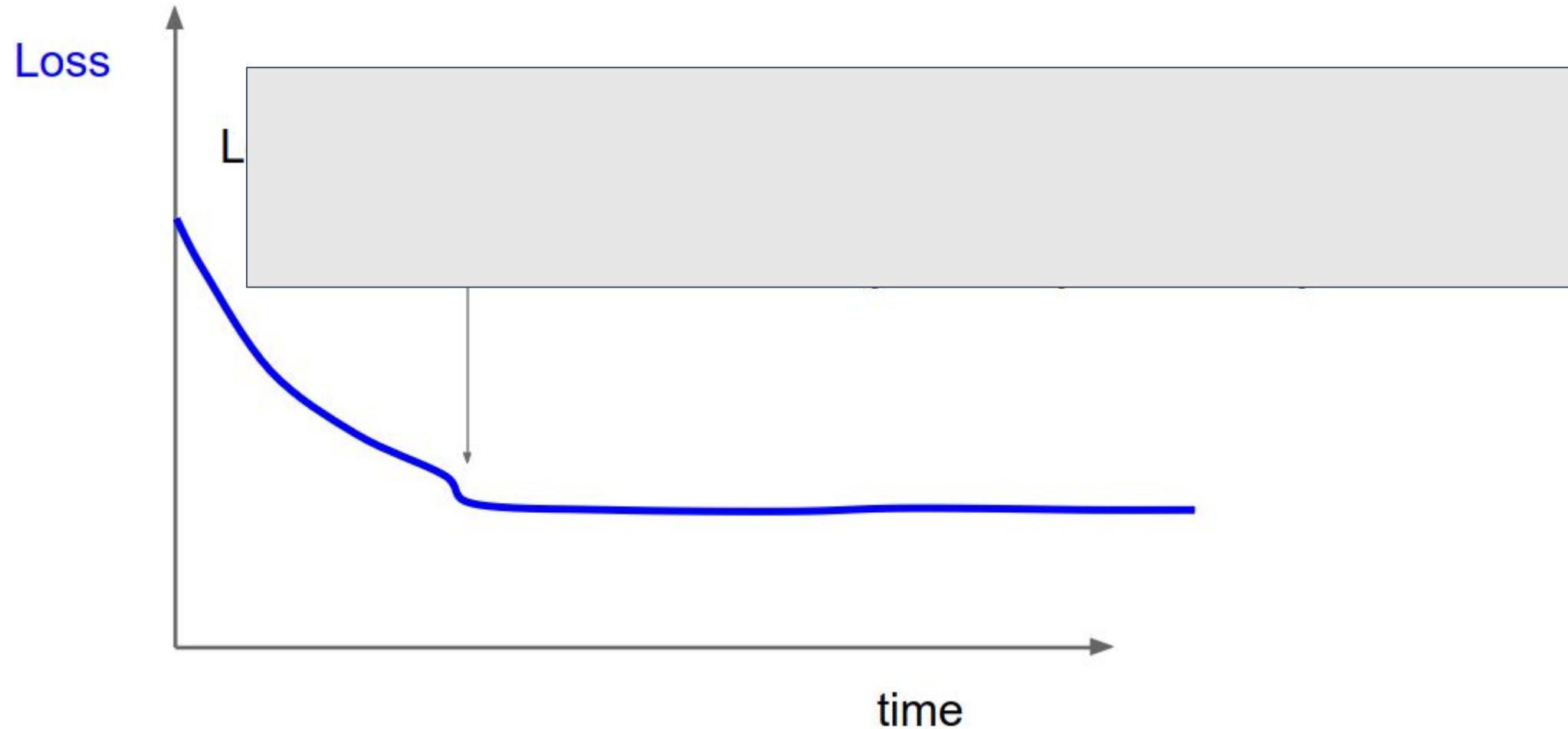
Michigan Tech



Michigan Tech



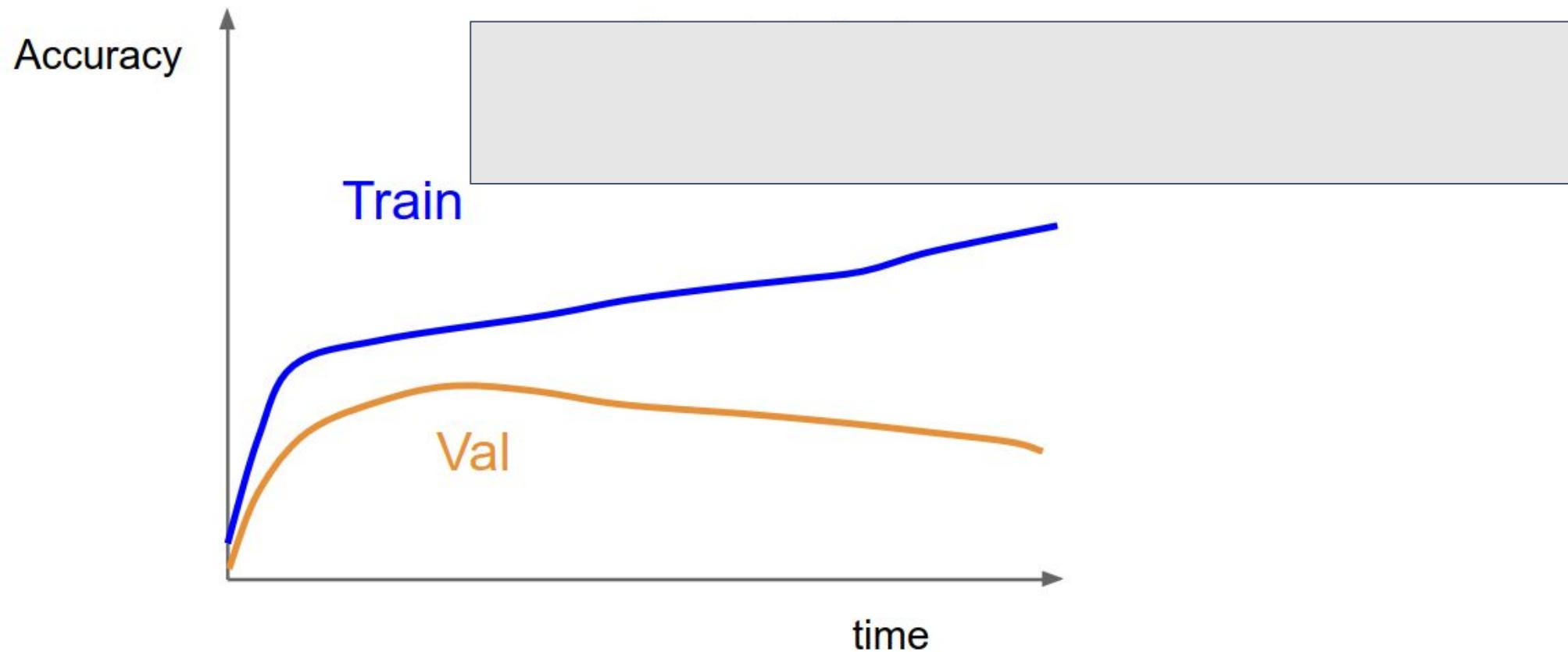
Michigan Tech



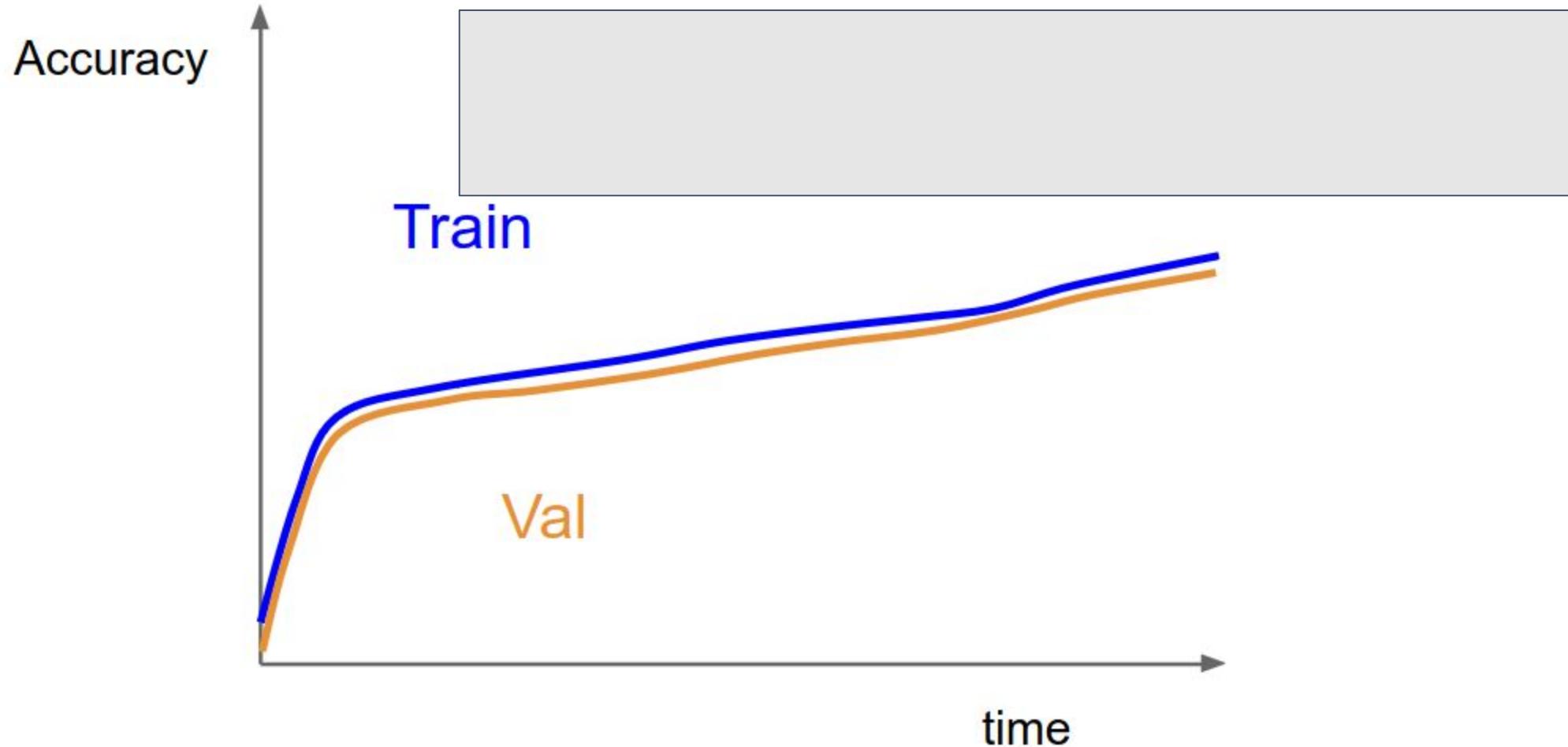
Michigan Tech



Michigan Tech



Michigan Tech



Michigan Tech

# Choosing Hyperparameters

**Step 1:** Check initial loss

**Step 2:** Overfit a small sample

**Step 3:** Find LR that makes loss go down

**Step 4:** Coarse grid, train for ~1-5 epochs

**Step 5:** Refine grid, train longer

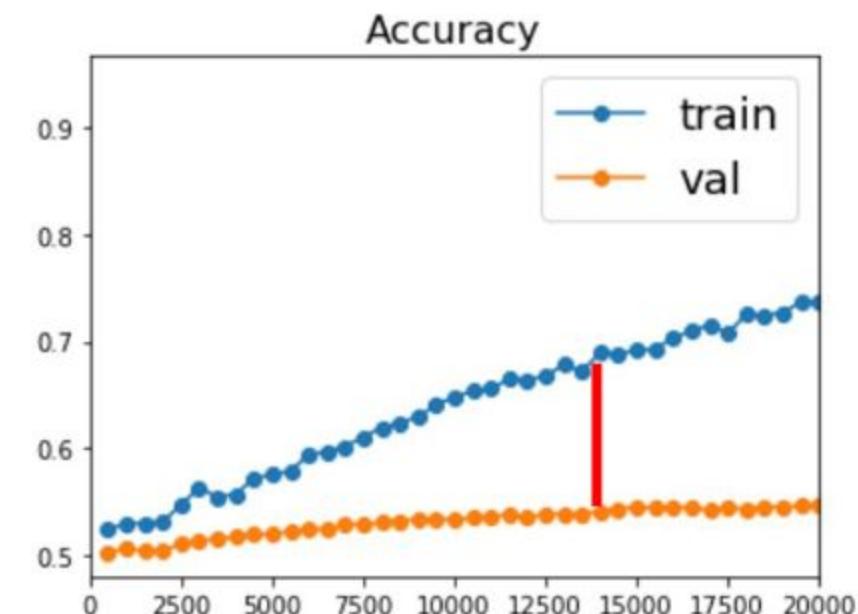
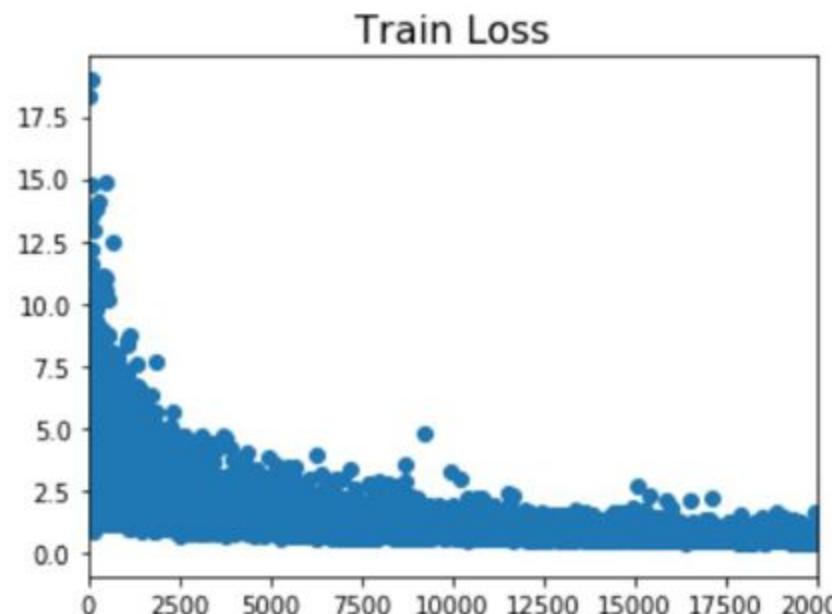
**Step 6:** Look at loss curves

**Step 7:** GOTO step 5



Michigan Tech

# How to improve single-model performance?



## Regularization



Michigan Tech

# Regularization: Add term to loss

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \boxed{\lambda R(W)}$$

In common use:

**L2 regularization**  $R(W) = \sum_k \sum_l W_{k,l}^2$  (Weight decay)

**L1 regularization**  $R(W) = \sum_k \sum_l |W_{k,l}|$

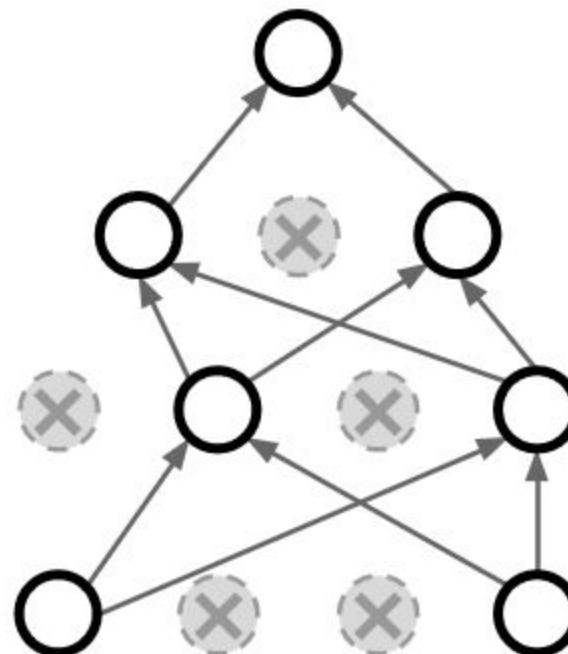
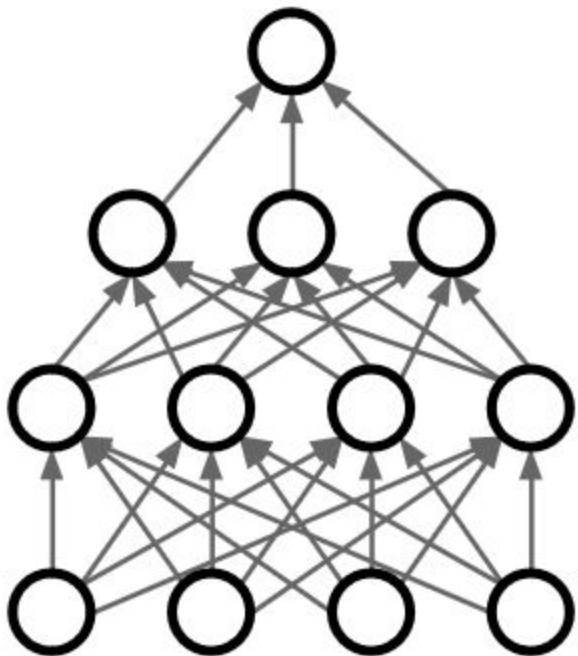
**Elastic net (L1 + L2)**  $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$



Michigan Tech

# Regularization: Dropout

In each forward pass, randomly set some neurons to zero  
Probability of dropping is a hyperparameter; 0.5 is common



Srivastava et al, "Dropout: A simple way to prevent neural networks from overfitting", JMLR 2014



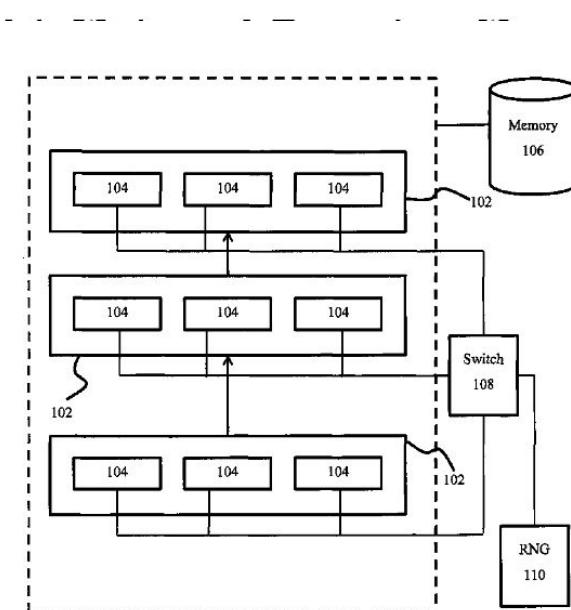
Michigan Tech



(57)

## ABSTRACT

A system for training a neural network. A switch is linked to feature detectors in at least some of the layers of the neural network. For each training case, the switch randomly selectively disables each of the feature detectors in accordance with a preconfigured probability. The weights from each training case are then normalized for applying the neural network to test data.



(12) **United States Patent**  
Hinton et al.

(10) **Patent No.:** US 9,406,017 B2  
(45) **Date of Patent:** Aug. 2, 2016

(54) **SYSTEM AND METHOD FOR ADDRESSING OVERFITTING IN A NEURAL NETWORK**

(71) Applicant: **Google Inc.**, Mountain View, CA (US)

(72) Inventors: **Geoffrey E. Hinton**, Toronto (CA);  
**Alexander Krizhevsky**, Toronto (CA);  
**Ilya Sutskever**, Mountain View, CA (US); **Nitish Srivastva**, Toronto (CA)

(73) Assignee: **Google Inc.**, Mountain View, CA (US)

(\*) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 369 days.

(21) Appl. No.: **14/015,768**

(22) Filed: **Aug. 30, 2013**

(65) **Prior Publication Data**

US 2014/0180986 A1 Jun. 26, 2014

### Related U.S. Application Data

(60) Provisional application No. 61/745,711, filed on Dec. 24, 2012.

(51) **Int. Cl.**  
**G06N 3/08** (2006.01)  
**G06N 3/04** (2006.01)

(52) **U.S. Cl.**  
CPC ..... **G06N 3/082** (2013.01); **G06N 3/0454** (2013.01)

(58) **Field of Classification Search**  
CPC ..... G06N 3/082  
See application file for complete search history.

(56) **References Cited**  
PUBLICATIONS

Castellano, Giovanna, Anna Maria Fanelli, and Marcello Pelillo. "An iterative pruning algorithm for feedforward neural networks." *Neural Networks, IEEE Transactions on* 8.3 (1997): 519-531.\*  
Rumelhart et al., "Learning representations by back-propagating errors," *Nature* 323:533-536, Oct. 9, 1986.  
Hinton, "Training Products of Experts by Minimizing Contrastive Divergence," *Neural Computation* 14(8):1771-1800, Aug. 2002.  
Ciresan et al., "Deep, Big, Simple Neural Nets for Handwritten Digit Recognition," *Neural Computation* 22(12):3207-3220, Dec. 2010.  
LeCun et al., "Gradient-based learning applied to document recognition," *Proceedings of the IEEE* 86(11):2278-2324, Nov. 1998.  
Hinton et al., "Reducing the Dimensionality of Data with Neural Networks," *Science* 313(5768):504-507, Jul. 28, 2006.  
Mohamed et al., "Acoustic Modeling Using Deep Belief Networks," *IEEE Transactions on Audio, Speech, and Language Processing*, 20(1):14-22, Jan. 2012.

Dahl et al., "Context-Dependent Pre-Trained Deep Neural Networks for Large-Vocabulary Speech Recognition," *IEEE Transactions on Audio, Speech, and Language Processing*, 20:30-42, Jan. 2012.  
Jaitly et al., "Application of Pretrained Deep Neural Networks to Large Vocabulary Conversational Speech Recognition," *Tech. Rep. 001*, Department of Computer Science, University of Toronto, Mar. 12, 2012, 11 pages.

Krizhevsky, "Learning multiple layers of features from tiny images," *Tech. Rep. 001*, Department of Computer Science, University of Toronto, Apr. 8, 2009, 60 pages.

(Continued)

*Primary Examiner* — Kakali Chaki

*Assistant Examiner* — Kevin W Figueroa

(74) *Attorney, Agent, or Firm* — Fish & Richardson P.C.

### ABSTRACT

A system for training a neural network. A switch is linked to feature detectors in at least some of the layers of the neural network. For each training case, the switch randomly selectively disables each of the feature detectors in accordance with a preconfigured probability. The weights from each training case are then normalized for applying the neural network to test data.

**24 Claims, 2 Drawing Sheets**



**Michigan Tech**

# Regularization: Dropout

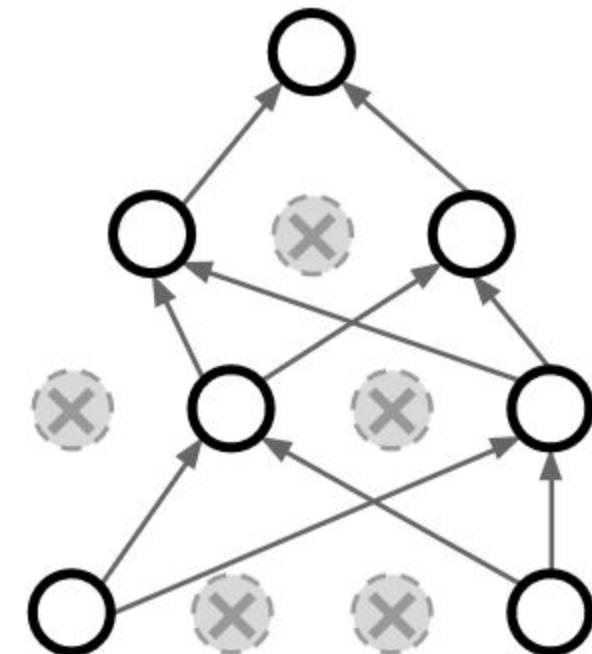
```
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    """ X contains the data """

    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)
```

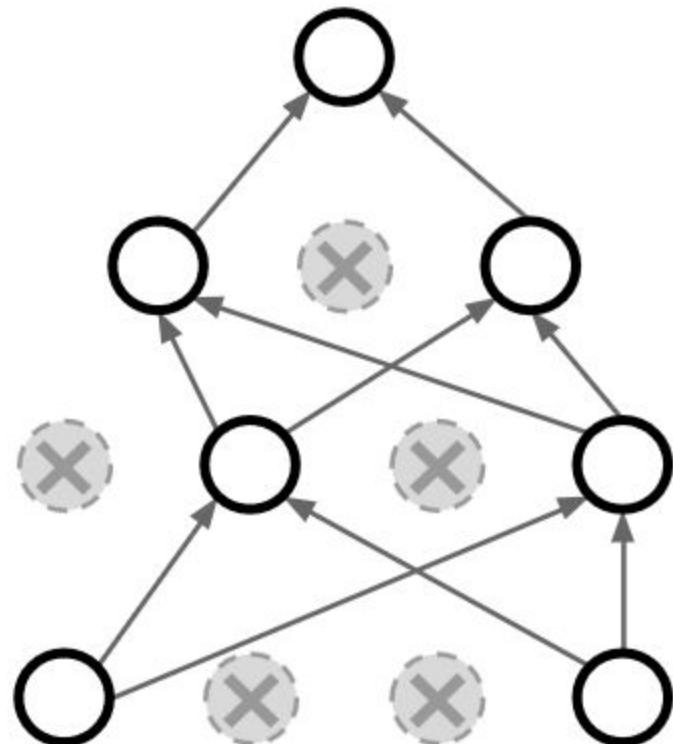
Example forward pass with a 3-layer network using dropout



Michigan Tech

# Regularization: Dropout

How can this possibly be a good idea?



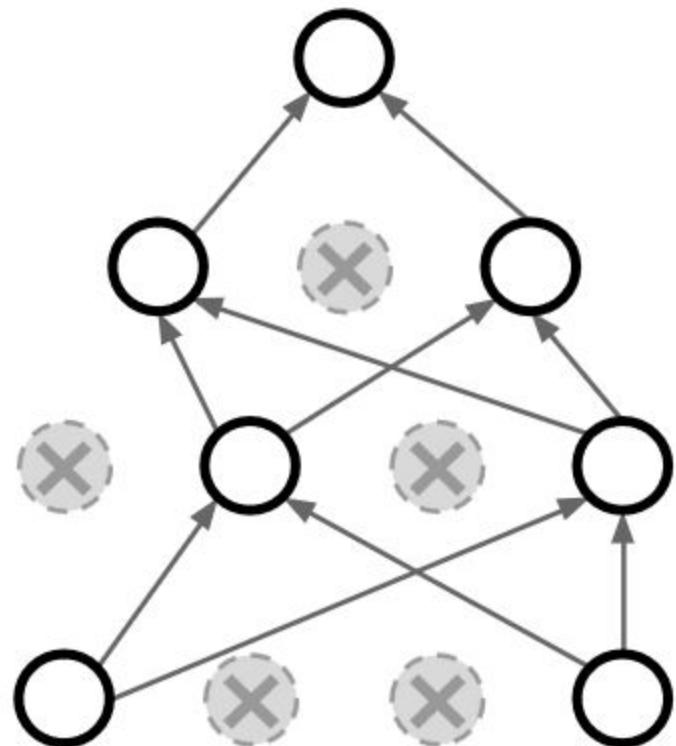
Forces the network to have a redundant representation;  
Prevents co-adaptation of features



Michigan Tech

# Regularization: Dropout

How can this possibly be a good idea?



Another interpretation:

Dropout is training a large **ensemble** of models (that share parameters).

Each binary mask is one model

An FC layer with 4096 units has  $2^{4096} \sim 10^{1233}$  possible masks!  
Only  $\sim 10^{82}$  atoms in the universe...



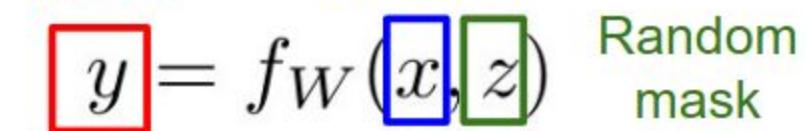
Michigan Tech

# Dropout: Test time

Dropout makes our output random!

$$y = f_W(x, z)$$

Output (label)      Input (image)      Random mask



Want to “average out” the randomness at test-time

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

But this integral seems hard ...



Michigan Tech

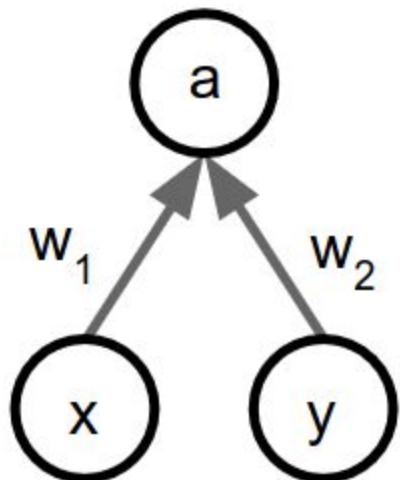
# Dropout: Test time

Want to approximate  
the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.

At test time we have:  $E[a] = w_1x + w_2y$



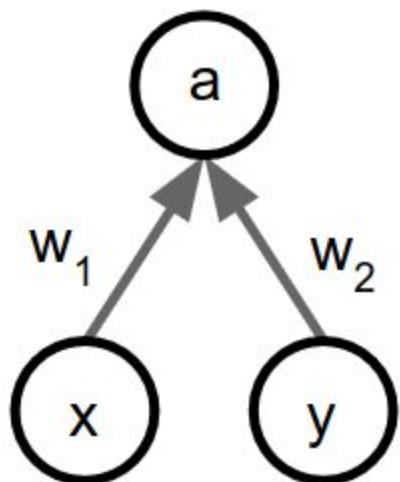
Michigan Tech

# Dropout: Test time

Want to approximate  
the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.



At test time we have:  $E[a] = w_1x + w_2y$

During training we have:

$$\begin{aligned} E[a] &= \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) \\ &\quad + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) \\ &= \frac{1}{2}(w_1x + w_2y) \end{aligned}$$



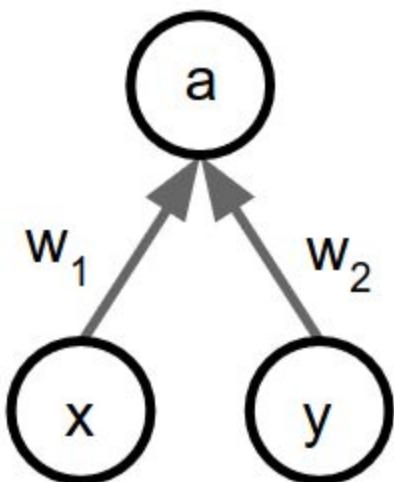
Michigan Tech

# Dropout: Test time

Want to approximate  
the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.



At test time we have:  $E[a] = w_1x + w_2y$

During training we have:

$$\begin{aligned} E[a] &= \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) \\ &\quad + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) \\ &= \frac{1}{2}(w_1x + w_2y) \end{aligned}$$

At test time, multiply  
by dropout probability



Michigan Tech

# Dropout: Test time

```
def predict(X):
    # ensembled forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations
    out = np.dot(W3, H2) + b3
```

At test time all neurons are active always

=> We must scale the activations so that for each neuron:  
output at test time = expected output at training time



Michigan Tech

# Dropout Summary

```
""" Vanilla Dropout: Not recommended implementation (see notes below) """  
  
p = 0.5 # probability of keeping a unit active. higher = less dropout  
  
def train_step(X):  
    """ X contains the data """  
  
    # forward pass for example 3-layer neural network  
    H1 = np.maximum(0, np.dot(W1, X) + b1)  
    U1 = np.random.rand(*H1.shape) < p # first dropout mask  
    H1 *= U1 # drop!  
    H2 = np.maximum(0, np.dot(W2, H1) + b2)  
    U2 = np.random.rand(*H2.shape) < p # second dropout mask  
    H2 *= U2 # drop!  
    out = np.dot(W3, H2) + b3  
  
    # backward pass: compute gradients... (not shown)  
    # perform parameter update... (not shown)  
  
def predict(X):  
    # ensembled forward pass  
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations  
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations  
    out = np.dot(W3, H2) + b3
```

drop in forward pass

scale at test time



Michigan Tech

# More common: “Inverted dropout”

```
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = (np.random.rand(*H1.shape) < p) / p # first dropout mask. Notice /p!
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = (np.random.rand(*H2.shape) < p) / p # second dropout mask. Notice /p!
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)

def predict(X):
    # ensembled forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) # no scaling necessary
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    out = np.dot(W3, H2) + b3
```

test time is unchanged!



Michigan Tech  
1885

# Regularization: A common pattern

**Training:** Add some kind  
of randomness

$$y = f_W(x, z)$$

**Testing:** Average out randomness  
(sometimes approximate)

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$



Michigan Tech

# Regularization: A common pattern

**Training:** Add some kind of randomness

$$y = f_W(x, z)$$

**Testing:** Average out randomness (sometimes approximate)

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

**Example:** Batch Normalization

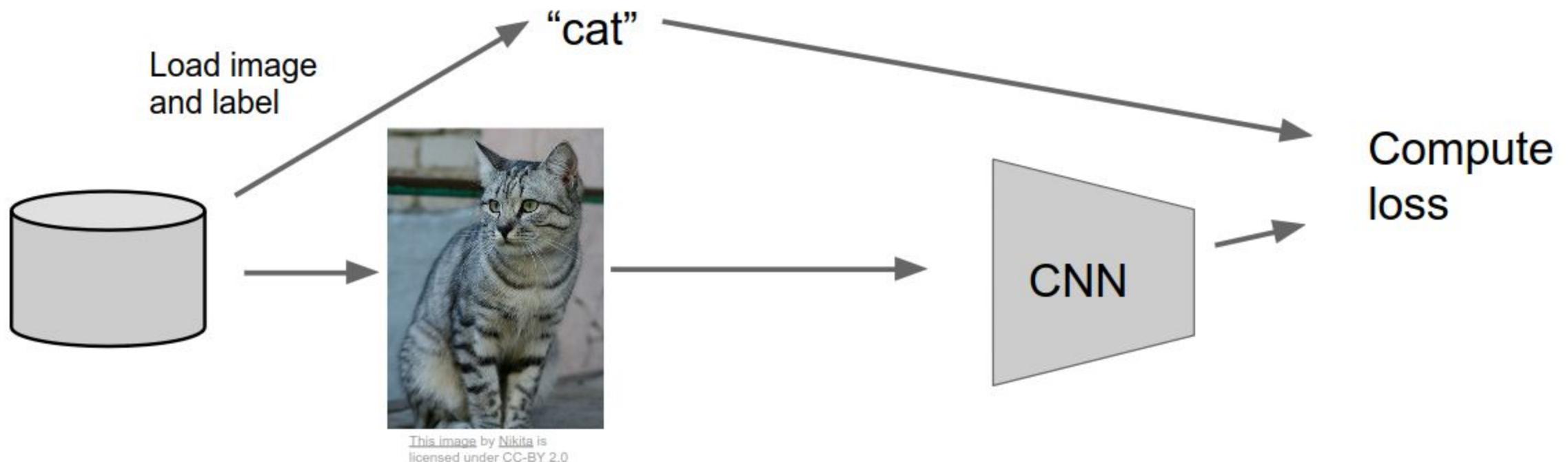
**Training:** Normalize using stats from random minibatches

**Testing:** Use fixed stats to normalize



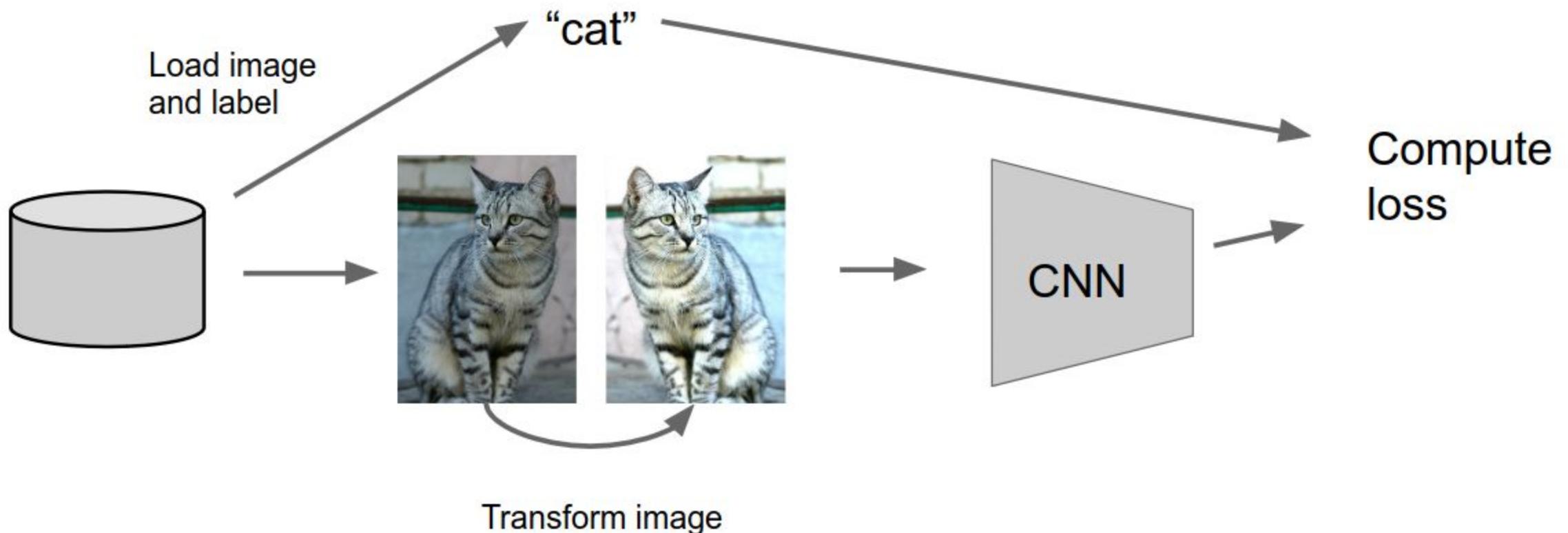
Michigan Tech

# Regularization: Data Augmentation



Michigan Tech

# Regularization: Data Augmentation



Michigan Tech

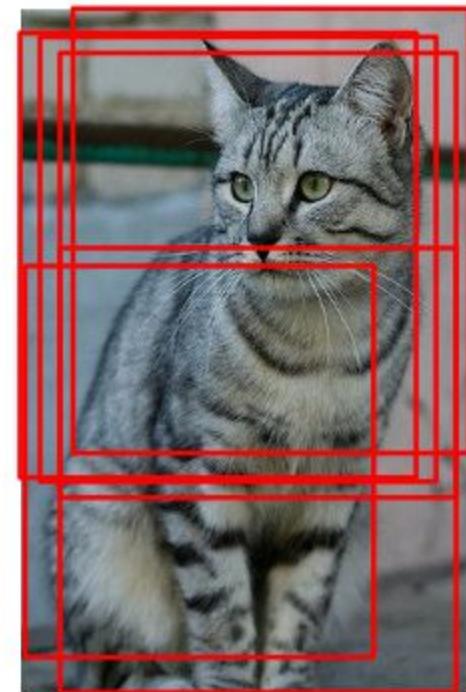
# Data Augmentation

## Random crops and scales

**Training:** sample random crops / scales

ResNet:

1. Pick random  $L$  in range [256, 480]
2. Resize training image, short side =  $L$
3. Sample random  $224 \times 224$  patch



Michigan Tech

# Data Augmentation

## Color Jitter

Simple: Randomize  
contrast and brightness



Michigan Tech

# Regularization: Cutout

**Training:** Set random image regions to zero

**Testing:** Use full image

## Examples:

Dropout

Batch Normalization

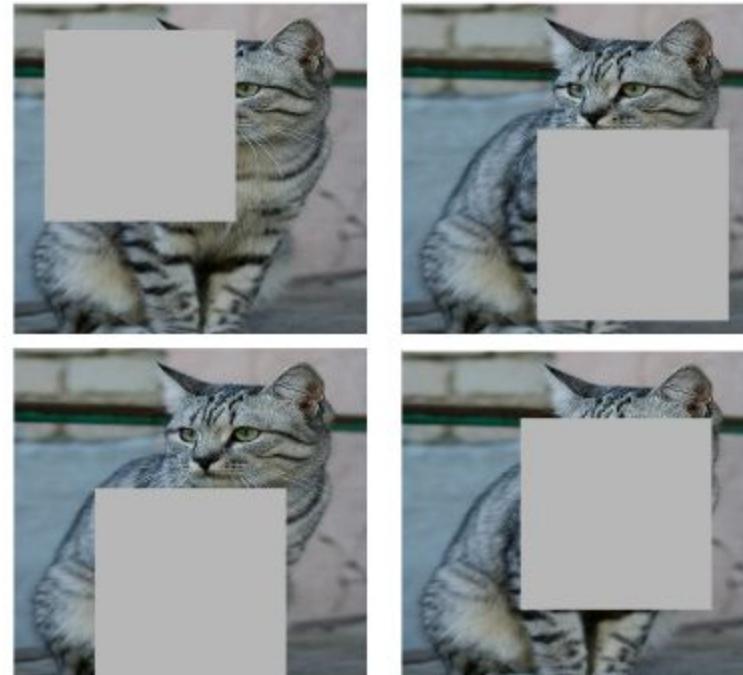
Data Augmentation

DropConnect

Fractional Max Pooling

Stochastic Depth

Cutout



Works very well for small datasets like CIFAR,  
less common for large datasets like ImageNet

DeVries and Taylor, "Improved Regularization of  
Convolutional Neural Networks with Cutout", arXiv 2017



Michigan Tech

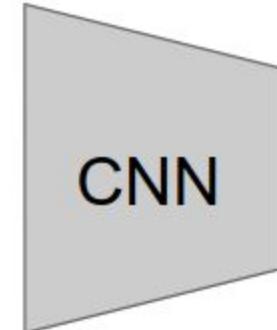
# Regularization: Mixup

**Training:** Train on random blends of images

**Testing:** Use original images

## Examples:

- Dropout
- Batch Normalization
- Data Augmentation
- DropConnect
- Fractional Max Pooling
- Stochastic Depth
- Cutout
- Mixup



Target label:  
cat: 0.4  
dog: 0.6

Randomly blend the pixels  
of pairs of training images,  
e.g. 40% cat, 60% dog

Zhang et al, "mixup: Beyond Empirical Risk Minimization", ICLR 2018



Michigan Tech

# Batch Normalization

[Ioffe and Szegedy, 2015]

“you want zero-mean unit-variance activations? just make them so.”

consider a batch of activations at some layer. To make each dimension zero-mean unit-variance, apply:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \text{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

this is a vanilla  
differentiable function...

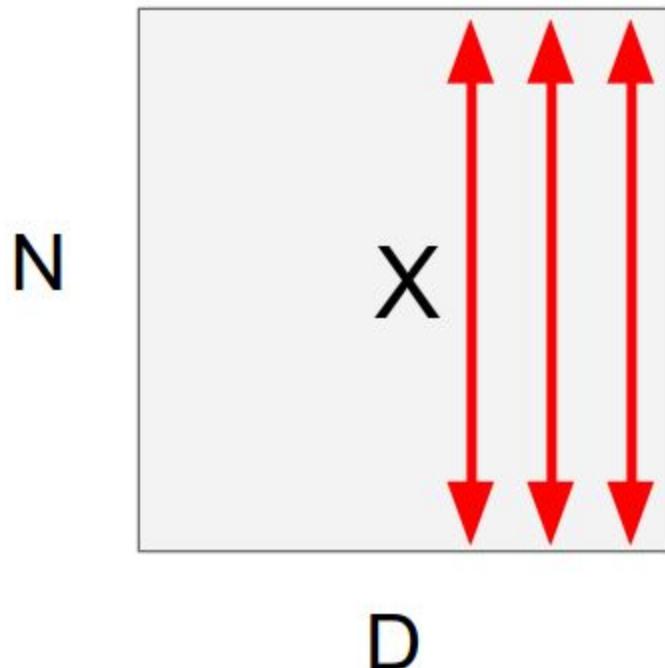


Michigan Tech

# Batch Normalization

[Ioffe and Szegedy, 2015]

**Input:**  $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,  
shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,  
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Normalized x,  
Shape is N x D



Michigan Tech

# Batch Normalization

[Ioffe and Szegedy, 2015]

**Input:**  $x : N \times D$

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,  
shape is D

**Learnable scale and shift parameters:**

$\gamma, \beta : D$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,  
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Normalized x,  
Shape is  $N \times D$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output,  
Shape is  $N \times D$

Learning  $\gamma = \sigma$ ,  
 $\beta = \mu$  will recover the  
identity function!



Michigan Tech

# Batch Normalization: Test-Time

**Input:**  $x : N \times D$

$\mu_j =$  (Running) average of  
values seen during training

Per-channel mean,  
shape is D

**Learnable scale and  
shift parameters:**

$\gamma, \beta : D$

$\sigma_j^2 =$  (Running) average of  
values seen during training

Per-channel var,  
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Normalized x,  
Shape is  $N \times D$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output,  
Shape is  $N \times D$

During testing batchnorm  
becomes a linear operator!

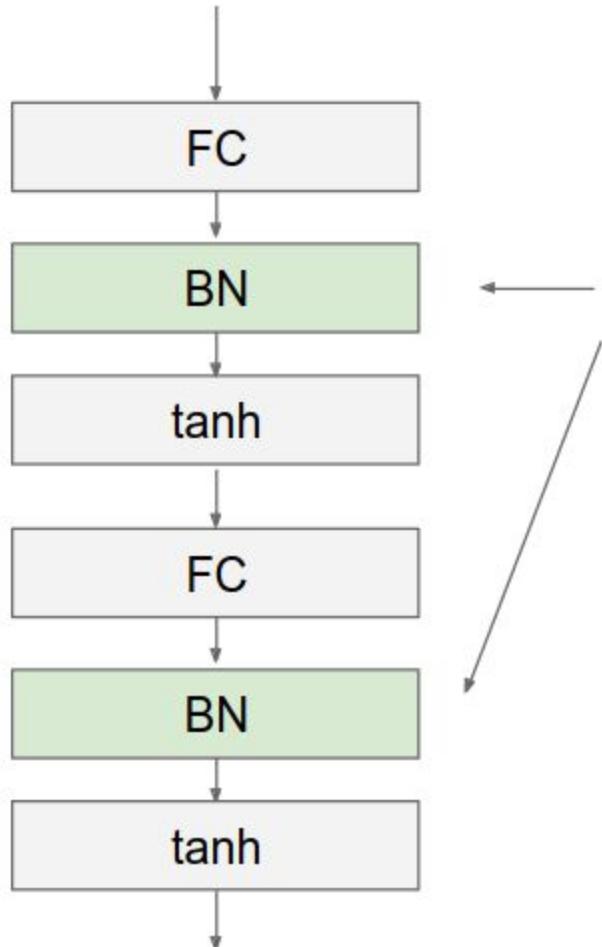
Can be fused with the previous  
fully-connected or conv layer



Michigan Tech

# Batch Normalization

[Ioffe and Szegedy, 2015]



Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

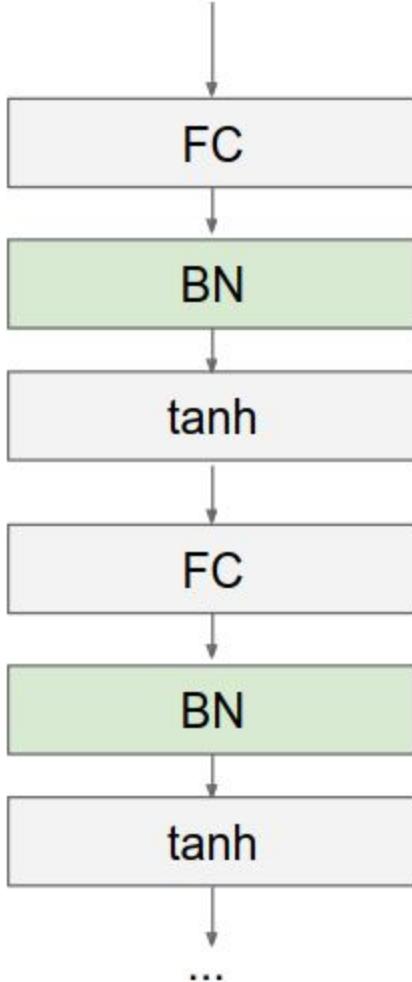
$$\hat{x}^{(k)} = \frac{x^{(k)} - \text{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$



Michigan Tech

# Batch Normalization

[Ioffe and Szegedy, 2015]



- Makes deep networks **much** easier to train!
- Improves gradient flow
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!
- **Behaves differently during training and testing: this is a very common source of bugs!**



Michigan Tech

# Layer Normalization

**Batch Normalization** for  
fully-connected networks

**Layer Normalization** for  
fully-connected networks  
Same behavior at train and test!  
Can be used in recurrent networks

$\mathbf{x}: N \times D$

Normalize



$\mu, \sigma: 1 \times D$

$\gamma, \beta: 1 \times D$

$$\mathbf{y} = \gamma(\mathbf{x} - \mu) / \sigma + \beta$$

$\mathbf{x}: N \times D$

Normalize



$\mu, \sigma: N \times 1$

$\gamma, \beta: 1 \times D$

$$\mathbf{y} = \gamma(\mathbf{x} - \mu) / \sigma + \beta$$

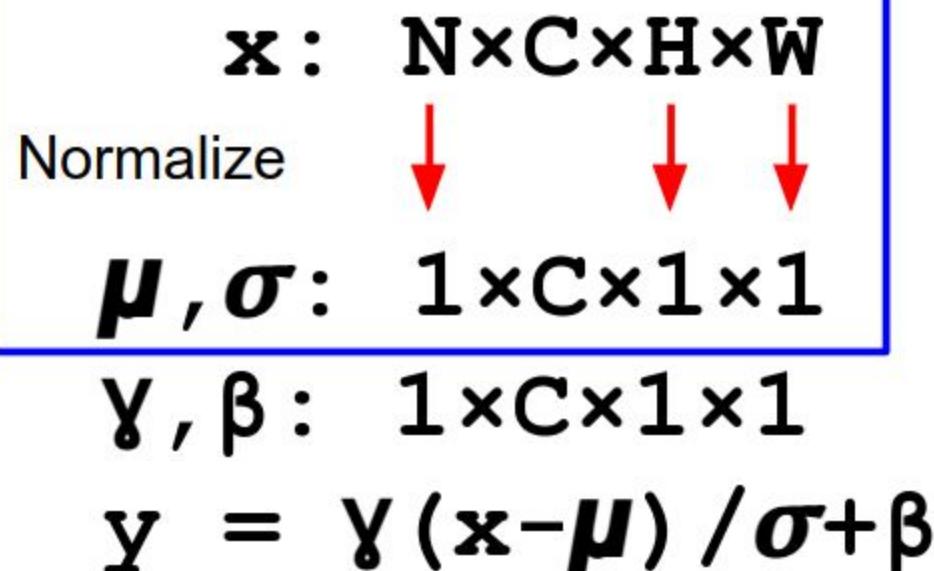
Ba, Kiros, and Hinton, "Layer Normalization", arXiv 2016



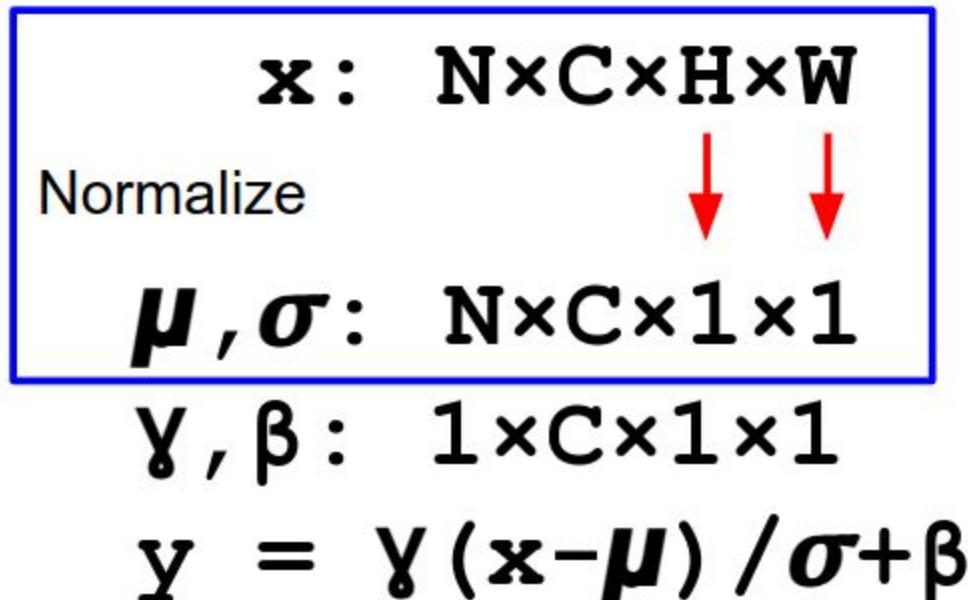
Michigan Tech

# Instance Normalization

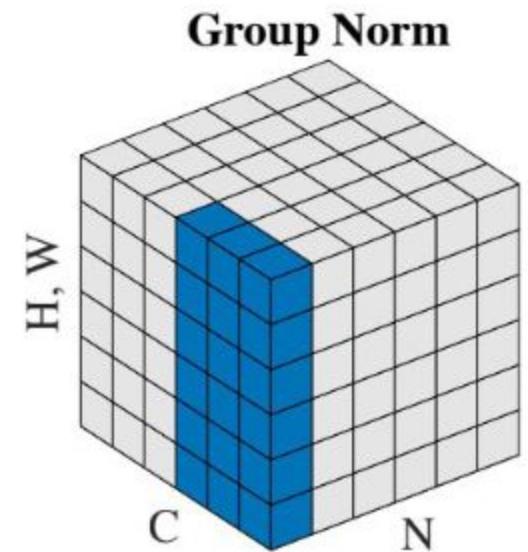
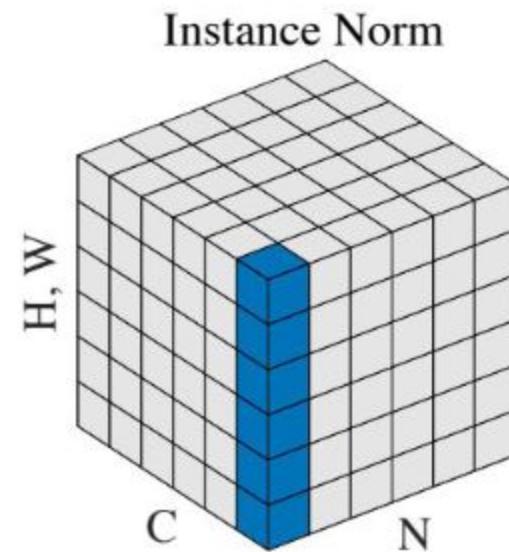
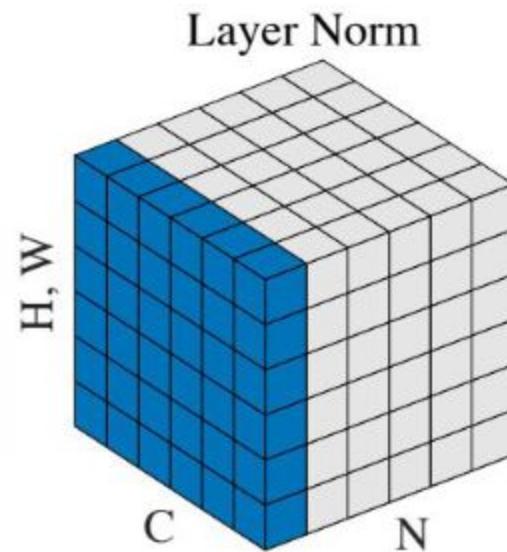
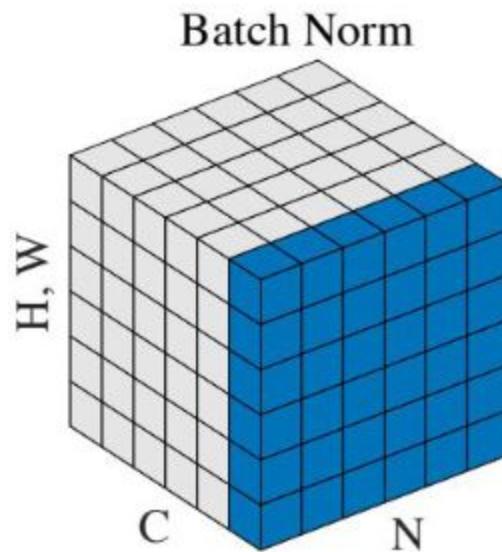
**Batch Normalization** for convolutional networks



**Instance Normalization** for convolutional networks  
Same behavior at train / test!



# Group Normalization



Wu and He, "Group Normalization", ECCV 2018



Michigan Tech

# Questions + Comments?



**Michigan Tech**

# Resources used

[http://cs231n.stanford.edu/slides/2023/lecture\\_2.pdf](http://cs231n.stanford.edu/slides/2023/lecture_2.pdf)

[http://cs231n.stanford.edu/slides/2023/lecture\\_3.pdf](http://cs231n.stanford.edu/slides/2023/lecture_3.pdf)

[http://cs231n.stanford.edu/slides/2023/lecture\\_4.pdf](http://cs231n.stanford.edu/slides/2023/lecture_4.pdf)



Michigan Tech