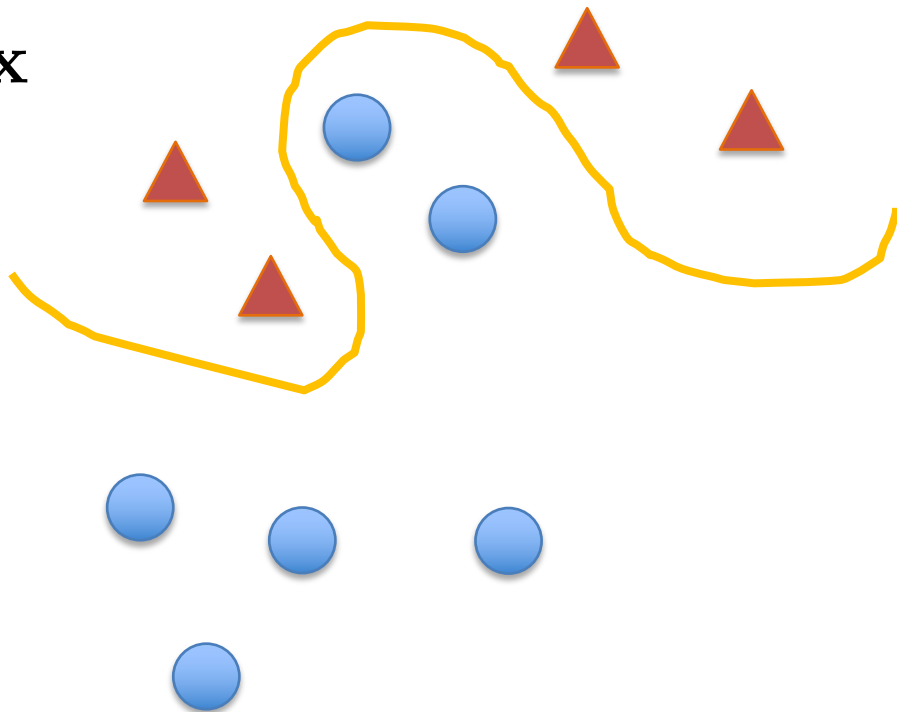
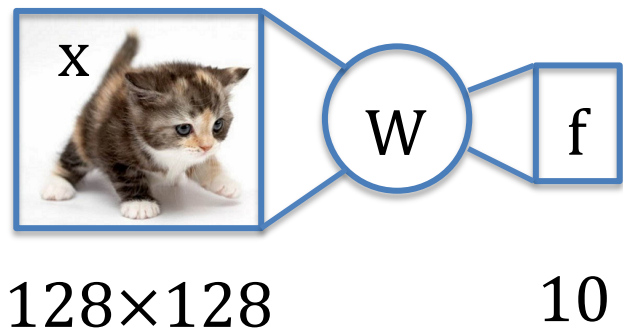


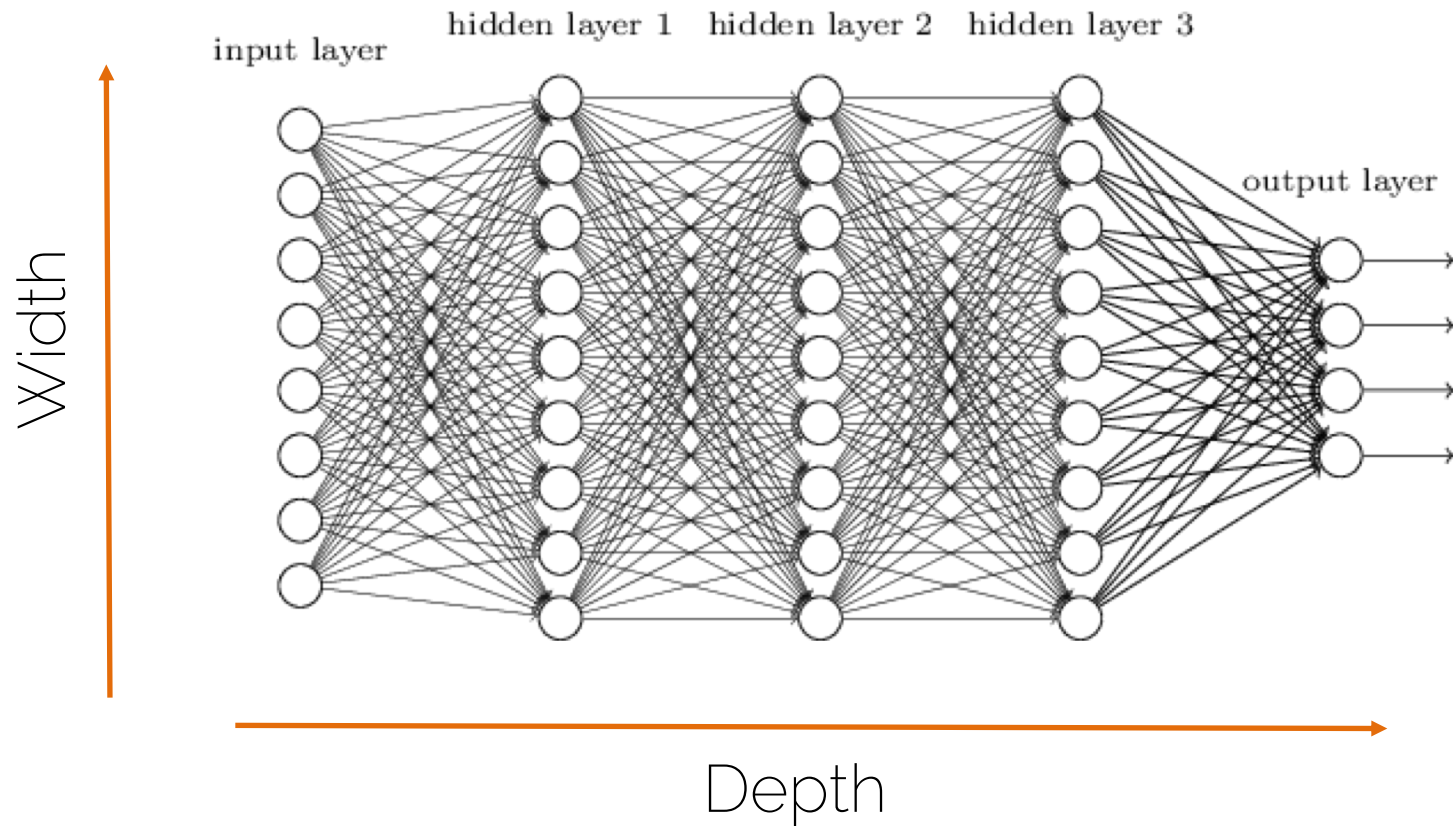
Lecture 7 Recap

Beyond linear

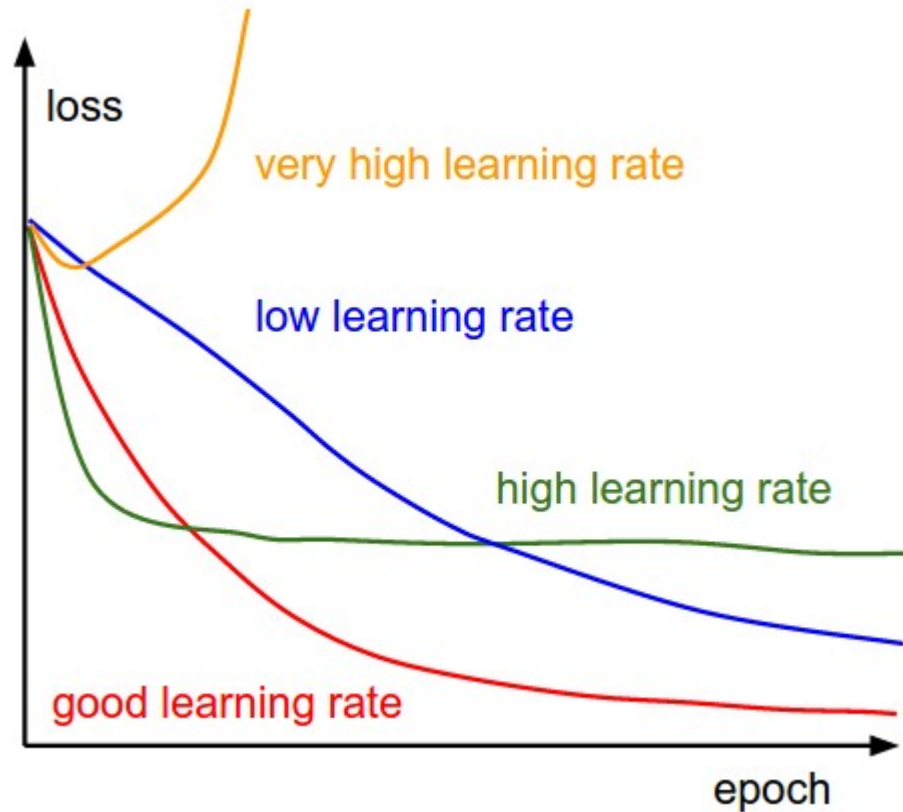
1-layer network: $f = \mathbf{W}\mathbf{x}$



Neural Network

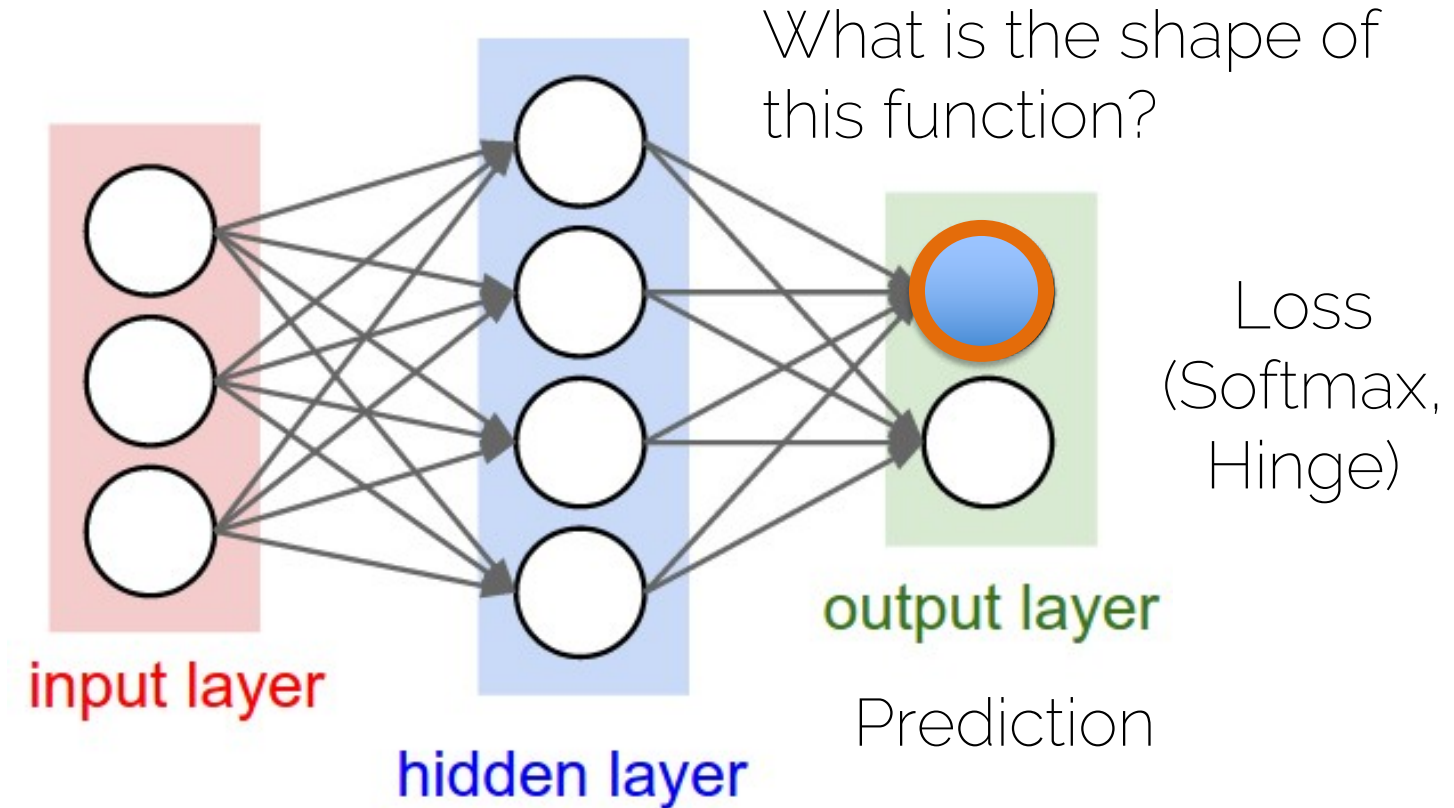


Optimization

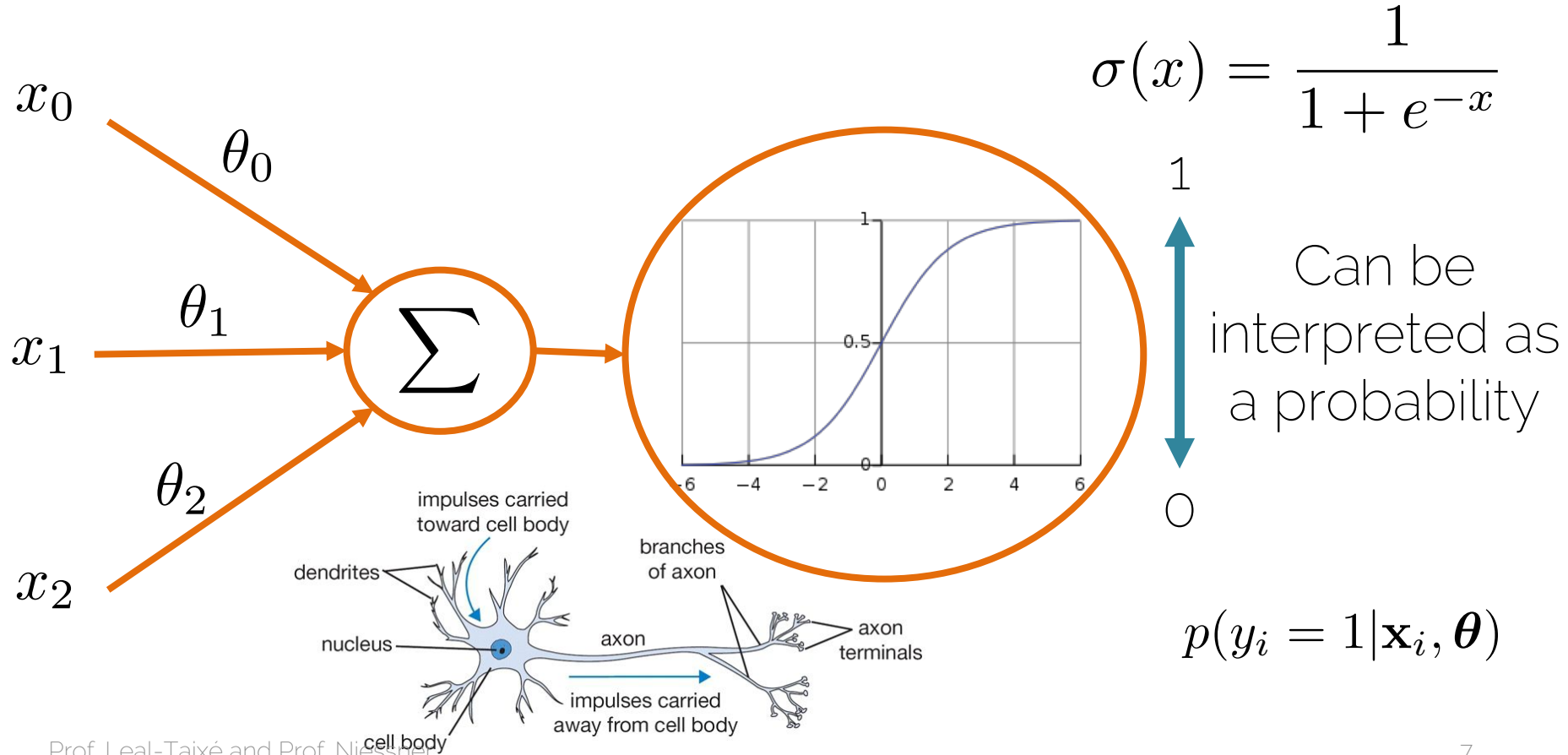


Loss functions

Neural networks

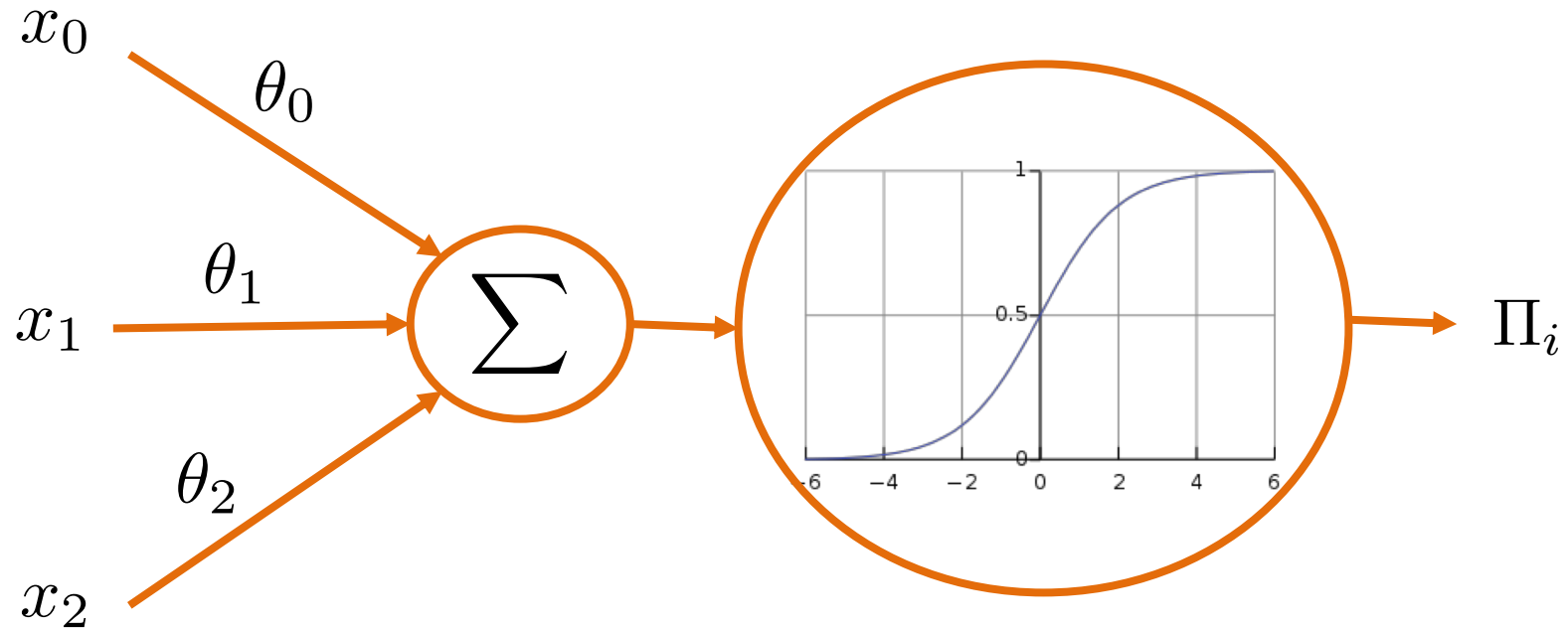


Sigmoid for binary predictions



Logitic regression

- Binary classification



Logistic regression

- Loss function

$$\mathcal{L}(\Pi_i, y_i) = y_i \log \Pi_i + (1 - y_i) \log(1 - \Pi_i)$$

One training sample

- Cost function

$$C(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n y_i \log \Pi_i + (1 - y_i) \log(1 - \Pi_i)$$

Minimization

$$\sigma(\mathbf{x}_i \boldsymbol{\theta})$$

Softmax regression

- Cost function for the binary case

$$C(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n y_i \log \Pi_i + (1 - y_i) \log(1 - \Pi_i)$$

- Extension to multiple classes

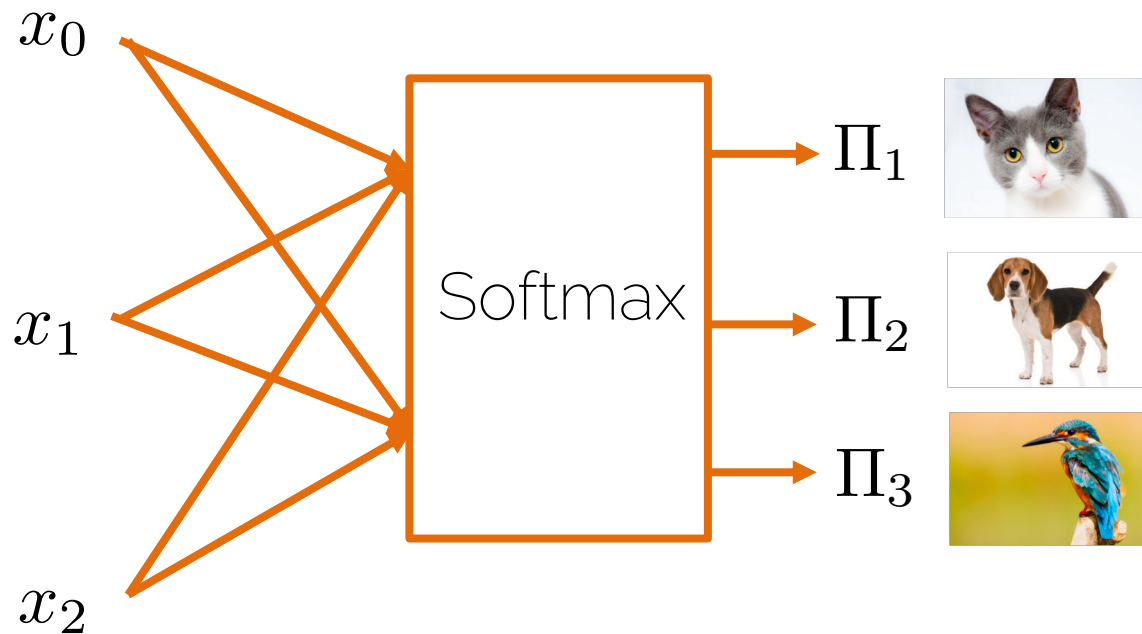
$$C(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n \sum_{c=1}^M y_{i,c} \log p_{i,c}$$

Probability given by our sigmoid function

Binary indicator whether c is the label for image i

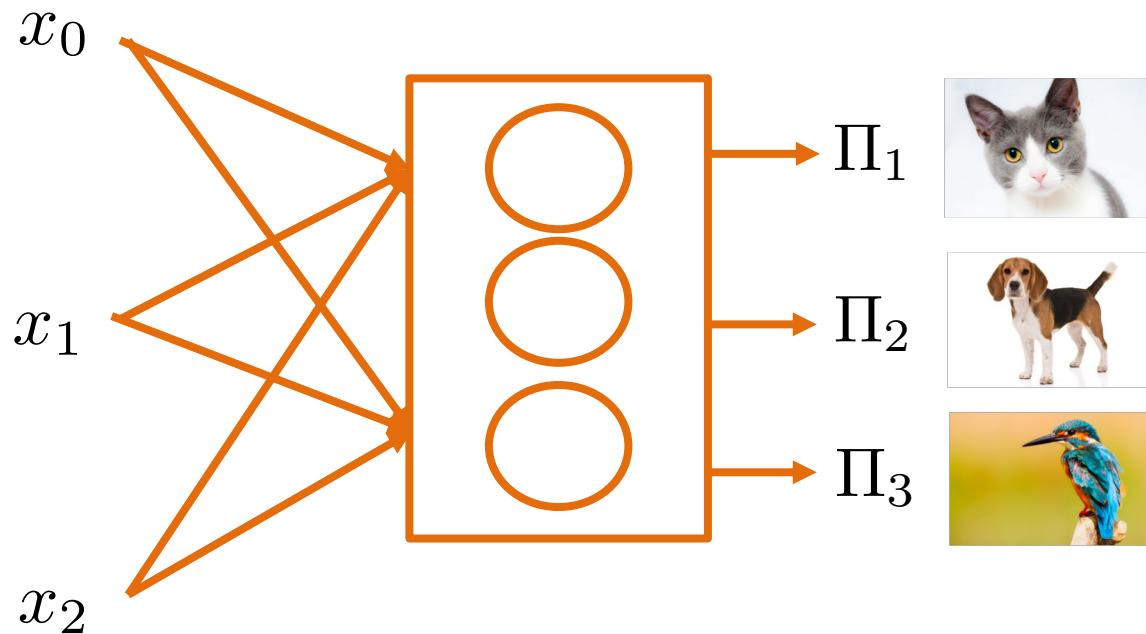
Softmax formulation

- What if we have multiple classes?



Softmax formulation

- Three neurons in the output layer for three classes



Softmax formulation

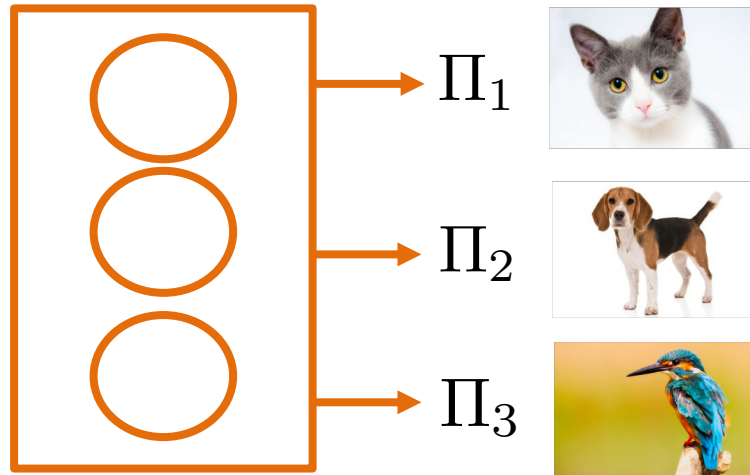
- What if we have multiple classes?

$$C(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n \sum_{c=1}^M y_{i,c} \log p_{i,c}$$

- You can no longer assign $p_{i,c}$ to Π_i as in the binary case, because all outputs need to sum to 1

$$\sum_c \Pi_{i,c}$$

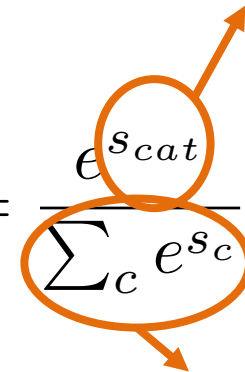
Softmax formulation



$$p(cat|\mathbf{X}_i) = \frac{e^{s_{cat}}}{\sum_c e^{s_c}}$$

$$p(dog|\mathbf{X}_i)$$

$$p(bird|\mathbf{X}_i)$$



Score for class cat given by all the layers of the network

Normalization

- Softmax takes M inputs (Scores) and outputs M probabilities (M is the number of classes)

Loss functions

- Softmax loss function

$$L_i = -\log \left(\frac{e^{s_{y_i}}}{\sum_k e^{s_k}} \right)$$

Evaluate the ground truth score for the image

Comes from Maximum Likelihood Estimate

- Hinge Loss (derived from the Multiclass SVM loss)

$$L_i = \sum_{k \neq y_i} \max(0, s_k - s_{y_i} + 1)$$

Loss functions

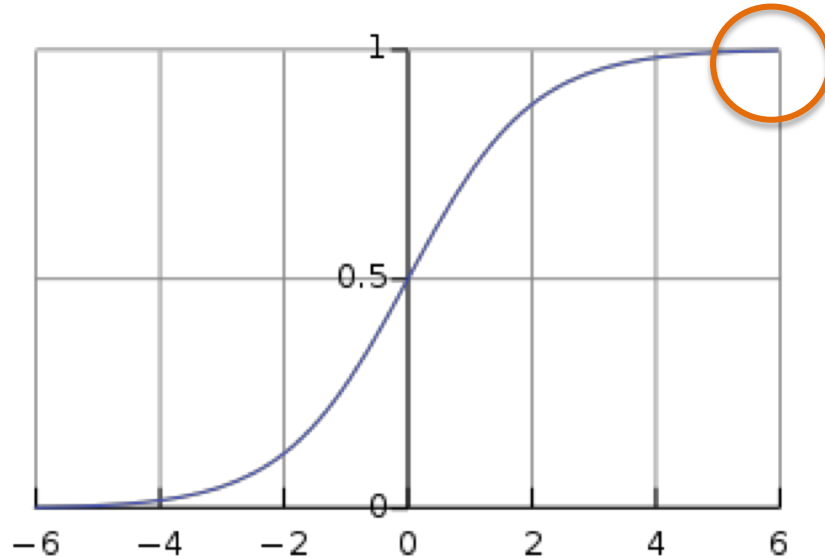
- Softmax loss function
 - Optimizes until the loss is zero
- Hinge Loss (derived from the Multiclass SVM loss)
 - Saturates whenever it has learned a class “well enough”

Activation functions

Sigmoid

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

Forward



$$x = 6$$

✗ Saturated neurons kill the gradient flow

$$\frac{\partial L}{\partial x} = \frac{\partial \sigma}{\partial x} \frac{\partial L}{\partial \sigma}$$

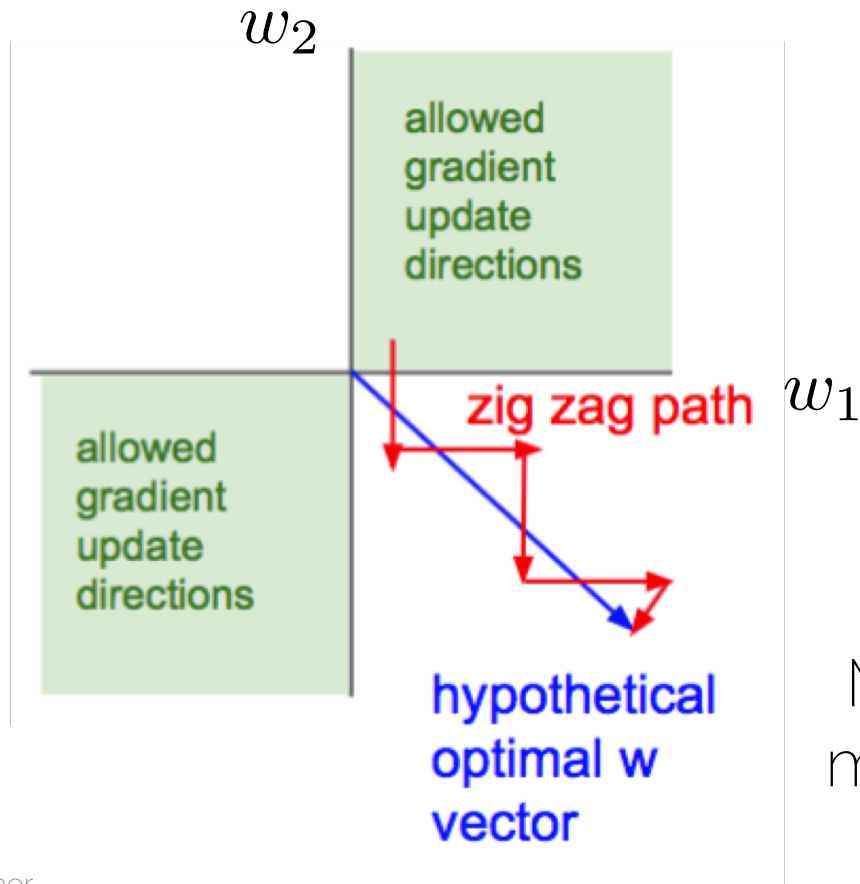


$$\frac{\partial \sigma}{\partial x}$$



$$\frac{\partial L}{\partial \sigma}$$

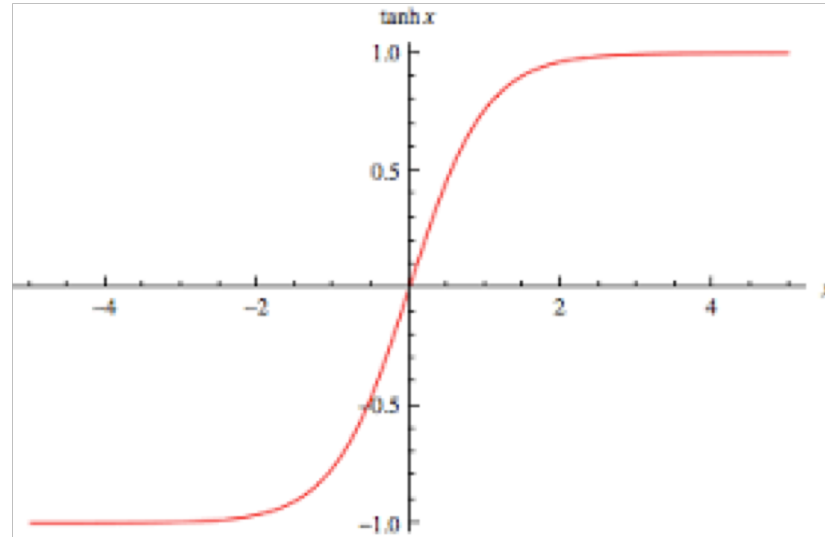
Problem of positive output



More on zero-mean data later

tanh

✗ Still saturates



✗ Still saturates

✓ Zero-centered

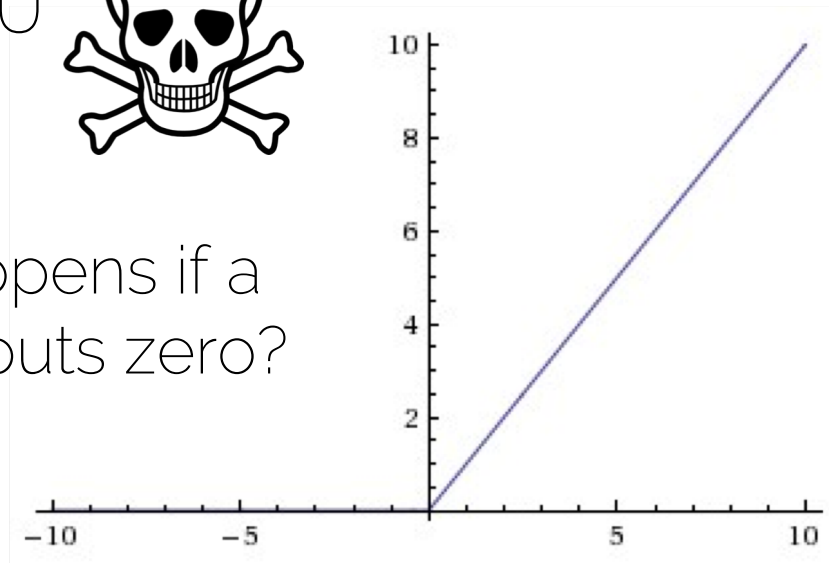
Rectified Linear Units (ReLU)



Dead ReLU



What happens if a ReLU outputs zero?



Large and consistent gradients

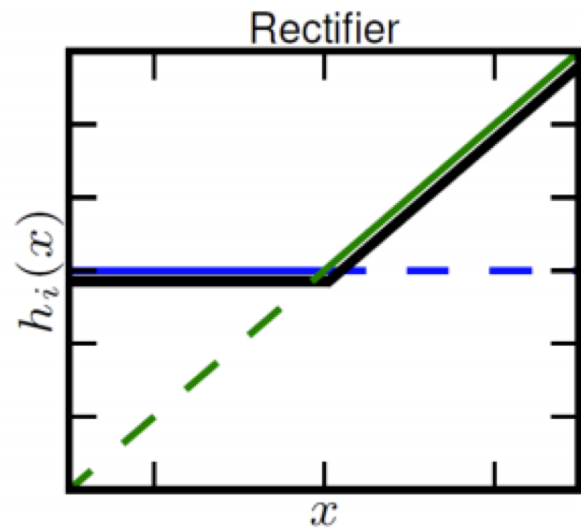


Fast convergence



Does not saturate

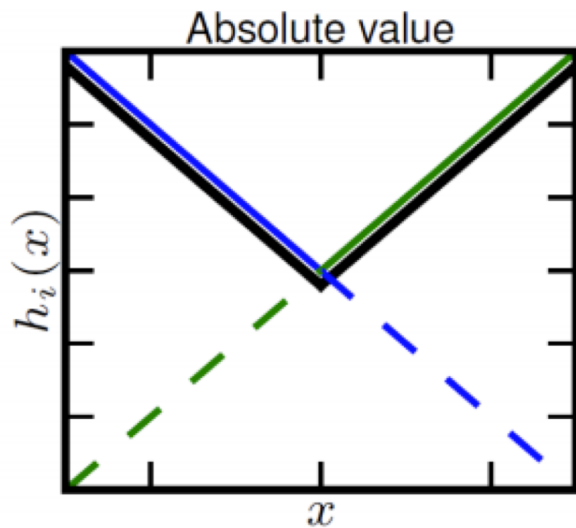
Maxout units



$k=2$



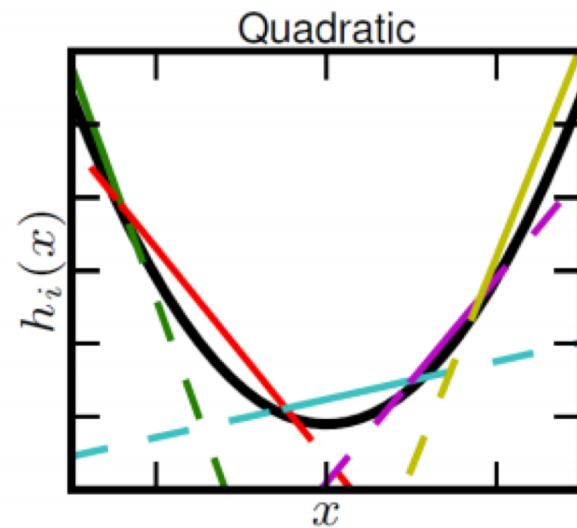
Generalization
of ReLUs



$k=2$



Linear
regimes



$k=5$



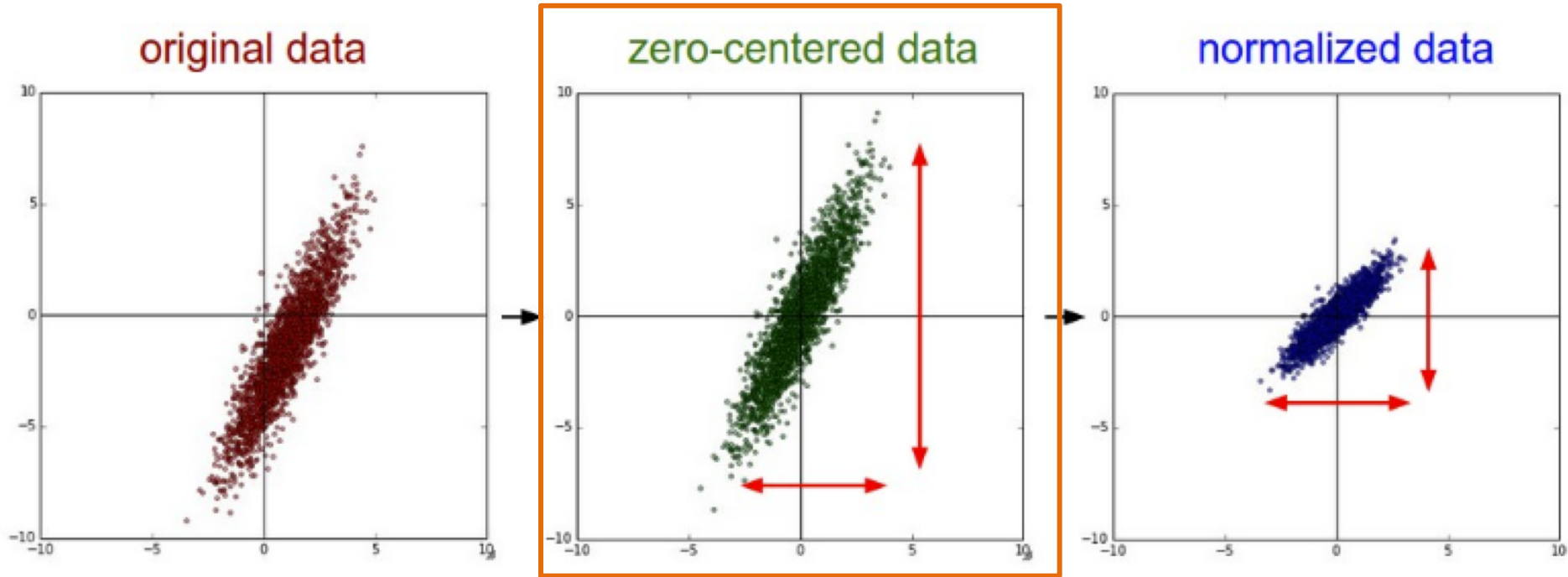
Does not
die

Does not
saturate



Increase of the number of parameters

Data pre-processing

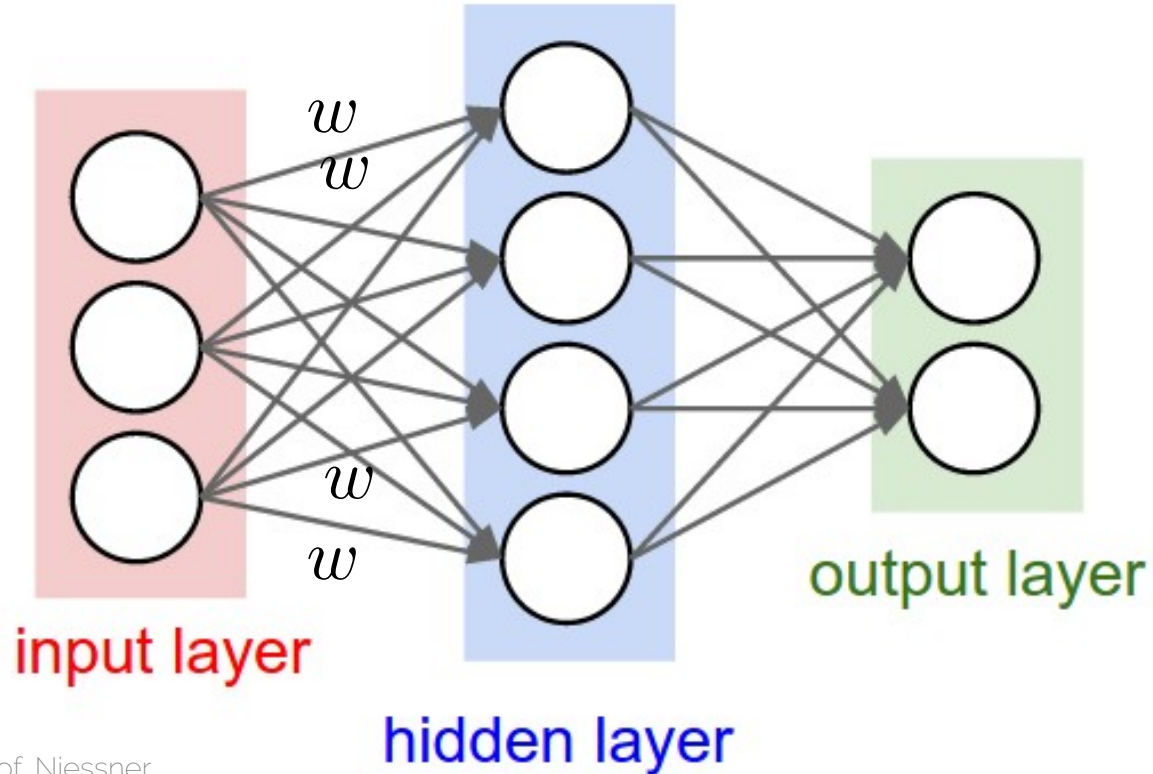


For images subtract the mean image (AlexNet) or per-channel mean (VGG-Net)

Weight initialization

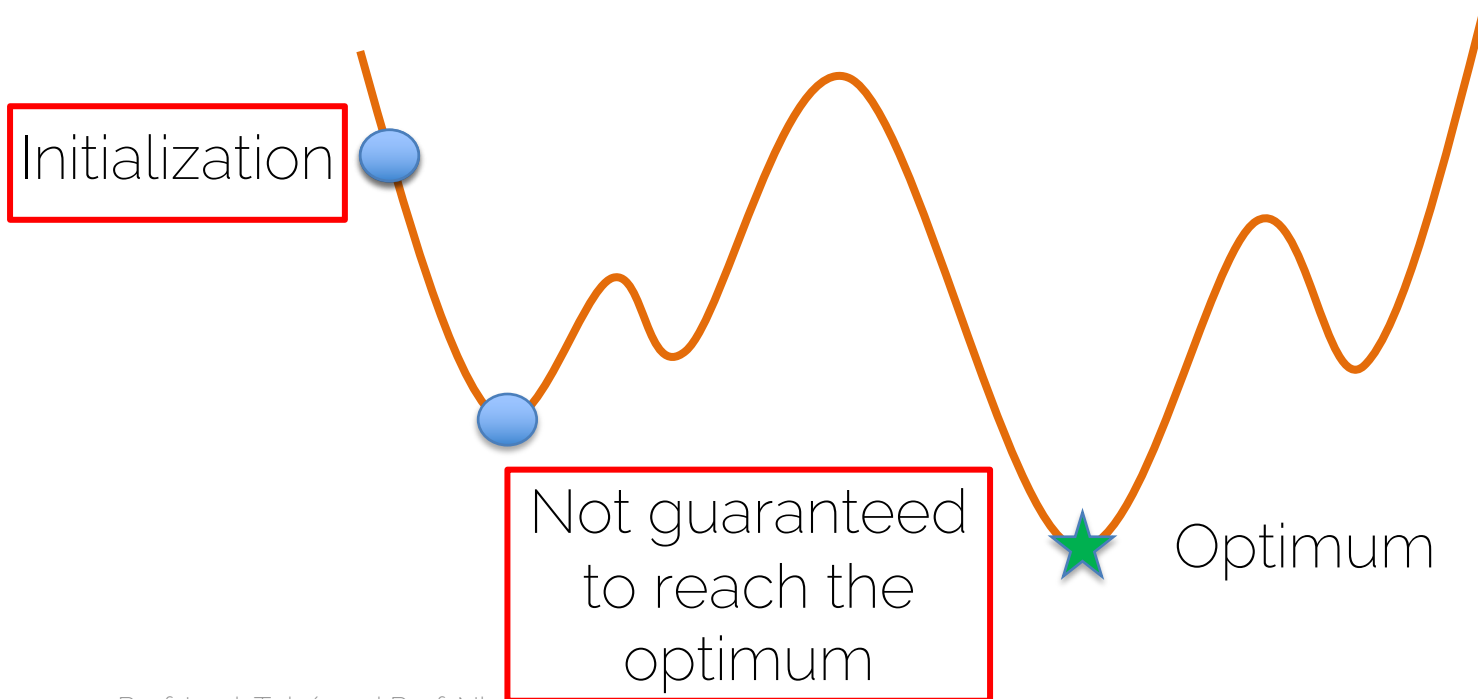
How do I start?

Forward



Initialization is extremely important

$$\mathbf{x}^* = \arg \min f(\mathbf{x})$$



How do I start?

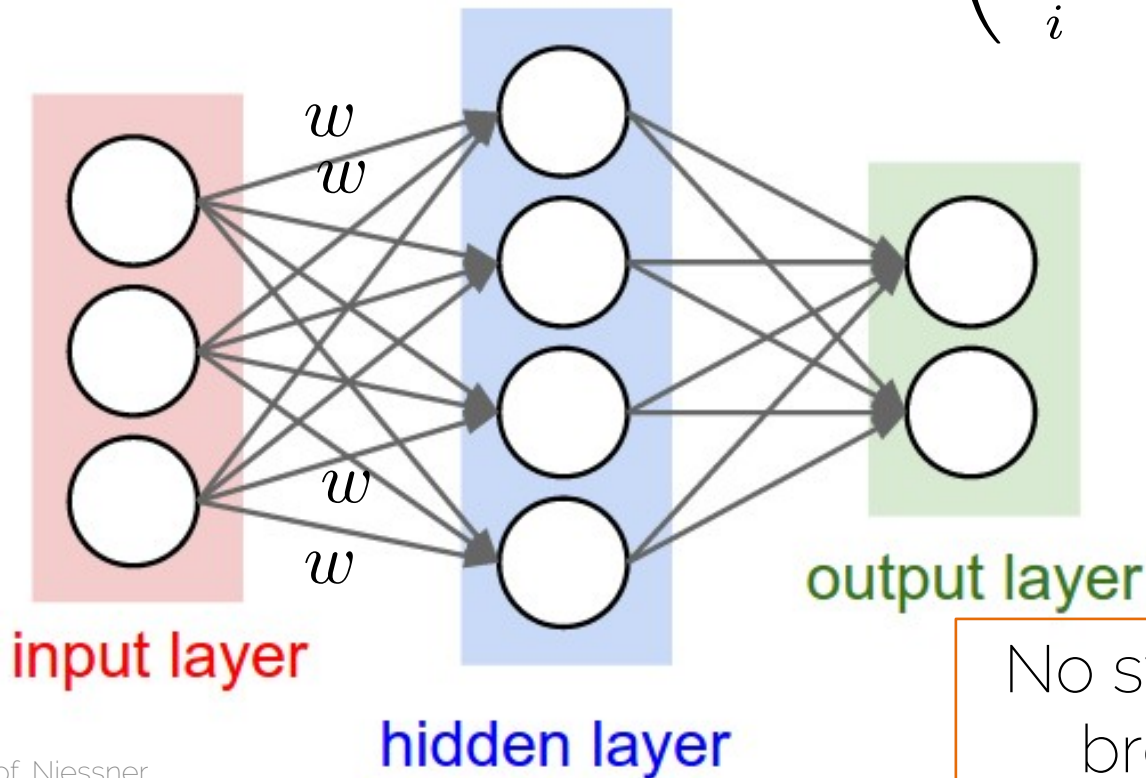
Forward



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

$w = 0$

What happens to the gradients?



No symmetry breaking₂₇

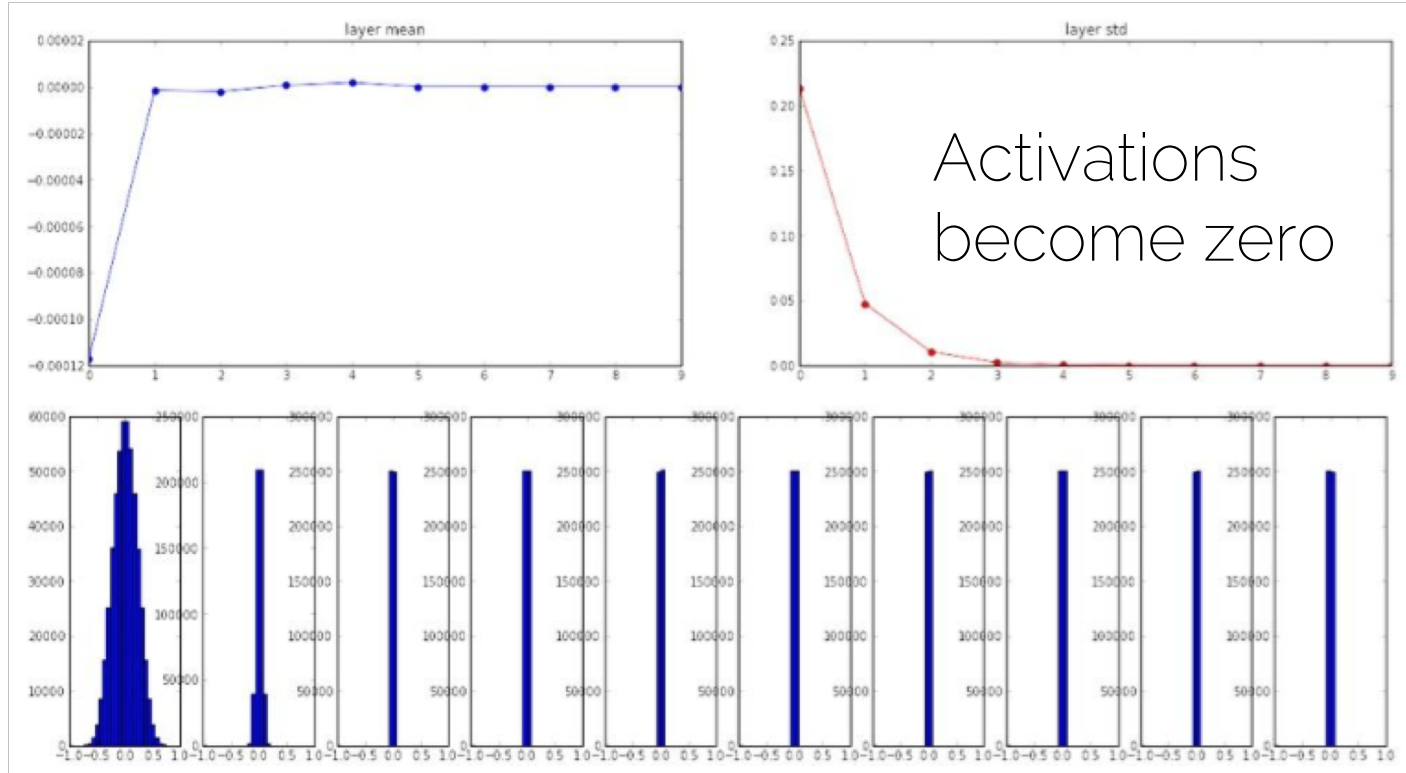
All weights to zero

- Elaborate: the hidden units are all going to compute the same function, gradients are going to be the same

Small random numbers

- Gaussian with zero mean and standard deviation 0.01
- Let us see what happens:
 - Network with 10 layers with 500 neurons each
 - Tanh as activation functions
 - Input unit Gaussian data

Small random numbers

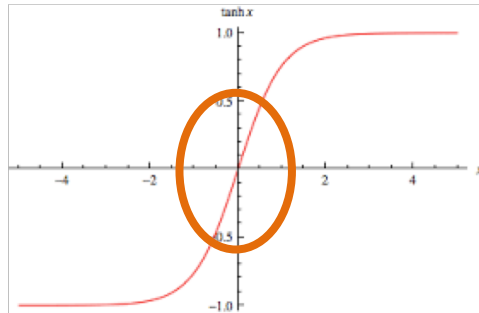
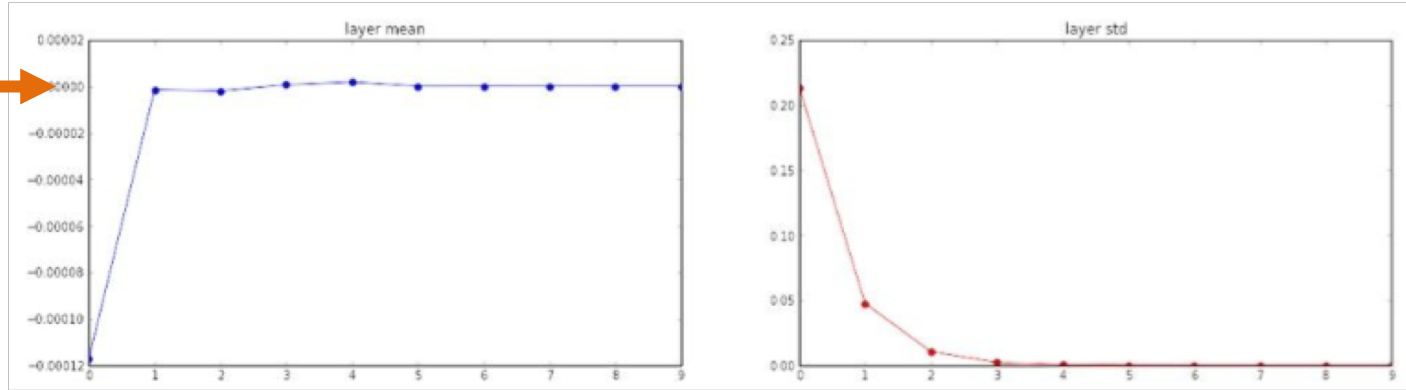


Input

Last layer

Forward

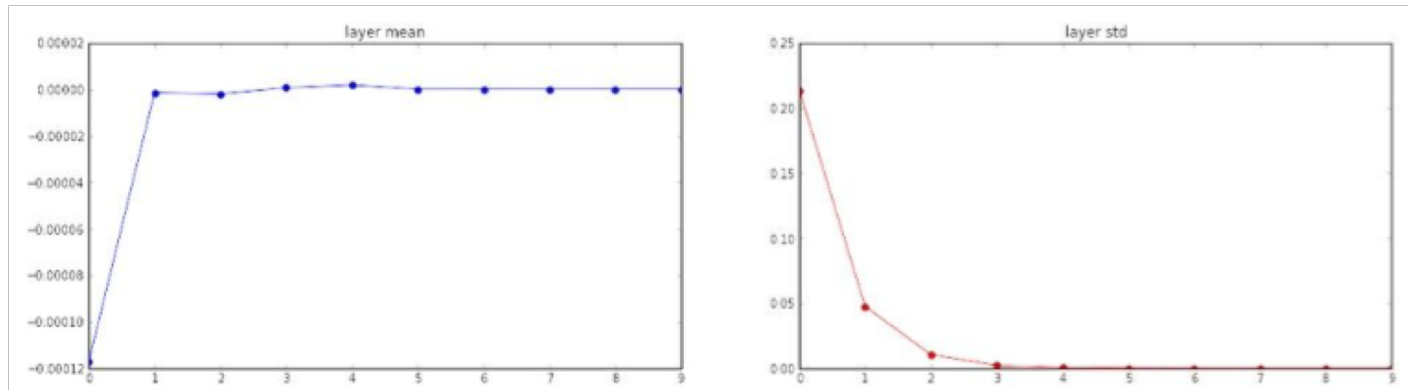
Small random numbers



$$f \left(\sum_i^{\text{small}} w_i x_i + b \right)$$

Forward

Small random numbers



1. Activation
function
gradient is ok

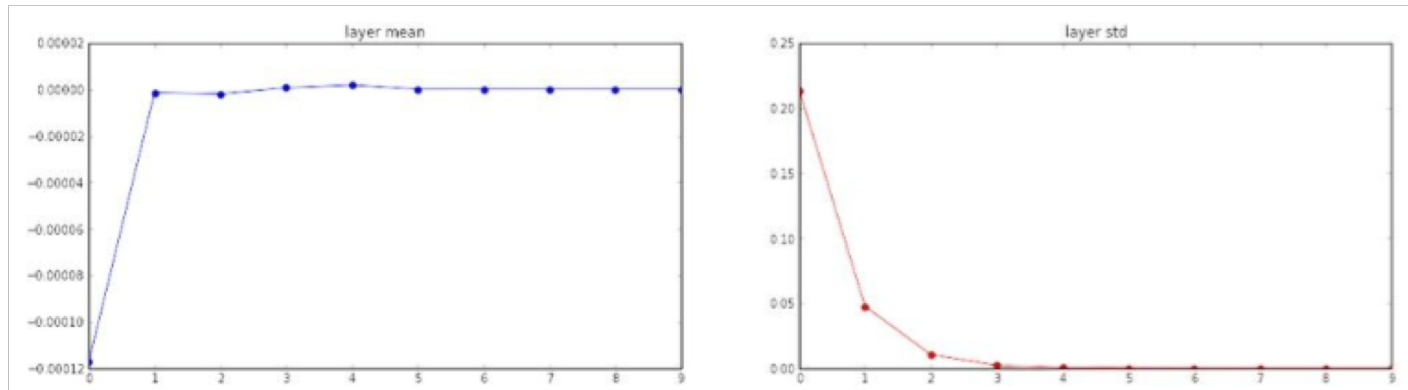
$$f \left(\sum_i w_i x_i + b \right)$$

An orange bracket is drawn above the summation term $\sum_i w_i x_i$. An orange circle is drawn around the variable x_i . An orange arrow points from the bracket to the text '2. Compute the gradients wrt the weights'.

2. Compute the
gradients wrt
the weights

Backward

Small random numbers



1. Activation
function
gradient is ok

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

An orange bracket is drawn above the summation term $\sum_i w_i x_i$. An orange circle is drawn around the variable x_i . An orange arrow points from the circle around x_i towards the text '2. Compute the gradients wrt the weights'.

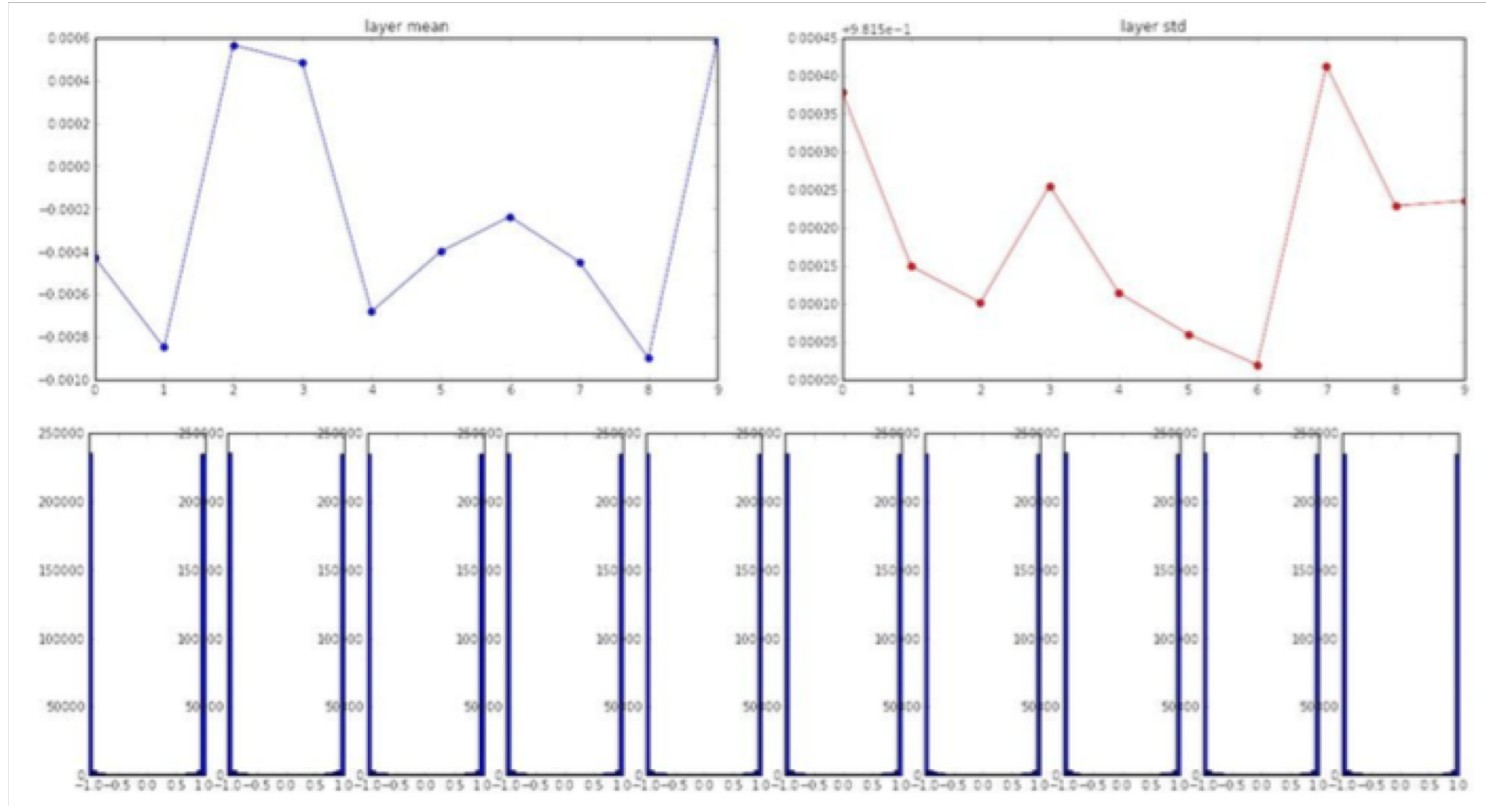
2. Compute the
gradients wrt
the weights

Gradients vanish

Big random numbers

- Gaussian with zero mean and standard deviation 1
- Let us see what happens:
 - Network with 10 layers with 500 neurons each
 - Tanh as activation functions
 - Input unit Gaussian data

Big random numbers



Everything
is saturated

How to solve this?

- Working on the initialization
- Working on the output generated by each layer

Xavier initialization

- Gaussian with zero mean, but what standard deviation?

$$\text{Var}(s) = \text{Var}\left(\sum_i^n w_i x_i\right) = \sum_i^n \text{Var}(w_i x_i)$$


Xavier initialization

- Gaussian with zero mean, but what standard deviation?

$$\begin{aligned}\text{Var}(s) &= \text{Var}\left(\sum_i^n w_i x_i\right) = \sum_i^n \text{Var}(w_i x_i) \quad \text{Independent} \\ &= \sum_i^n \left[\cancel{E(w_i)^2} \text{Var}(x_i) + \cancel{E(x_i)^2} \text{Var}(w_i) + \text{Var}(x_i) \text{Var}(w_i) \right] \\ &\quad \text{Zero mean}\end{aligned}$$

Xavier initialization

- Gaussian with zero mean, but what standard deviation?

$$\begin{aligned}\text{Var}(s) &= \text{Var}\left(\sum_i^n w_i x_i\right) = \sum_i^n \text{Var}(w_i x_i) \\ &= \sum_i^n [E(w_i)]^2 \text{Var}(x_i) + E[(x_i)]^2 \text{Var}(w_i) + \text{Var}(x_i) \text{Var}(w_i) \\ &= \sum_i^n \text{Var}(x_i) \text{Var}(w_i) = (n \text{Var}(w)) \text{Var}(x)\end{aligned}$$


Identically distributed

Xavier initialization

- Gaussian with zero mean, but what standard deviation?

$$\begin{aligned}\text{Var}(s) &= \text{Var}\left(\sum_i^n w_i x_i\right) = \sum_i^n \text{Var}(w_i x_i) \\ &= \sum_i^n [E(w_i)]^2 \text{Var}(x_i) + E[(x_i)]^2 \text{Var}(w_i) + \text{Var}(x_i) \text{Var}(w_i) \\ &= \sum_i^n \text{Var}(x_i) \text{Var}(w_i) = (n \text{Var}(w)) \text{Var}(x)\end{aligned}$$

Variance gets multiplied by the number of inputs

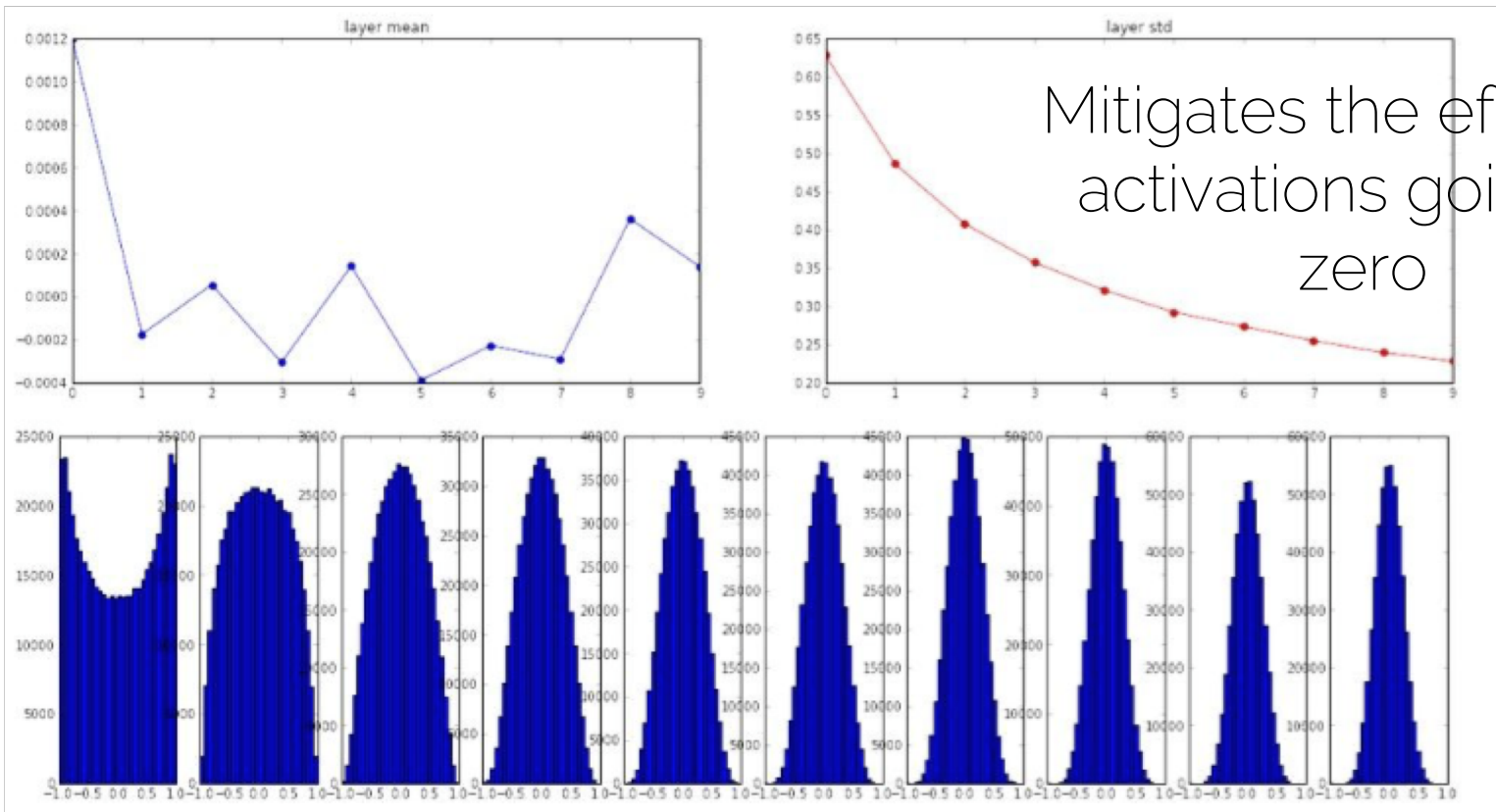
Xavier initialization

- How to ensure the variance of the output is the same as the input?

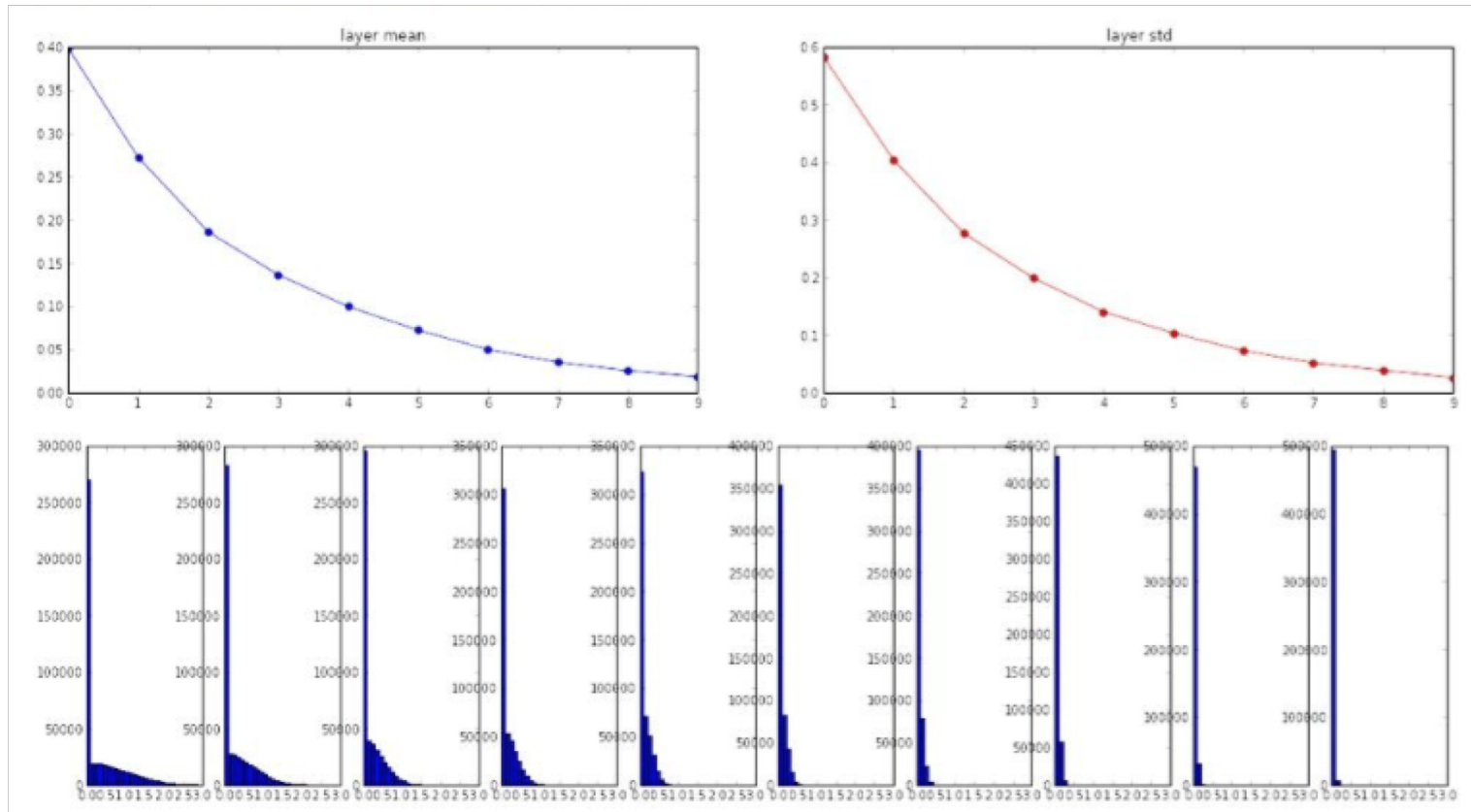
$$\frac{(n \text{Var}(w)) \text{Var}(x)}{1}$$

$$\text{Var}(w) = \frac{1}{n}$$

Xavier initialization

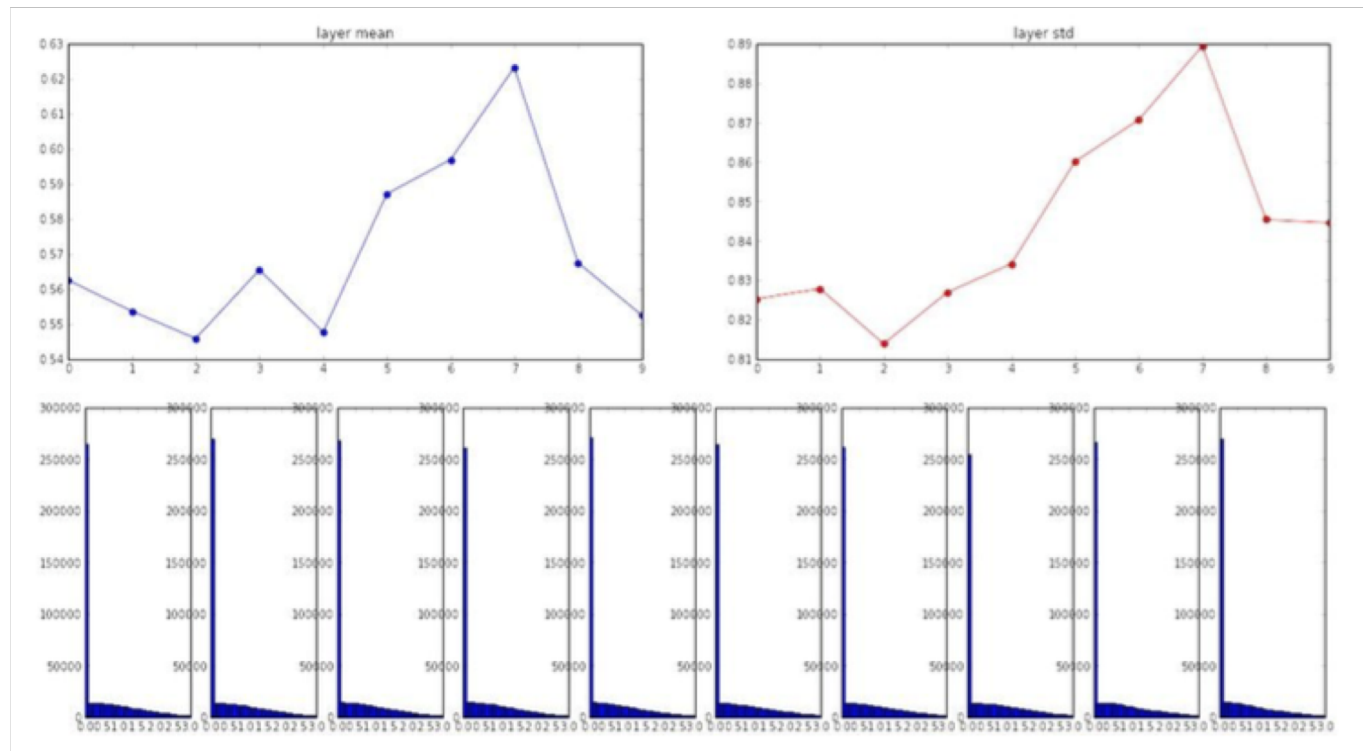


Xavier initialization with ReLU



ReLU kills half of the data

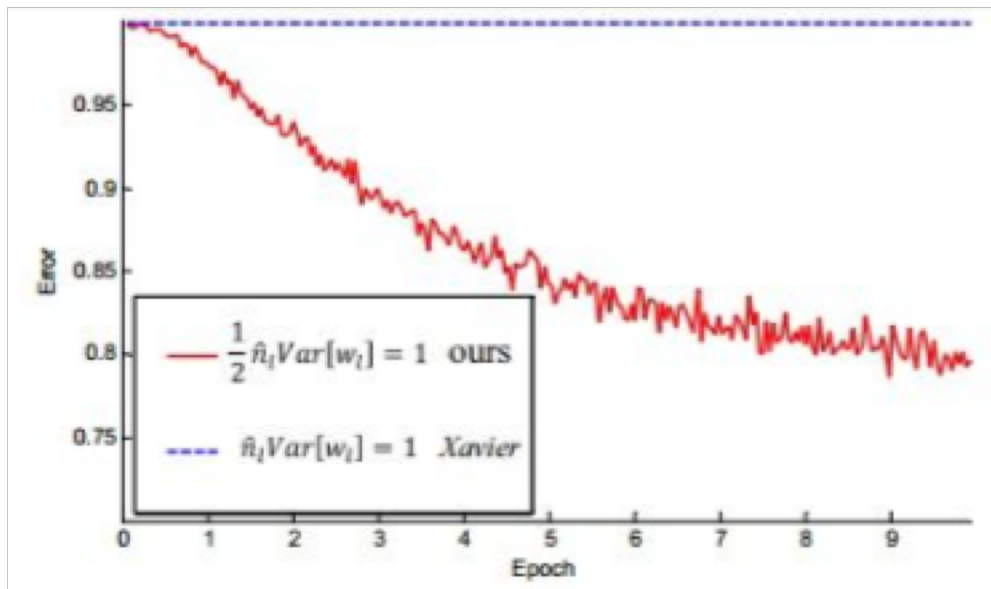
$$\text{Var}(w) = \frac{2}{n}$$



ReLU kills half of the data

$$\text{Var}(w) = \frac{2}{n}$$

It makes a huge difference!



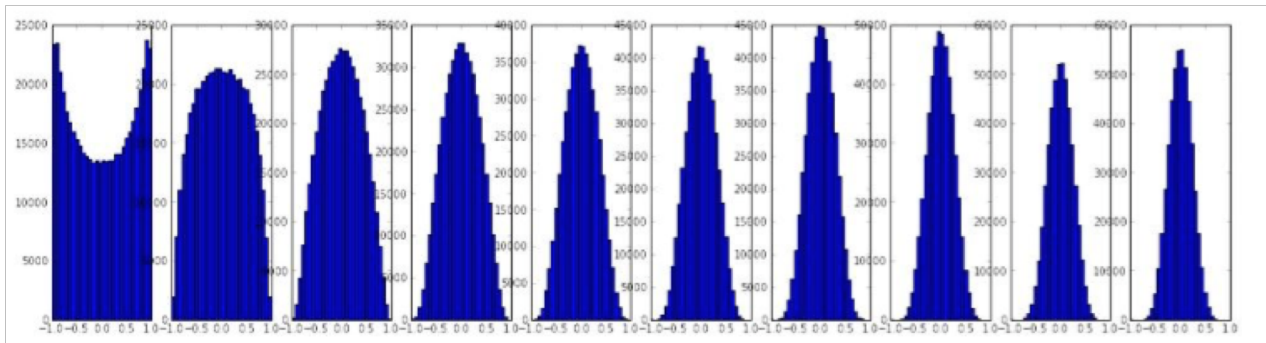
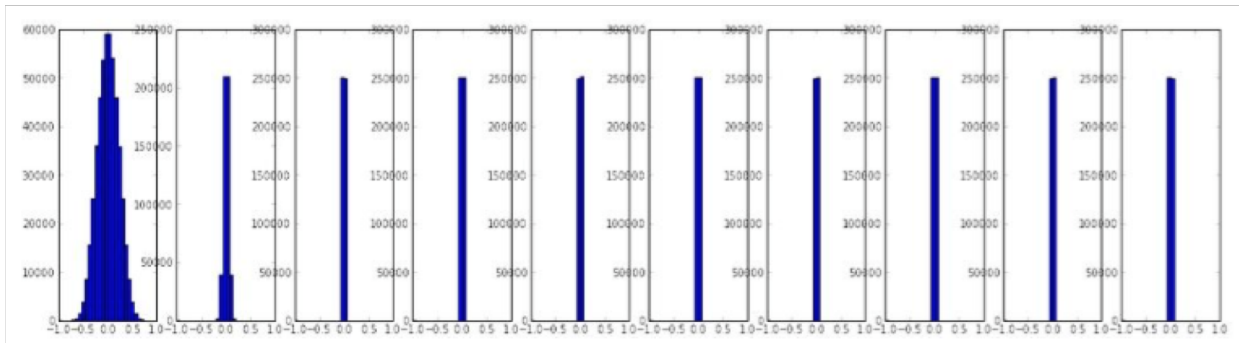
Tips and tricks

- Use ReLU and Xavier/2 initialization

Batch normalization

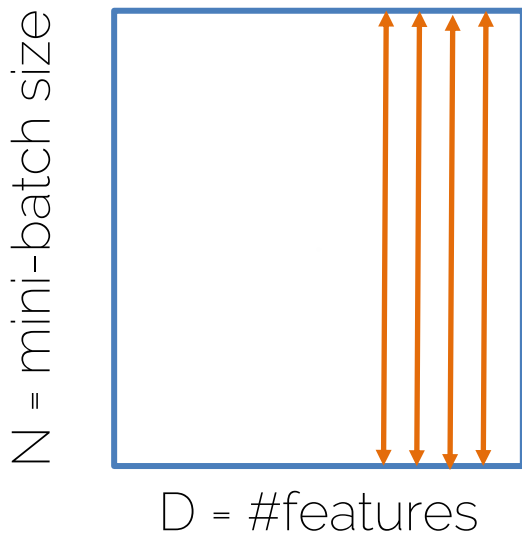
Our goal

- All we want is that our activations do not die out



Batch normalization

- Wish: unit Gaussian activations (in our example)
- Solution: let's do it



Mean of your mini-batch examples over feature k

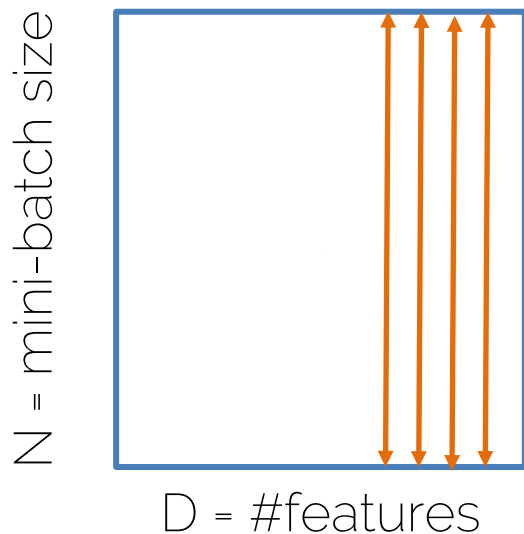
$$\hat{x}^{(k)} = \frac{x^{(k)} - \text{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

An orange arrow points from the text "Mean of your mini-batch examples over feature k" to the term $\text{E}[x^{(k)}]$ in the numerator of the equation.

Batch normalization

- In each dimension of the features, you have a unit gaussian (in our example)

Mean of your mini-batch examples over feature k



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

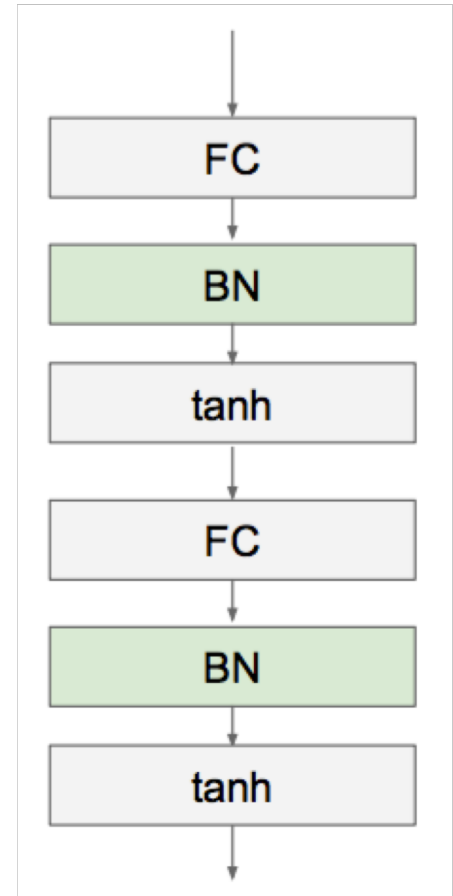
An orange arrow points to the $x^{(k)}$ term in the numerator of the equation.

Batch normalization

- In each dimension of the features, you have a unit gaussian (in our example)
- For NN in general → BN normalizes the mean and variance of the inputs to your activation functions

BN layer

- A layer to be applied after Fully Connected (or Convolutional) layers and **before** non-linear activation functions
- Is it a good idea to have all unit Gaussians before tanh? This normalization might not be the best for the network!



Batch normalization

- 1. Normalize

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

Differentiable function so
we can backprop through
it...

- 2. Allow the network to change the range

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

These parameters will be
optimized during backprop

Batch normalization

- 1. Normalize

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

- 2. Allow the network to change the range

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

backprop

The network *can* learn to undo the normalization

$$\gamma^{(k)} = \sqrt{\text{Var}[x^{(k)}]}$$

$$\beta^{(k)} = \mathbf{E}[x^{(k)}]$$

Batch normalization

- Is it ok to treat dimensions separately? Shown empirically that even if features are not decorrelated, convergence is still faster with this method
- You can set all biases of the layers before BN to zero, because they will be cancelled out by BN anyway

BN: train vs test time

- Train time: mean and variance is taken over the mini-batch

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

- Test-time: what happens if we can just process one image at a time?
 - No chance to compute a meaningful mean and variance

BN: train vs test time

Training

- Compute mean and variance from mini-batch 1
- Compute mean and variance from mini-batch 2
- Compute mean and variance from mini-batch 3

Testing

- Compute mean and variance by running an exponentially weighted averaged across training mini-batches

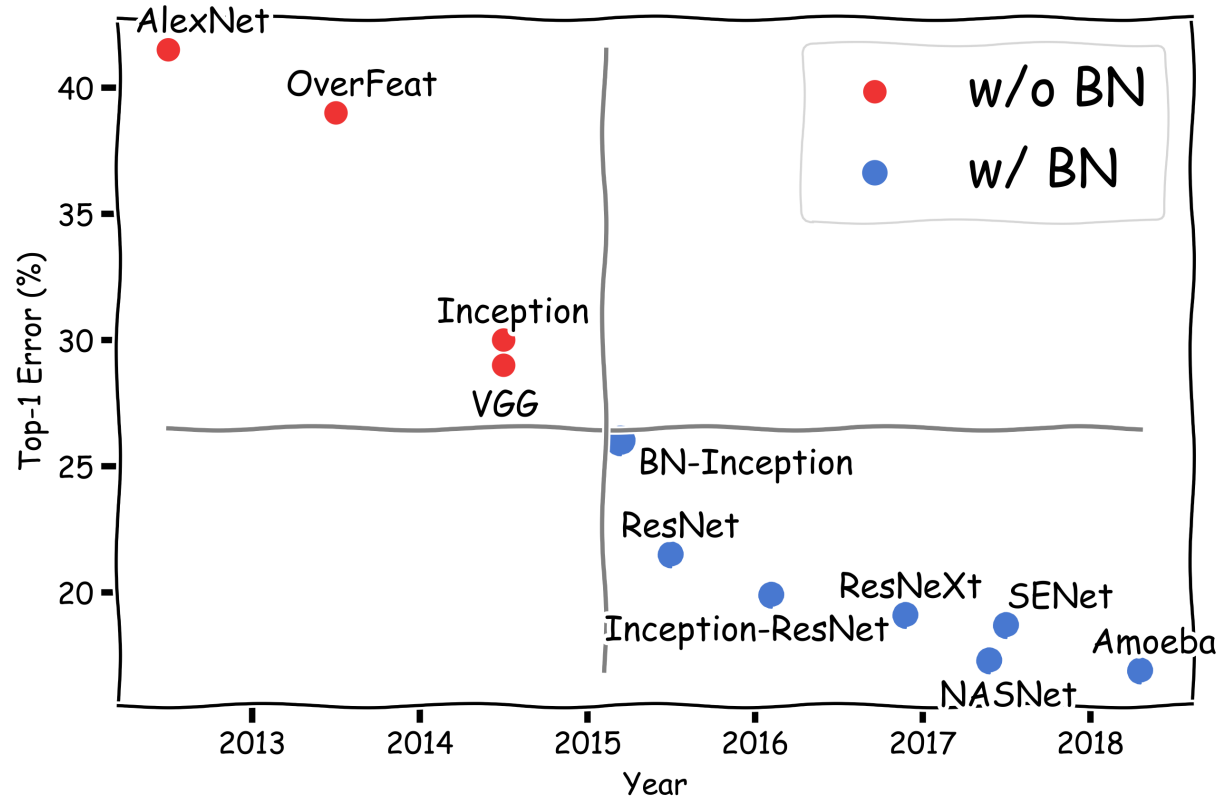
μ_{test}

σ_{test}^2

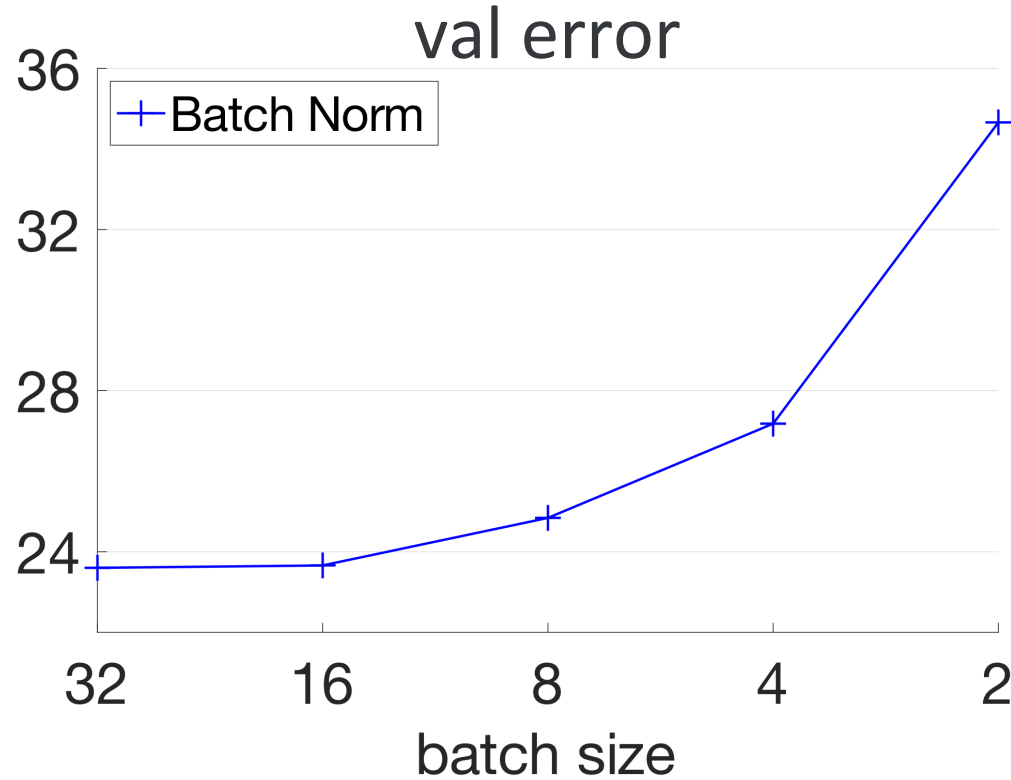
BN: what do you get?

- Very deep nets are much easier to train → more stable gradients
- A much larger range of hyperparameters works similarly when using BN

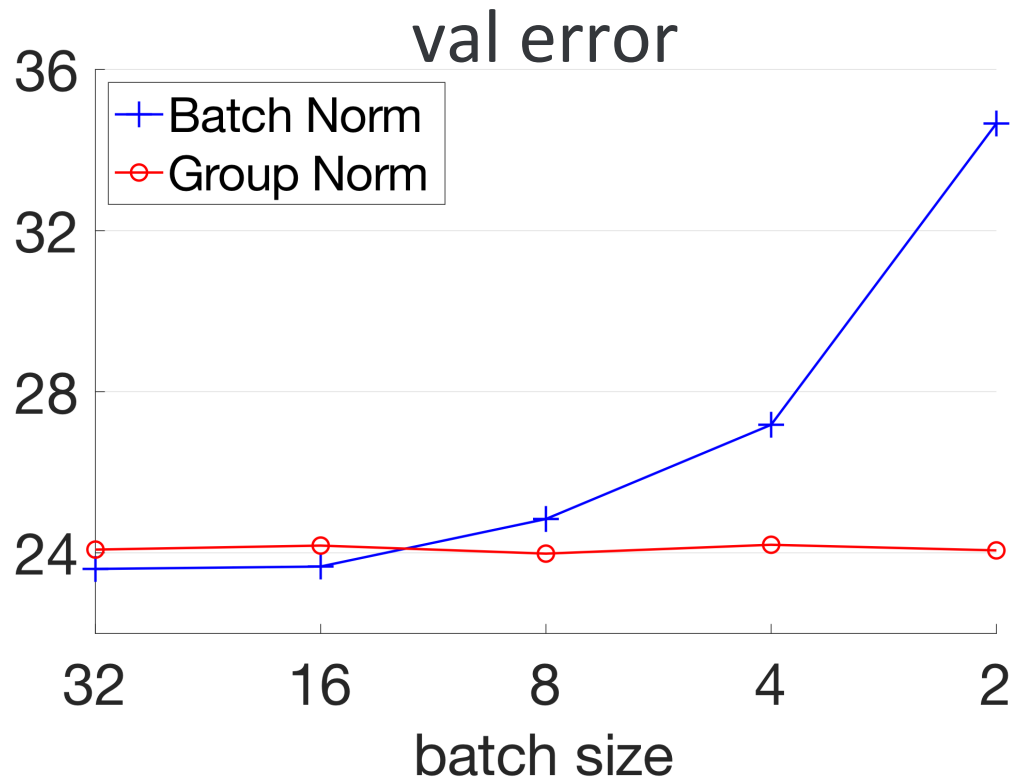
BN: a milestone



BN: drawbacks

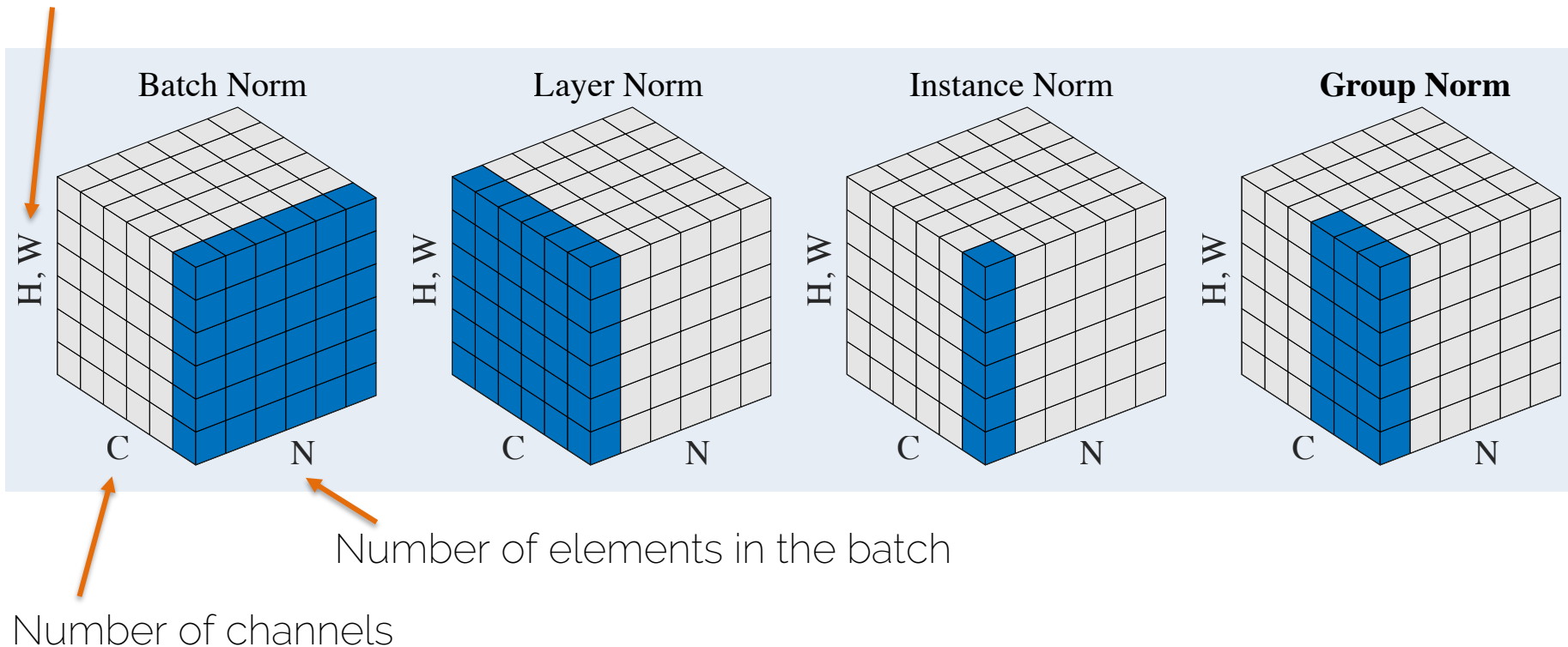


Other normalizations



Other normalizations

Image size



Regularization

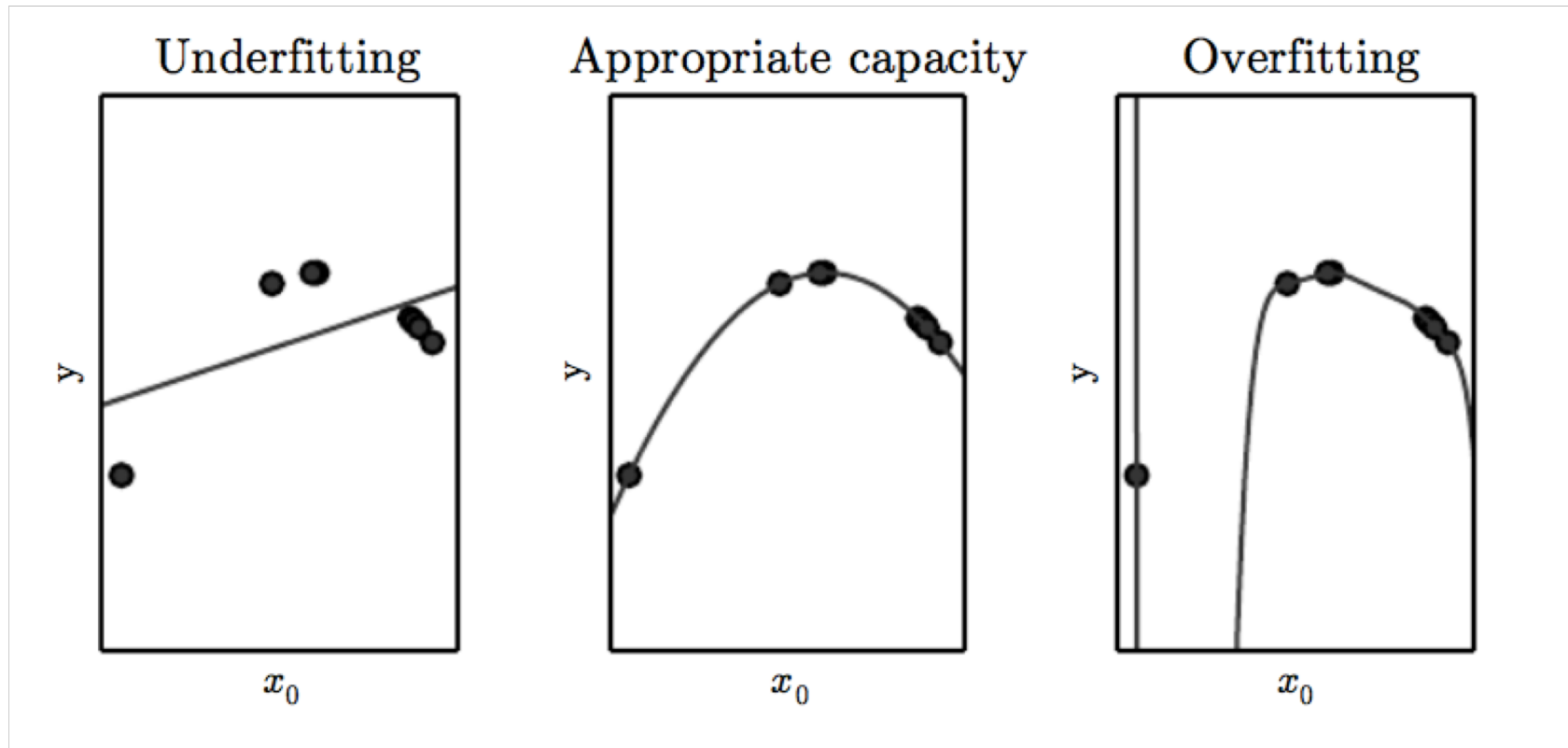
Regularization

- Any strategy that aims to

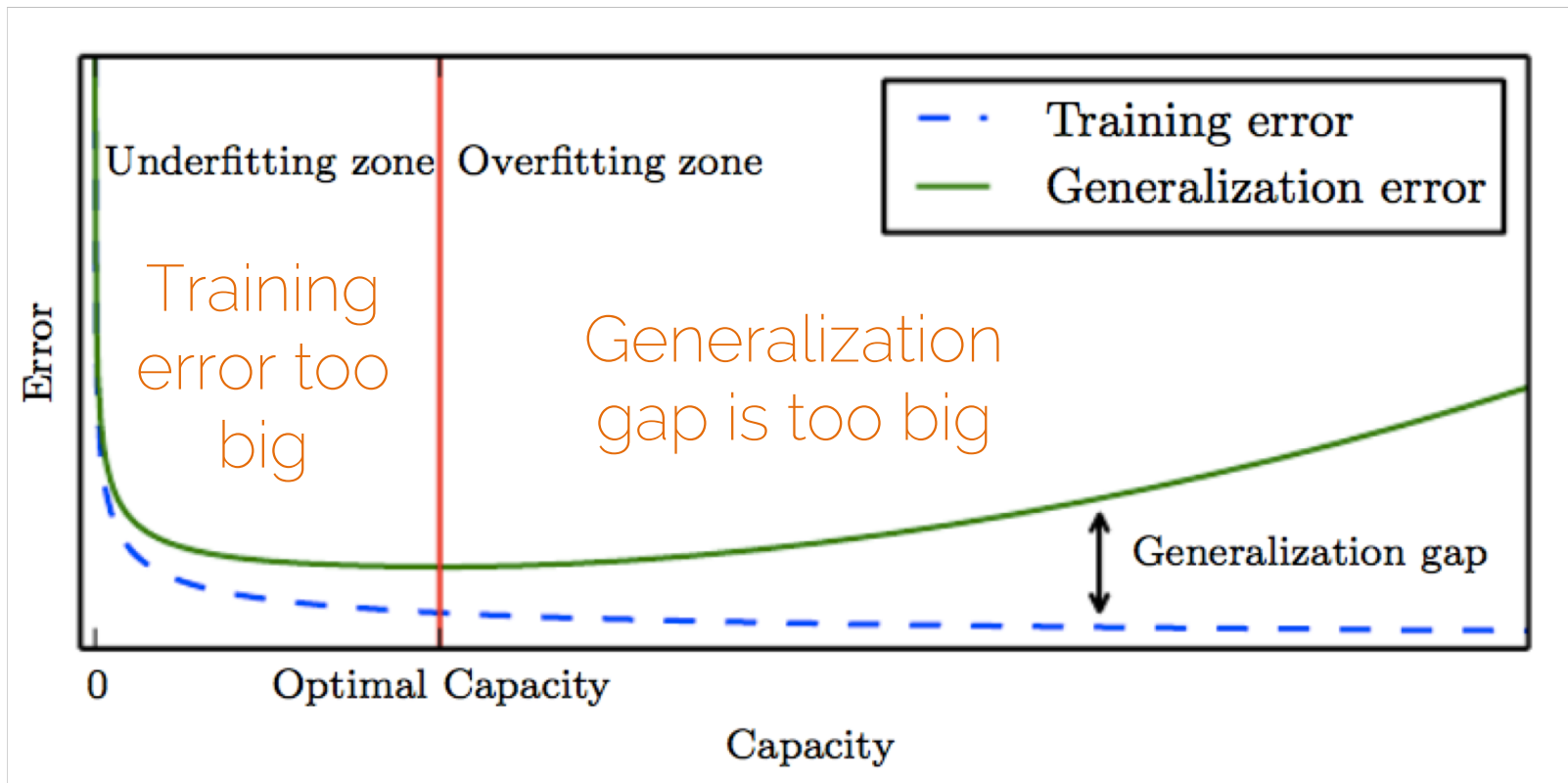
Lower
validation error

Increasing
training error

Overfitting and underfitting



Overfitting and underfitting



Weight decay

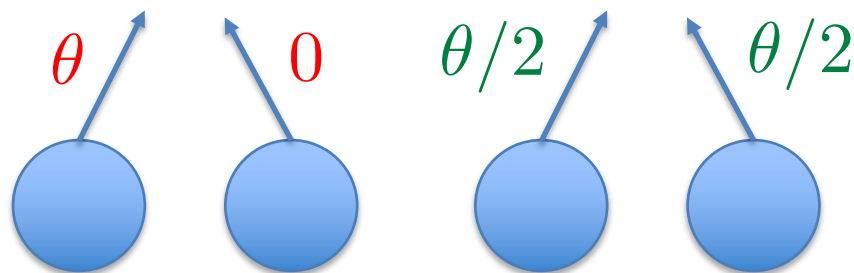
- L^2 regularization

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \epsilon \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_k, \mathbf{x}^i, \mathbf{y}^i) - \lambda \boldsymbol{\theta}_k^T \boldsymbol{\theta}_k$$

Learning rate

Gradient

- Penalizes large weights
- Improves generalization



Data augmentation

- A classifier has to be invariant to a wide variety of transformations

All

Images

Videos

News

Shopping

More

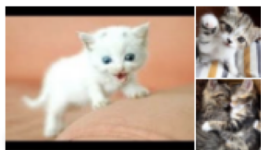
Settings

Tools

SafeSearch ▼



Cute



And Kittens



Clipart



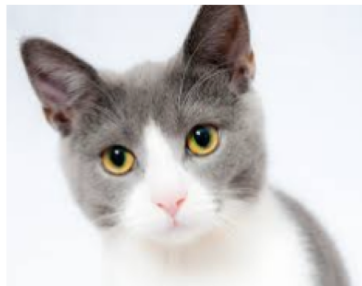
Drawing



Cute Baby



White Cats And Kittens



Pose

Appearance

Illumination

Data augmentation

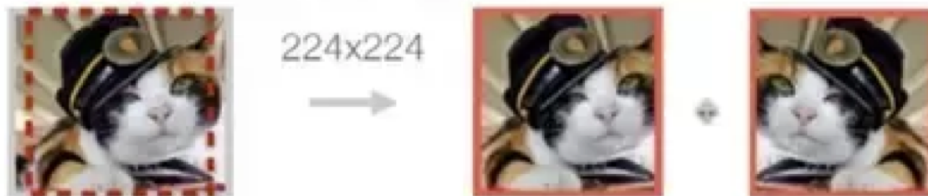
- A classifier has to be invariant to a wide variety of transformations
- Helping the classifier: generate fake data simulating plausible transformations

Data augmentation

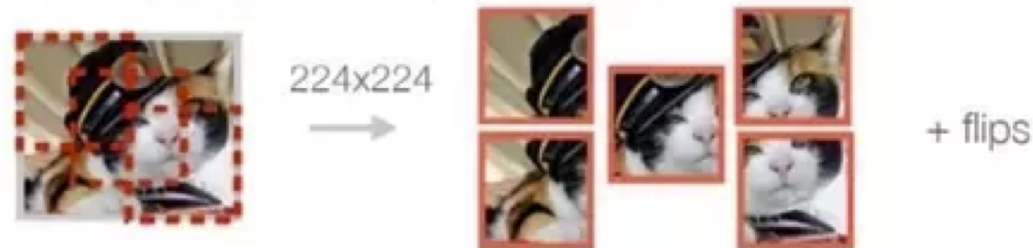
a. No augmentation (= 1 image)



b. Flip augmentation (= 2 images)

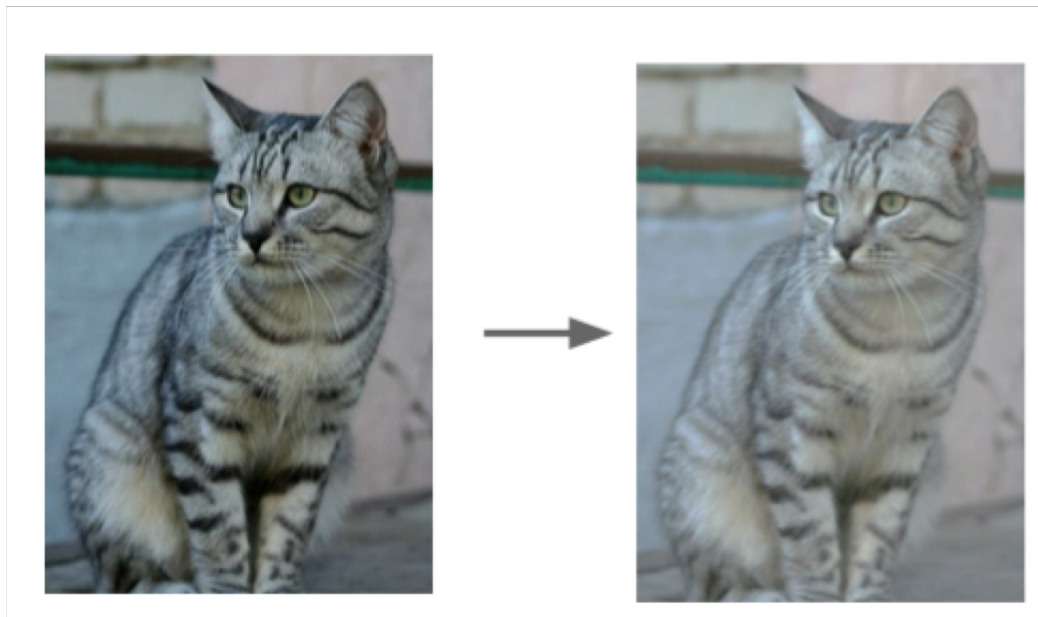


c. Crop+Flip augmentation (= 10 images)



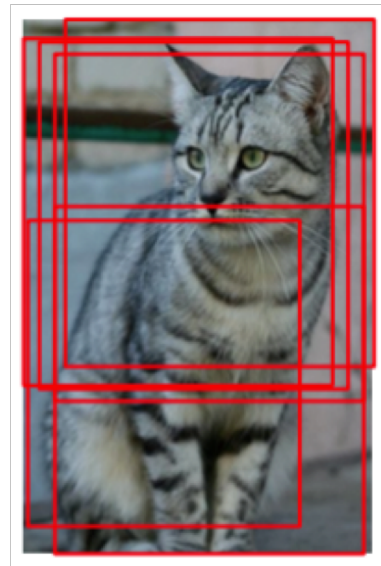
Data augmentation: random crops

- Random brightness and contrast changes



Data augmentation: random crops

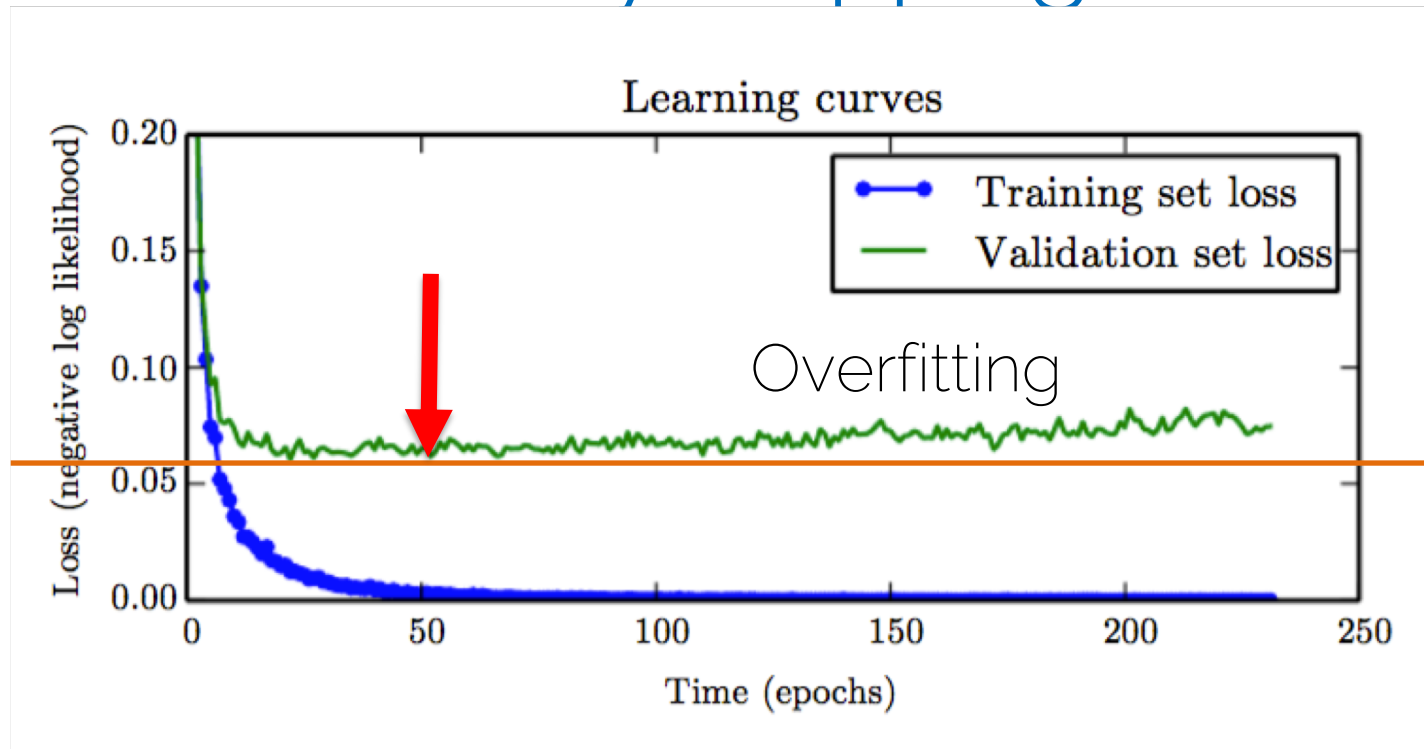
- Training: random crops
 - Pick a random L in $[256, 480]$
 - Resize training image, short side L
 - Randomly sample crops of 224×224
- Testing: fixed set of crops
 - Resize image at N scales
 - 10 fixed crops of 224×224 : 4 corners + center + flips



Data augmentation

- When comparing two networks make sure to use the same data augmentation!
- Consider data augmentation a part of your network design

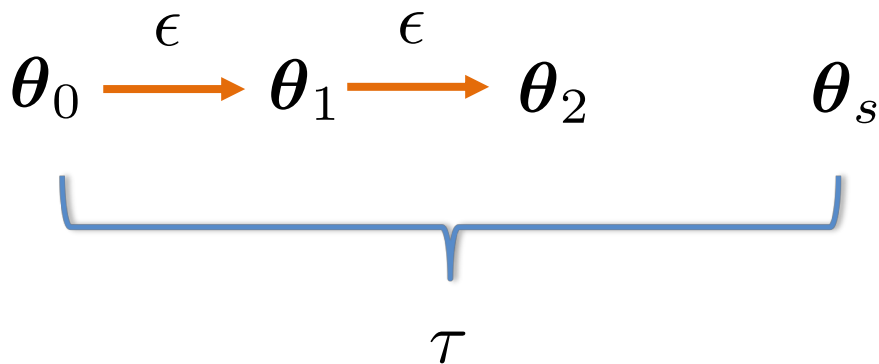
Early stopping



Training time is also a hyperparameter

Early stopping

- Easy form of regularization



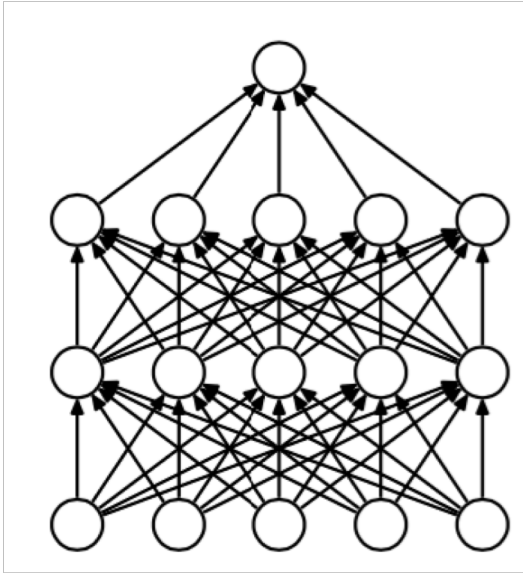
θ^*
Overfitting

Bagging and ensemble methods

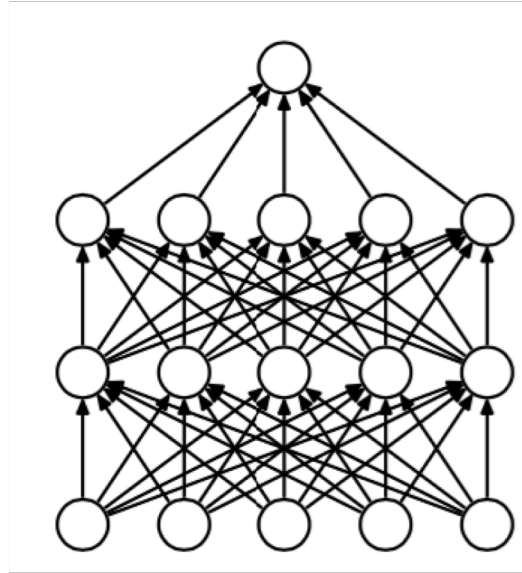
- Train three models and average their results
- Change a different algorithm for optimization or change the objective function
- If errors are uncorrelated, the expected combined error will decrease linearly with the ensemble size

Bagging and ensemble methods

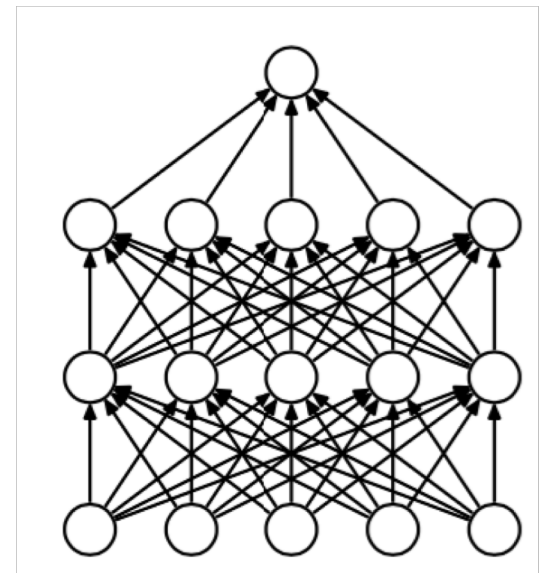
- Bagging: uses k different datasets



Training Set 1



Training Set 2

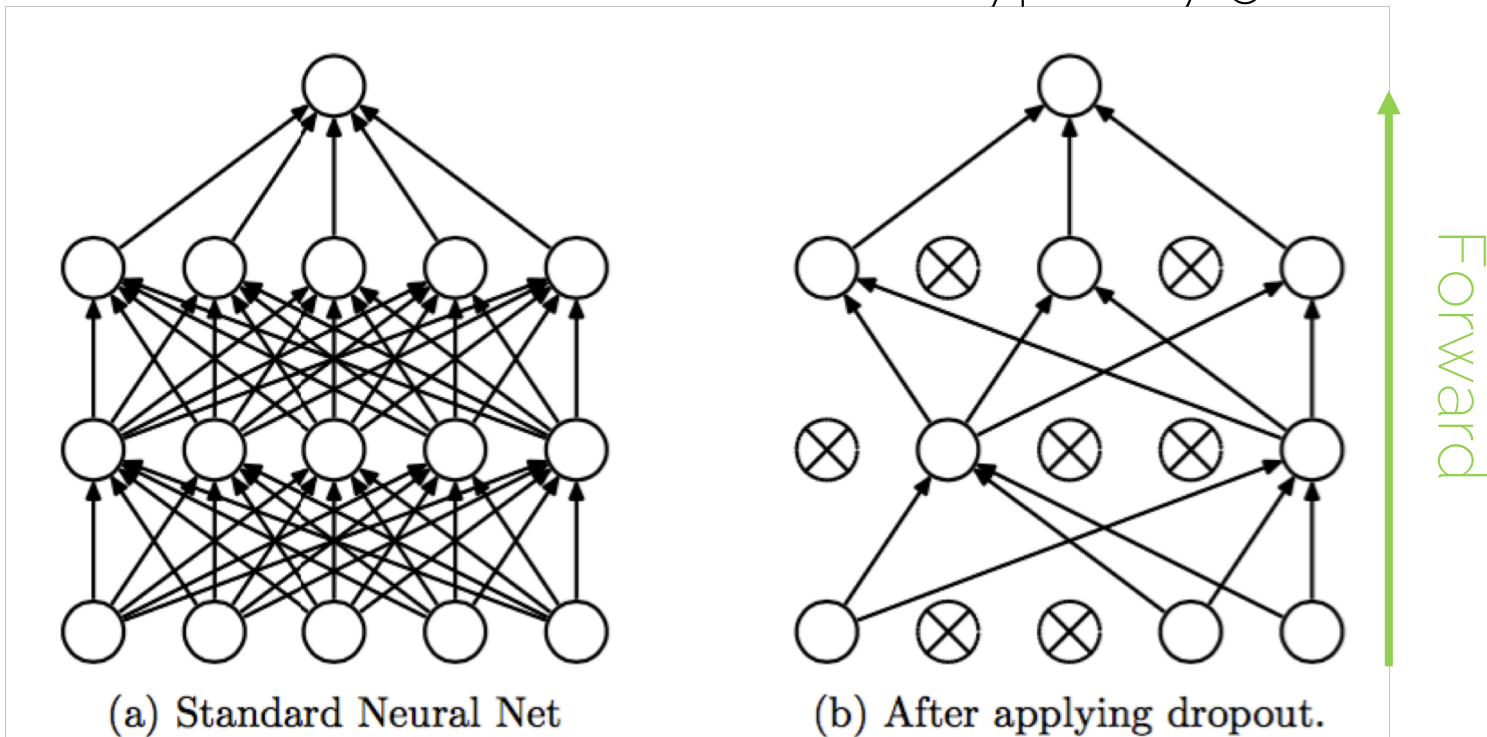


Training Set 3

Dropout

Dropout

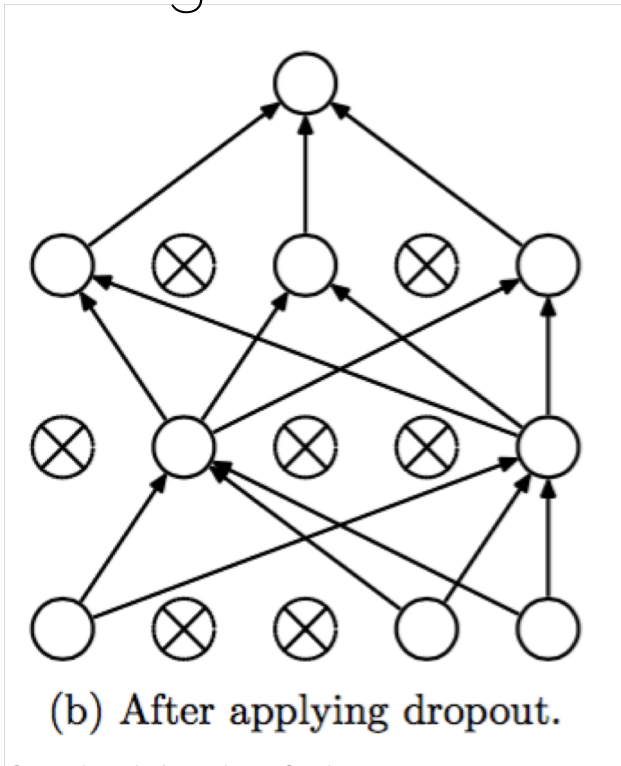
- Disable a random set of neurons (typically 50%)



Dropout: intuition

- Using half the network = half capacity

Redundant
representations



Furry

Has two eyes

Has a tail

Has paws

Has two ears

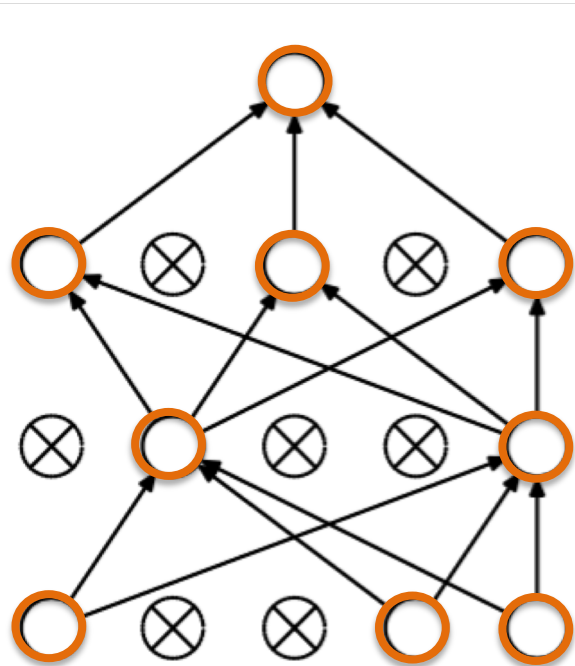


Dropout: intuition

- Using half the network = half capacity
 - Redundant representations
 - Base your scores on more features
- Consider it as model ensemble

Dropout: intuition

- Two models in one

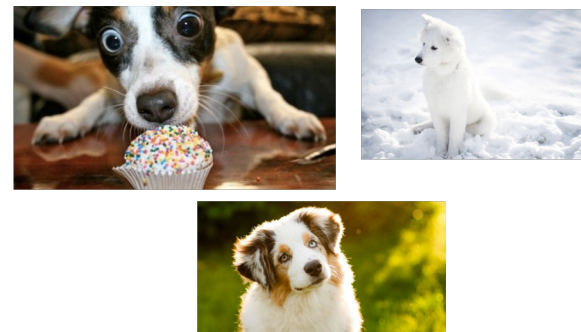


(b) After applying dropout.

○ Model 1



☐ Model 2



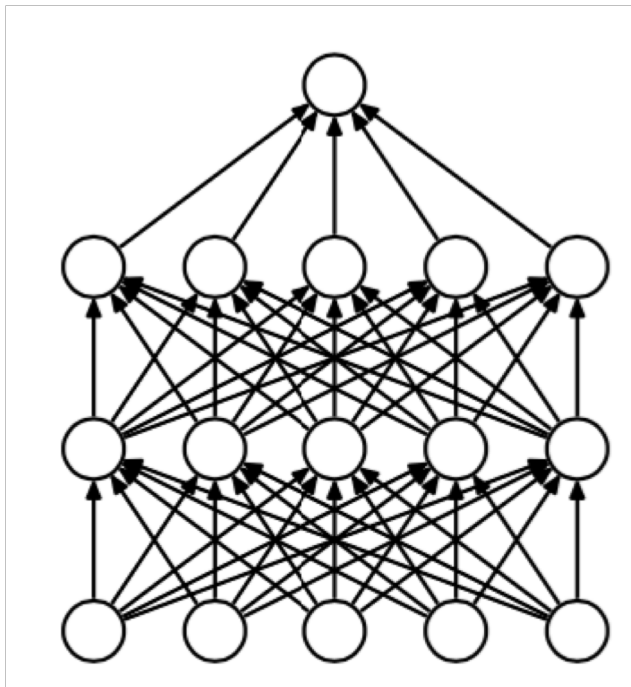
Dropout: intuition

- Using half the network = half capacity
 - Redundant representations
 - Base your scores on more features
- Consider it as two models in one
 - Training a large ensemble of models, each on different set of data (mini-batch) and with SHARED parameters

Reducing co-adaptation between neurons

Dropout: test time

- All neurons are “turned on” – no dropout

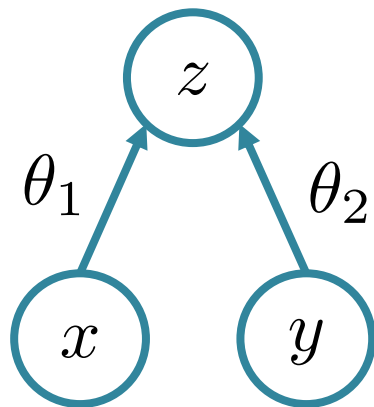


Conditions at train
and test time are
not the same

Dropout: test time

Dropout
probability
 $p=0.5$

- Test: $z = \theta_1 x + \theta_2 y$



Weight scaling
inference rule

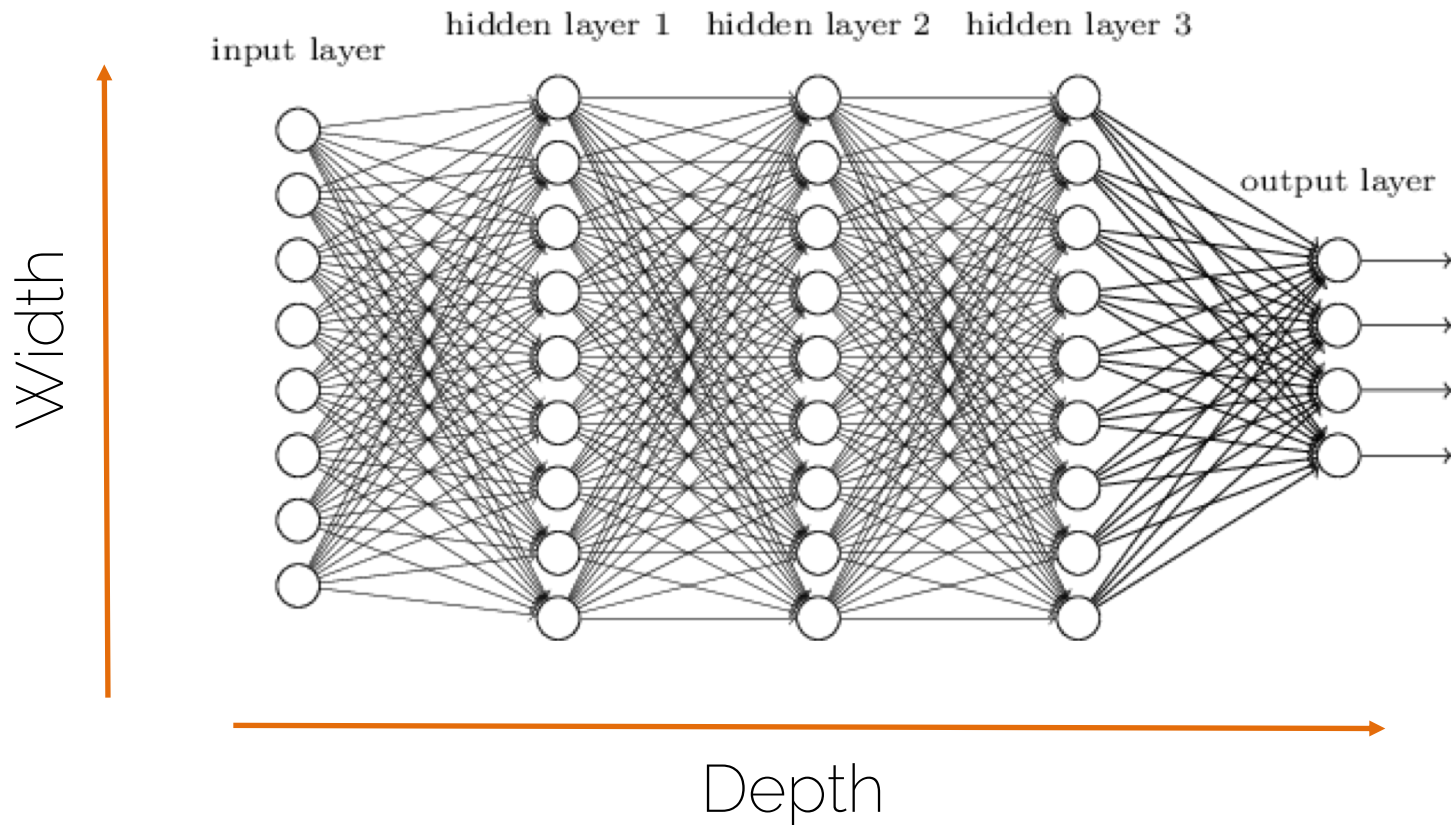
- Train:
$$\begin{aligned} \mathbb{E}[z] &= \frac{1}{4} (\theta_1 0 + \theta_2 0 \\ &\quad + \theta_1 x + \theta_2 0 \\ &\quad + \theta_1 0 + \theta_2 y \\ &\quad + \theta_1 x + \theta_2 y) \\ &= \frac{1}{2} (\theta_1 x + \theta_2 y) \end{aligned}$$

Dropout: verdict

- Efficient bagging method with parameter sharing
- Use it!
- Dropout reduces the effective capacity of a model → larger models, more training time

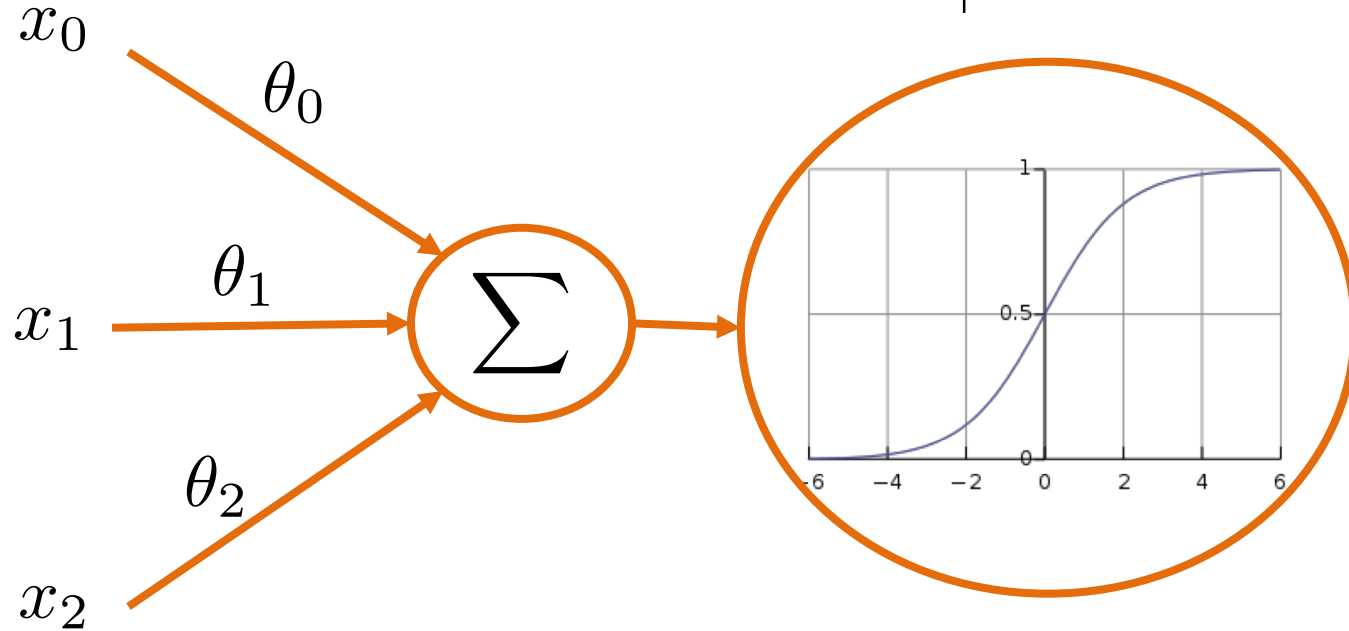
Recap

What do we know so far?



What do we know so far?

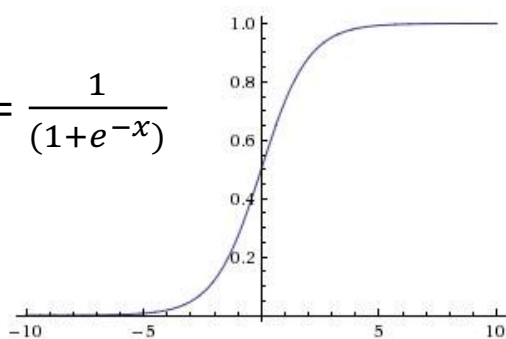
Concept of a 'Neuron'



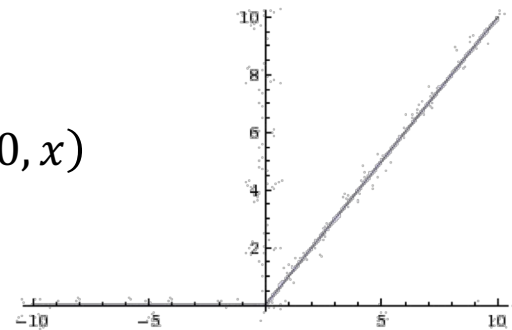
What do we know so far?

Activation Functions (non-linearities)

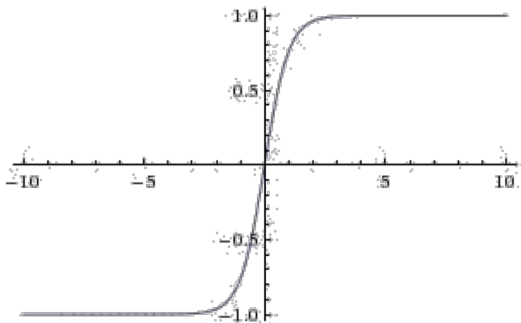
Sigmoid: $\sigma(x) = \frac{1}{1+e^{-x}}$



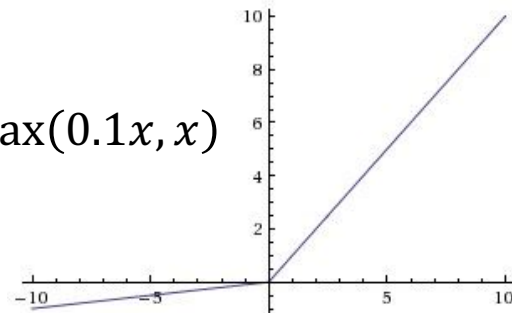
ReLU: $\max(0, x)$



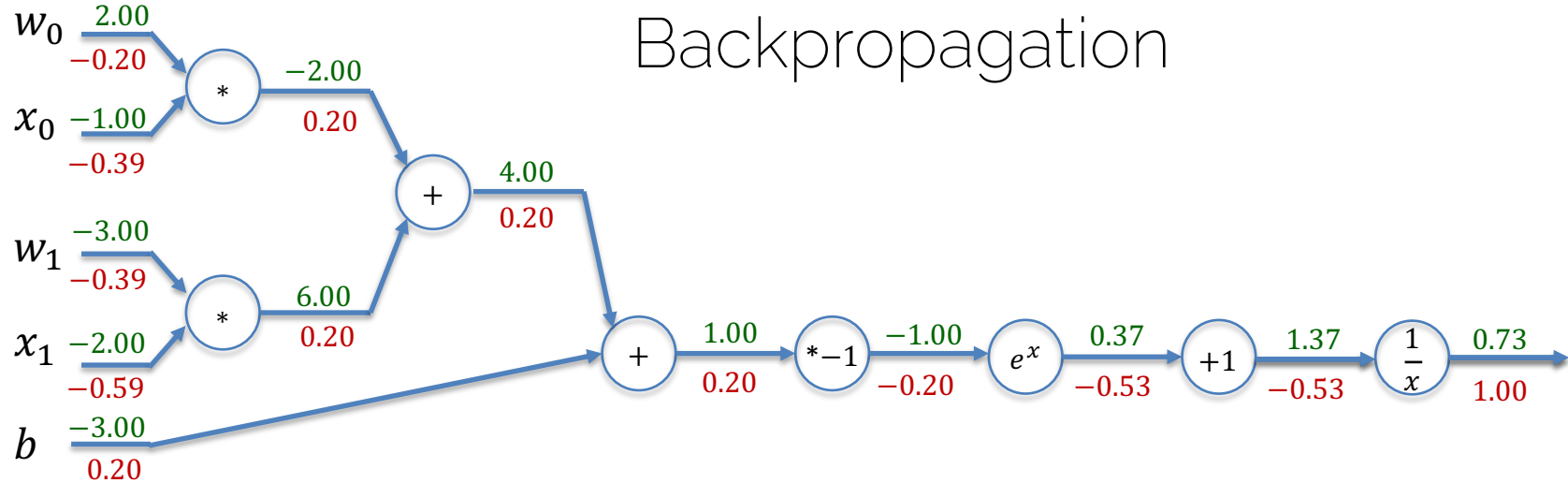
tanh: $\tanh(x)$



Leaky ReLU: $\max(0.1x, x)$

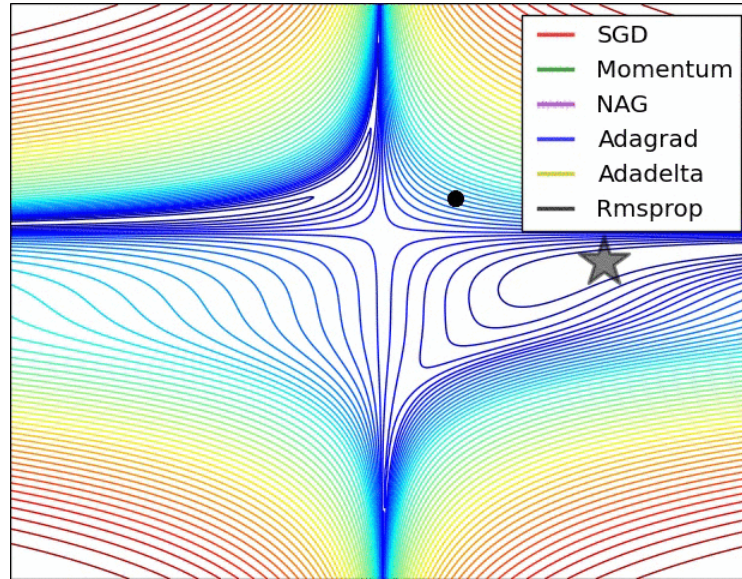


What do we know so far?



What do we know so far?

SGD Variations (Momentum, etc.)



What do we know so far?

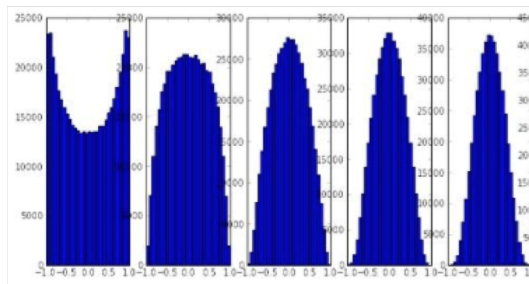
Data Augmentation



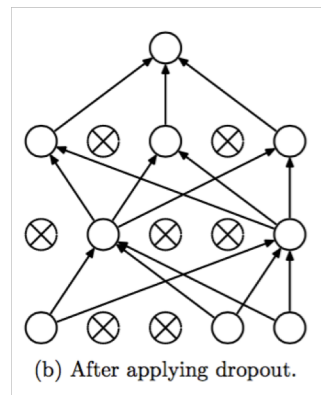
Batch-Norm

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

Weight Initialization (e.g., Xavier/2)



Dropout



Weight Regularization

e.g., L^2 -reg: $R^2(W) = \sum_{i=1}^N w_i^2$

Why not only more Layers?

- We can not make networks arbitrarily complex
 - Why not just go deeper and get better?
 - No structure!!
 - It's just brute force!
 - Optimization becomes hard
 - Performance plateaus / drops!

Administrative Things

- Happy holidays!
- Tuesday December 18th: solution to Exercise 2, introduction to exercise 3 and introduction to PyTorch!
- Thursday January 10th: Starting with CNN