

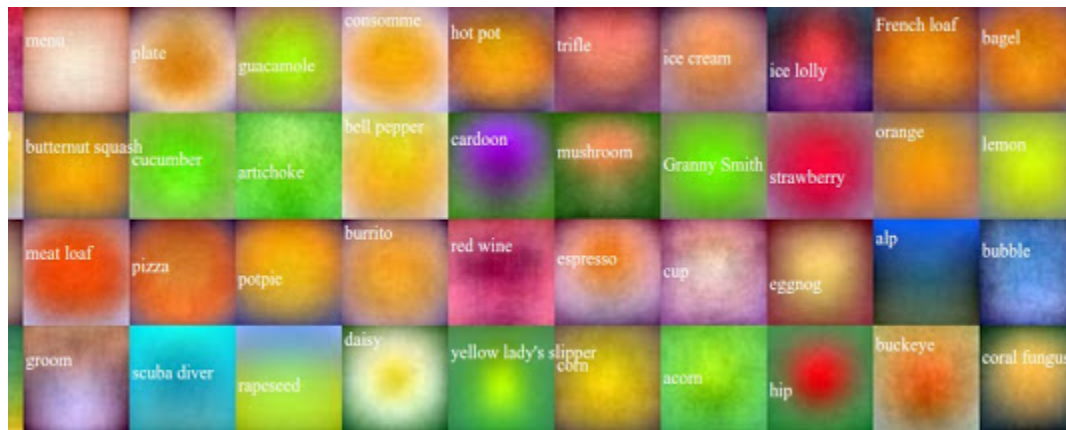
Introduction to Neural Networks

Beyond linear

- Linear score function $f = Wx$



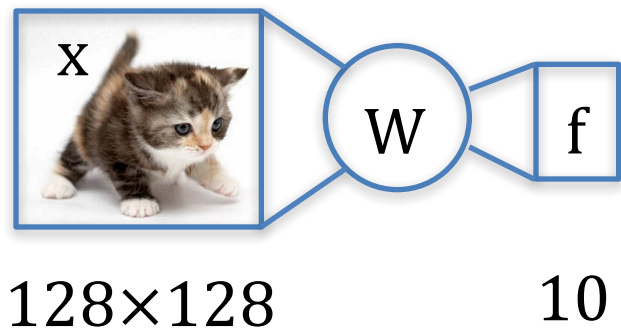
On CIFAR-10



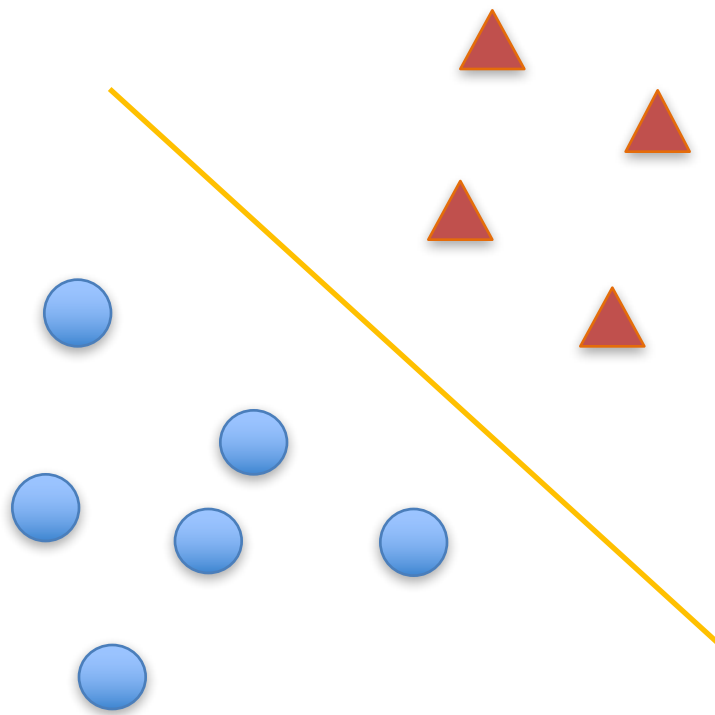
On ImageNet

Beyond linear

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



LINEAR
TRANSFORMATION



Beyond linear

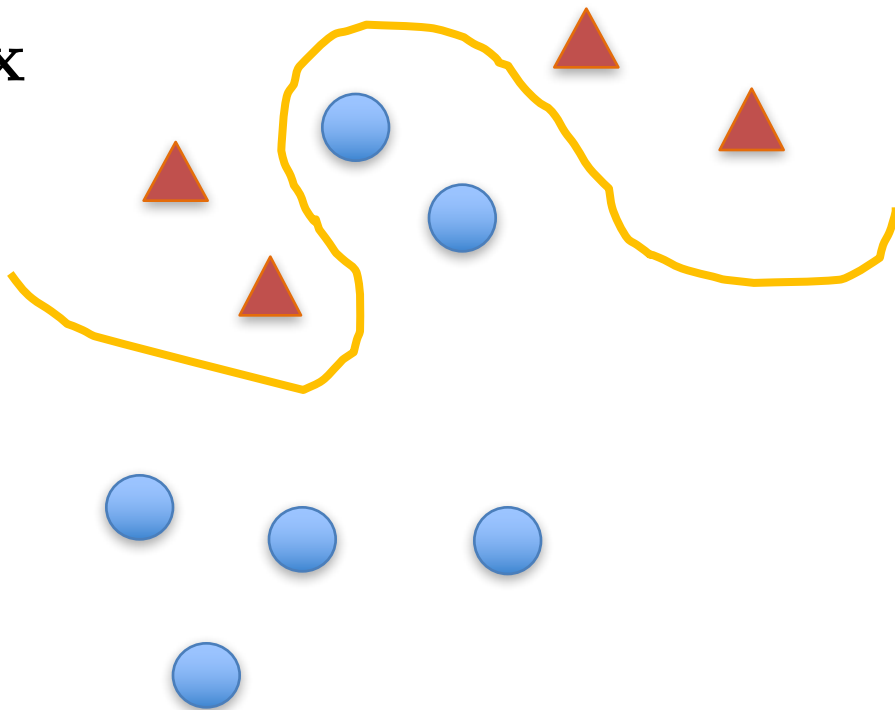
1-layer network: $f = \mathbf{W}\mathbf{x}$



128×128

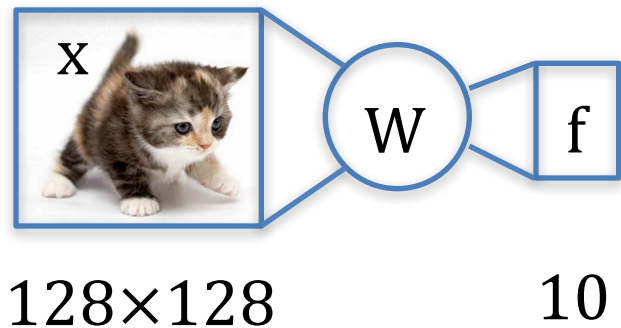
10

LINEAR
TRANSFORMATION



Kernel trick

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

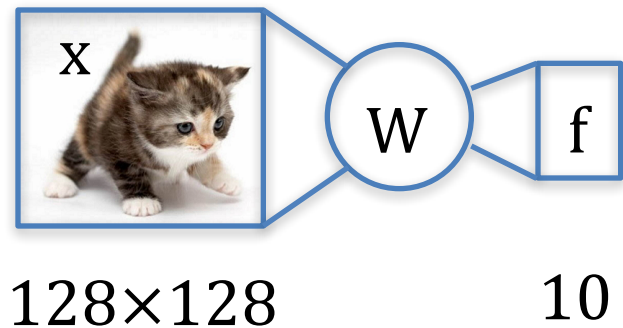


$$f = \mathbf{W}\phi(\mathbf{x})$$

↑
kernel

Neural networks

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



$$f = \mathbf{W}\phi(\mathbf{x}; \boldsymbol{\theta})$$

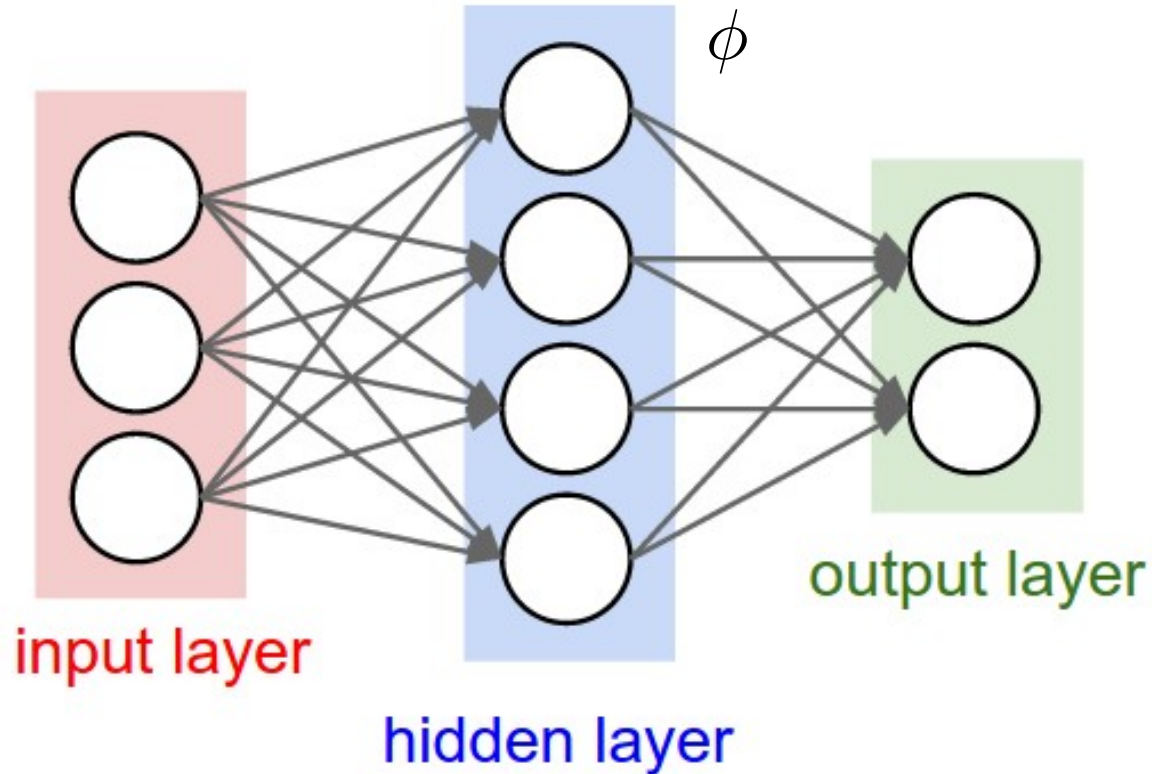
↑
kernel

↑
parameters

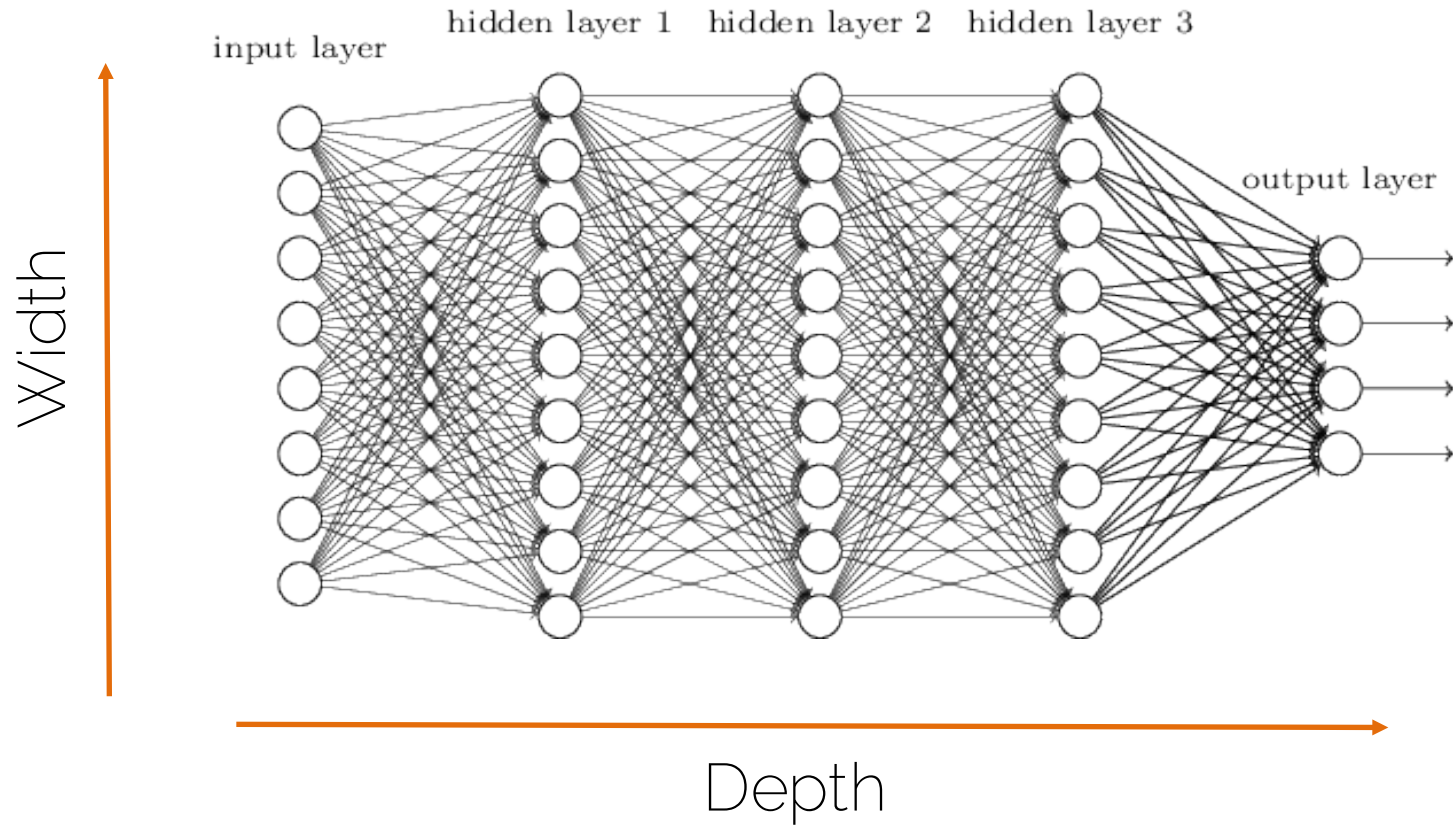
From the broad family of functions ϕ we learn the best representation by learning the parameters $\boldsymbol{\theta}$

Neural Network

Also SVM
is in this
category



Neural Network



Neural Network

- Linear score function $f = Wx$
- Neural network is a nesting of 'functions'
 - 2-layers: $f = W_2 \max(0, W_1 x)$
 - 3-layers: $f = W_3 \max(0, W_2 \max(0, W_1 x))$
 - 4-layers: $f = W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x)))$
 - 5-layers: $f = W_5 \sigma(W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x))))$
 - ... up to hundreds of layers

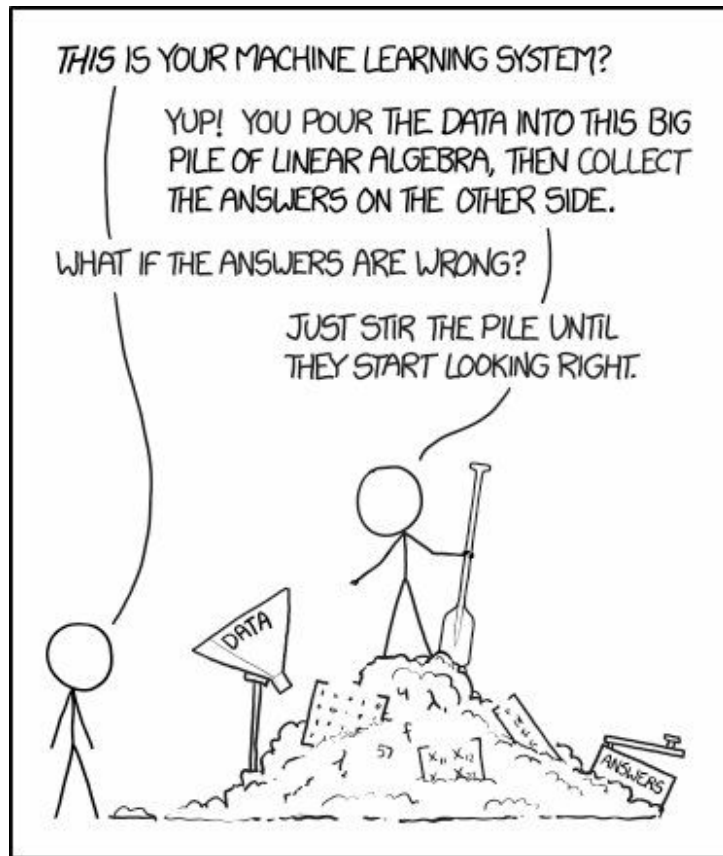
Playing around with networks

- <http://cs.stanford.edu/people/karpathy/convnetjs/index.html>

Neural Network

- Problems of going deeper...
- The impact of small decisions (architecture, activation functions...)
- Is my network training correctly?

A typical Deep Learner day

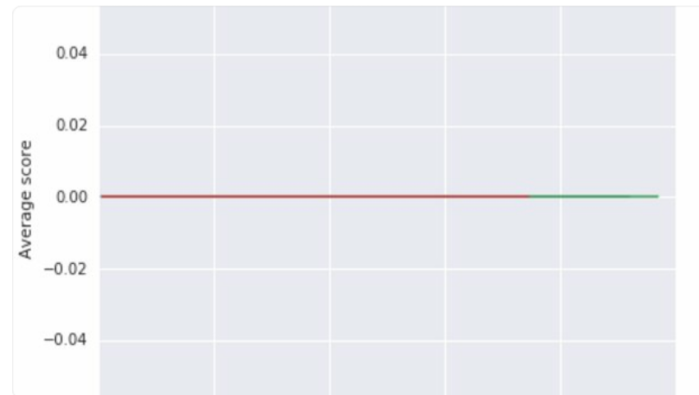


♥ A Andrej Karpathy i Ian Goodfellow els agrada



Oriol Vinyals @OriolVinyalsML · 9h

A typical training curve in Montezuma's Revenge (note: there are several random seeds which overlap) 🤪 #nips #rl #exploration



← 2

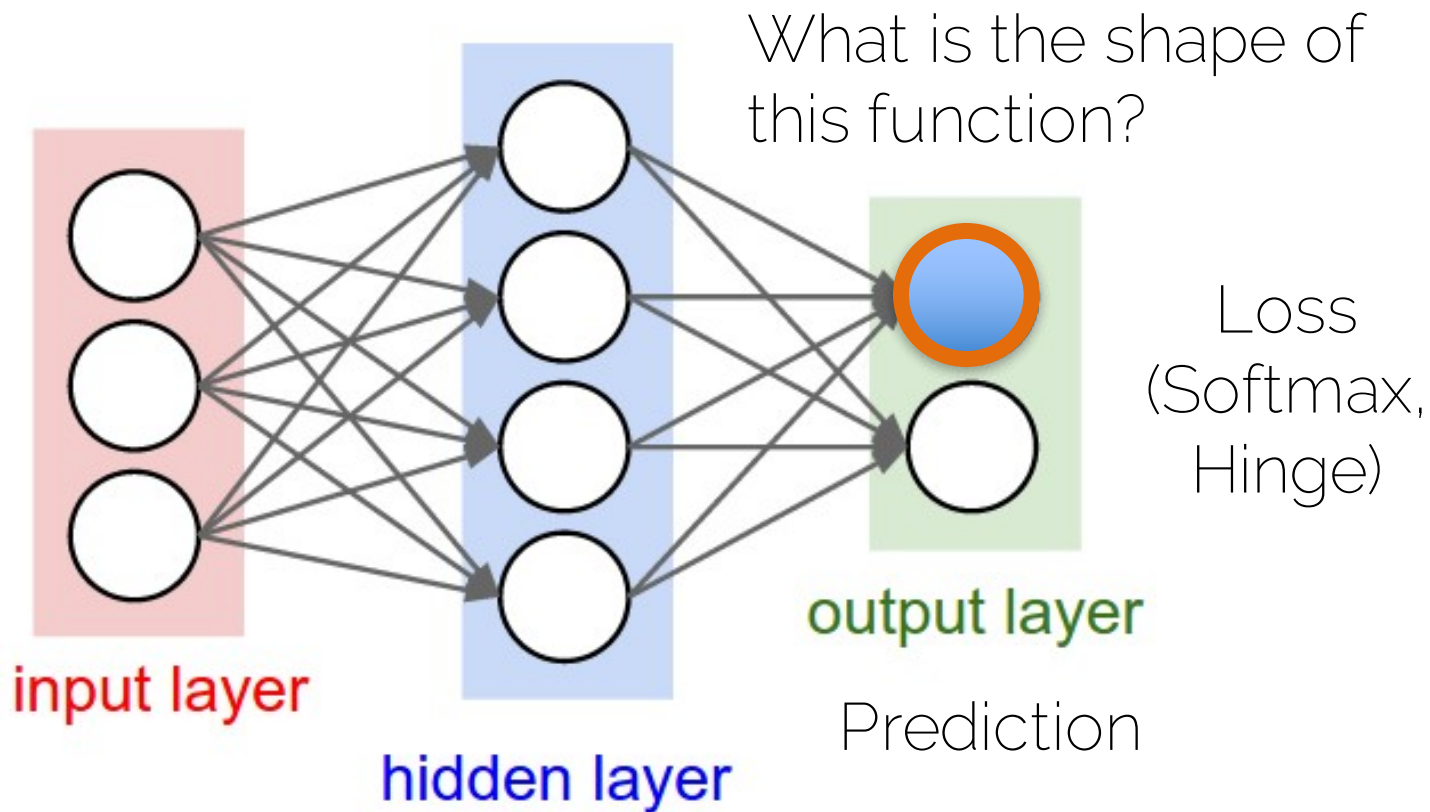
↺ 9

♥ 63

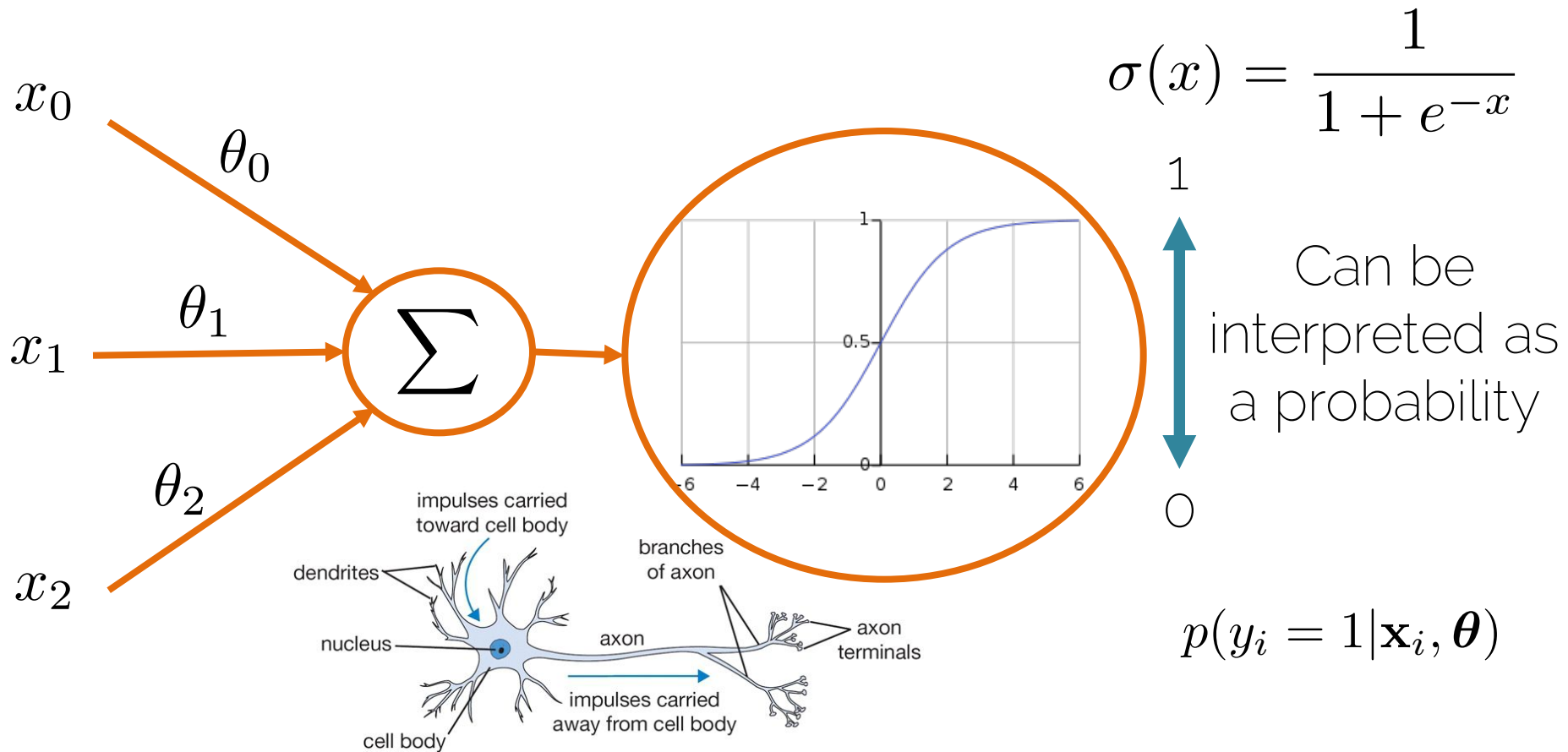


Output functions

Neural networks




Sigmoid for binary predictions



Logistic regression

- Probability of a binary output

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \prod_{i=1}^n \text{Ber}(y_i | \text{sigm}(\mathbf{x}_i, \boldsymbol{\theta}))$$


Model for
coins

$$p(x|\phi) = \phi^x (1 - \phi)^{1-x} = \begin{cases} \phi & \text{if } x = 1 \\ 1 - \phi & \text{if } x = 0 \end{cases}$$

Logistic regression

- Probability of a binary output

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) &= \prod_{i=1}^n \text{Ber}(y_i | \text{sigm}(\mathbf{x}_i, \boldsymbol{\theta})) \\ &= \prod_{i=1}^n \underbrace{\left[\frac{1}{1 + e^{-\mathbf{x}_i \boldsymbol{\theta}}} \right]}_{\Pi_i}^{y_i} \underbrace{\left[1 - \frac{1}{1 + e^{-\mathbf{x}_i \boldsymbol{\theta}}} \right]}_{\Pi_i}^{1-y_i} \end{aligned}$$

$$p(x|\phi) = \phi^x (1 - \phi)^{1-x}$$

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

Logistic regression

- Probability of a binary output

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \prod_{i=1}^n [\Pi_i]^{y_i} [1 - \Pi_i]^{1-y_i}$$

- Maximum Likelihood

$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})$$

Logistic regression

- Probability of a binary output $\Pi_i = \frac{1}{1 + e^{-\mathbf{x}_i \boldsymbol{\theta}}}$

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \prod_{i=1}^n [\Pi_i]^{y_i} [1 - \Pi_i]^{1-y_i}$$

$$C(\boldsymbol{\theta}) = -\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})$$

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

Referred to as cross-entropy

Logistic regression

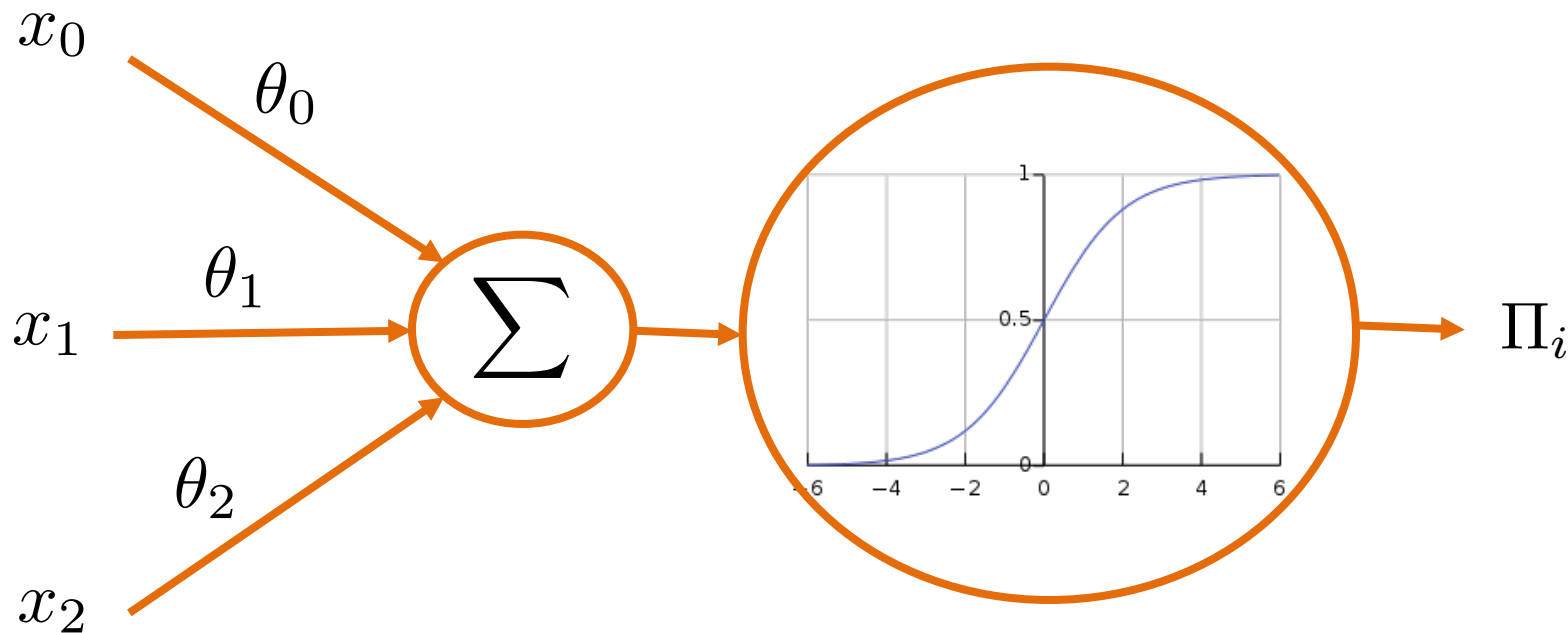
- Optimize using gradient descent
- Saturation occurs only when the model already has the right answer

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

Referred to as cross-entropy

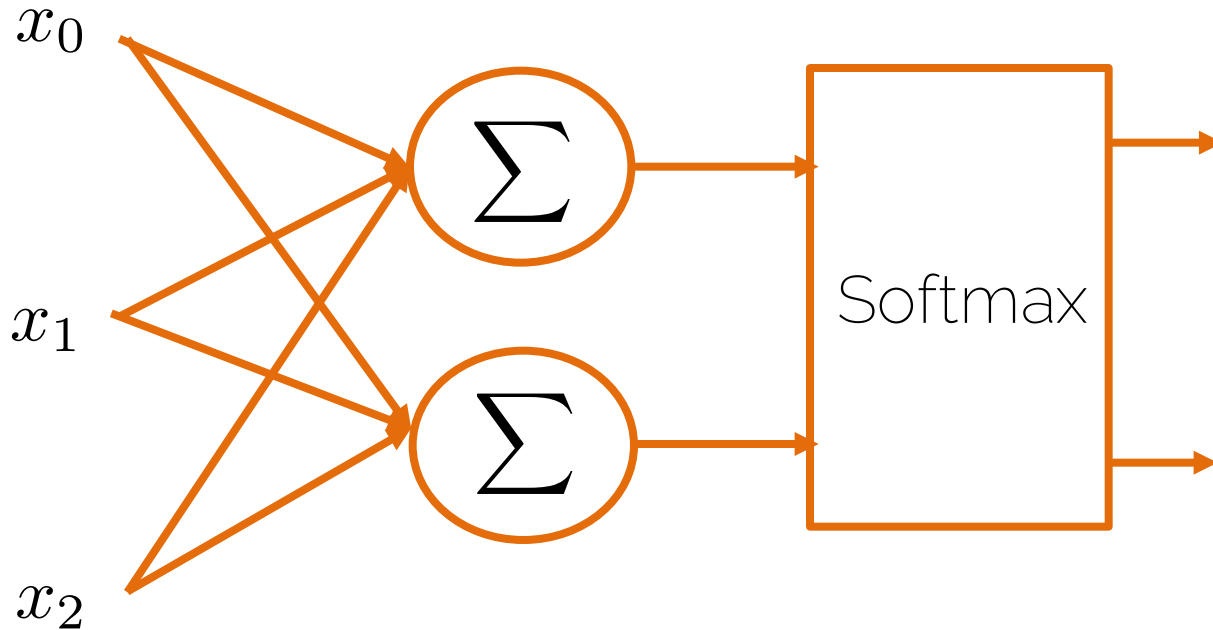
Softmax formulation

- What if we have multiple classes?



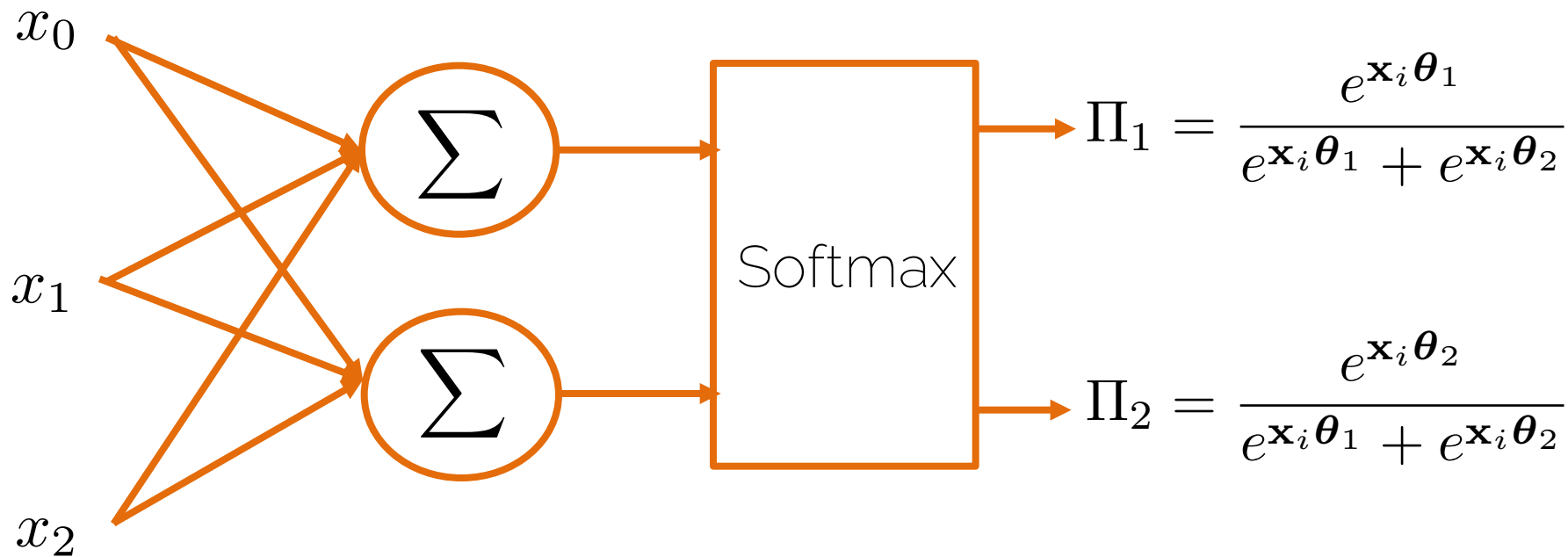
Softmax formulation

- What if we have multiple classes?



Softmax formulation

- What if we have multiple classes?



Softmax formulation

- Softmax

$$p(y_i | \mathbf{x}, \boldsymbol{\theta}) = \frac{e^{\mathbf{x}\boldsymbol{\theta}_i}}{\sum_{k=1}^n e^{\mathbf{x}\boldsymbol{\theta}_k}}$$

exp

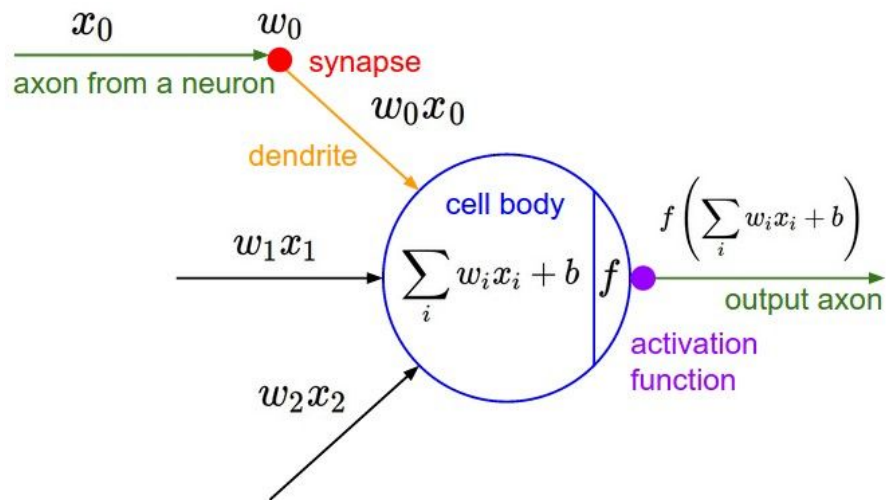
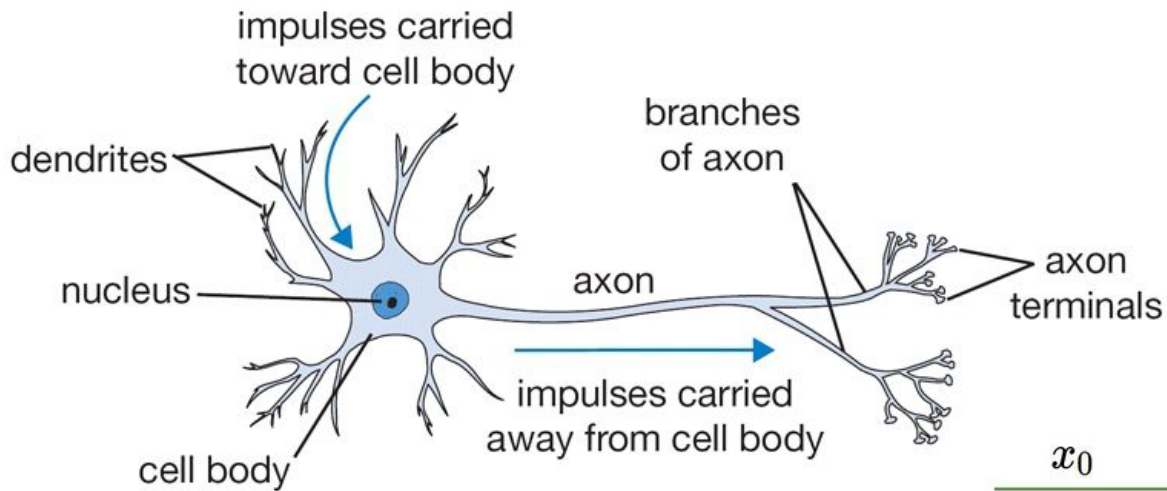
normalize

- Softmax loss (ML)

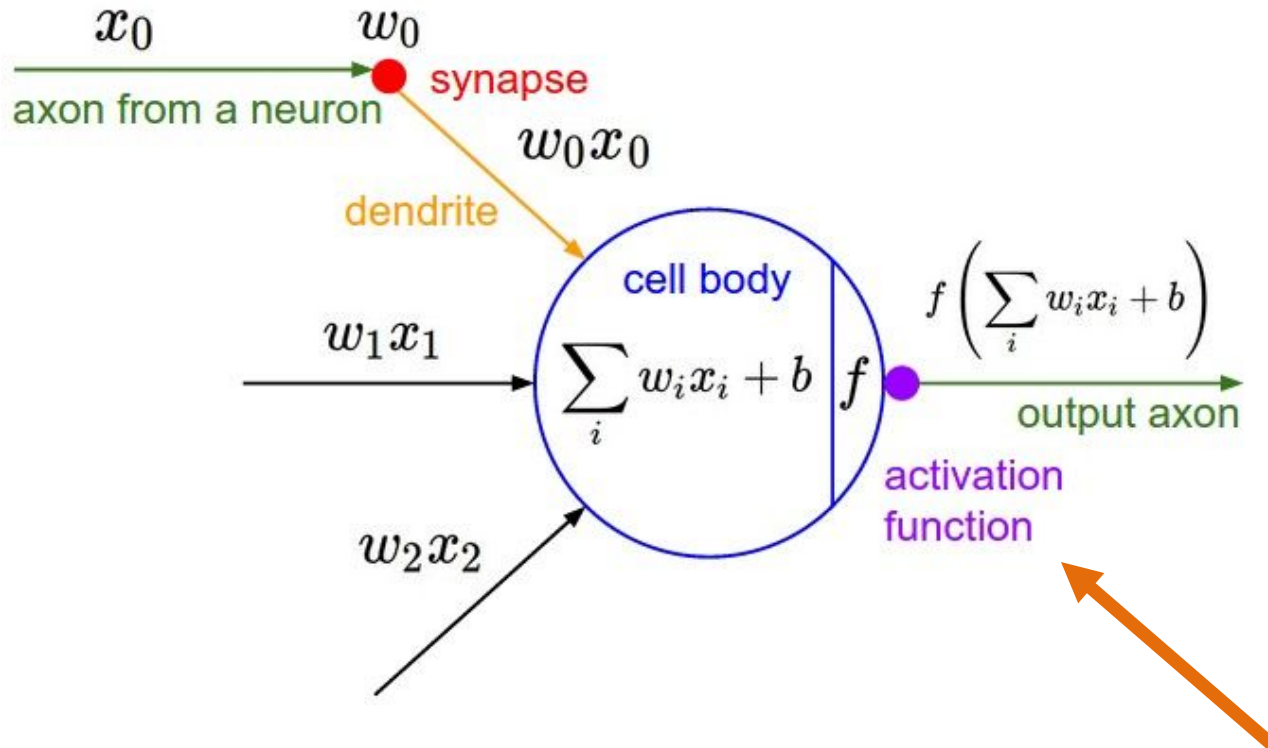
$$L_i = -\log \left(\frac{e^{s_{y_i}}}{\sum_k e^{s_k}} \right)$$

Activation functions

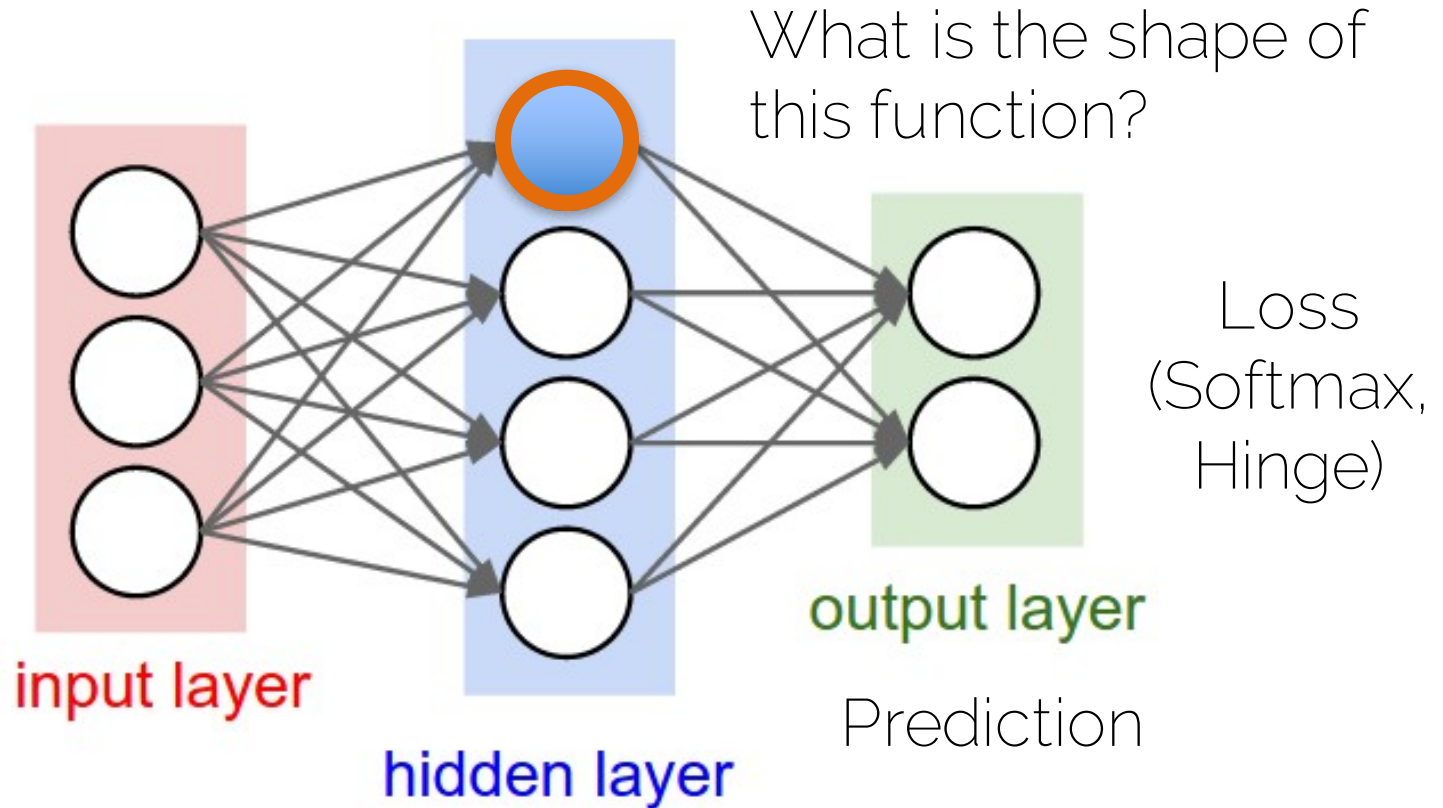
Neurons



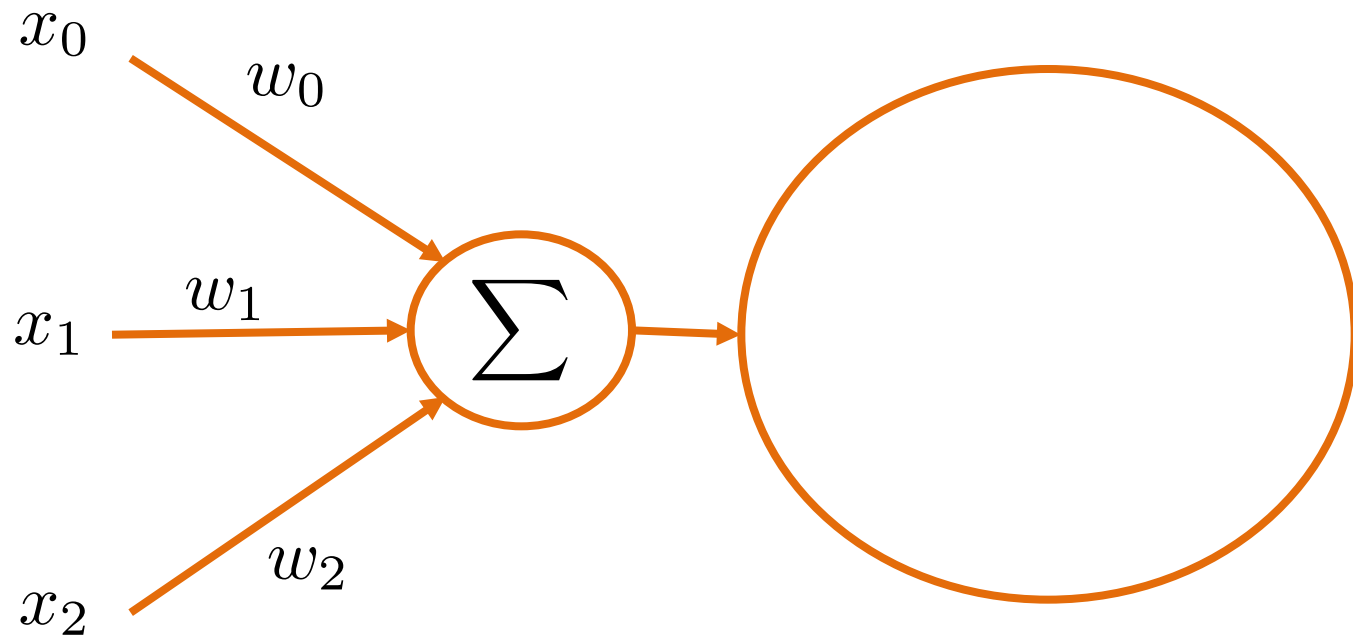
Neurons



Neural networks

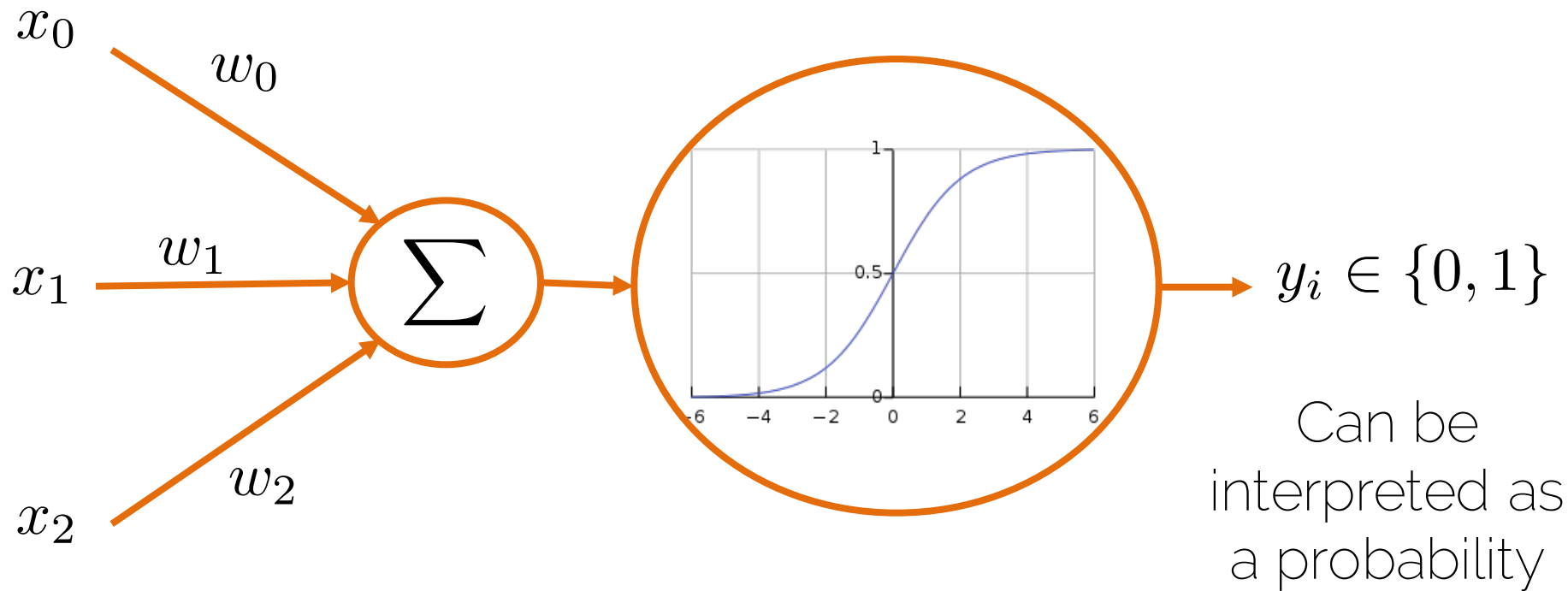


Activation functions or hidden units



Sigmoid

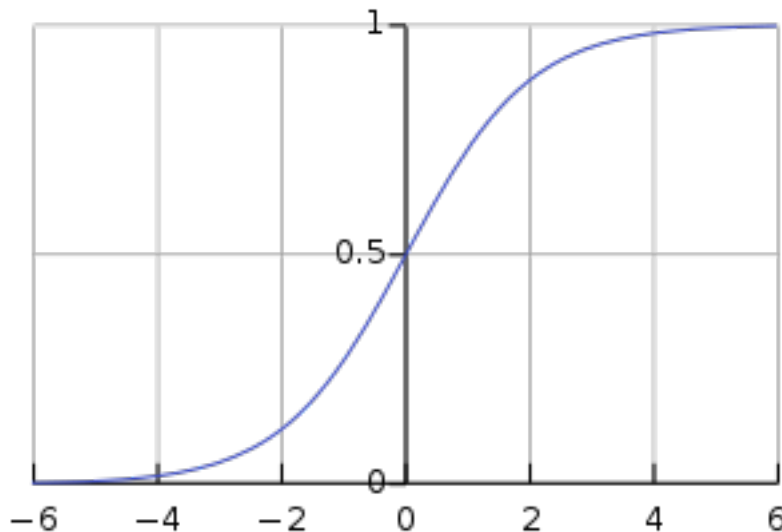
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



Sigmoid

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

Forward



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



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

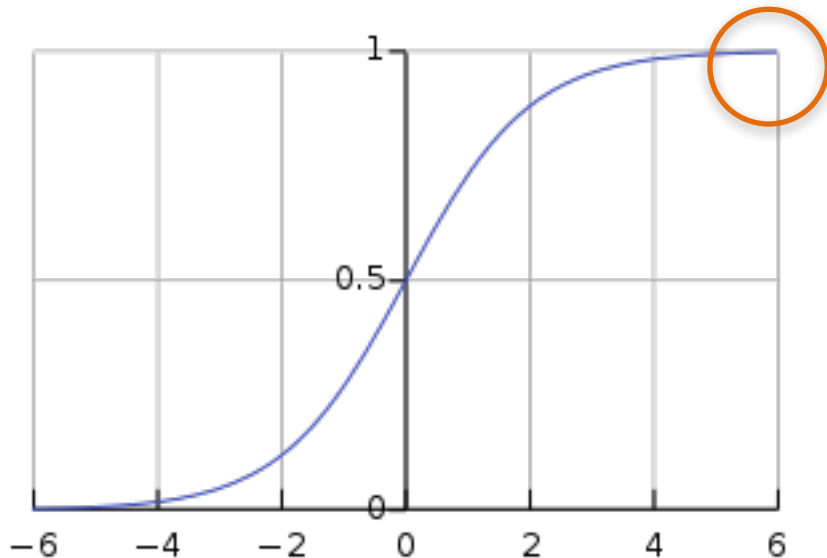


$$\frac{\partial L}{\partial \sigma}$$

Sigmoid

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

Forward



$$x = 6$$

✗ Saturated neurons kill the gradient flow

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



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

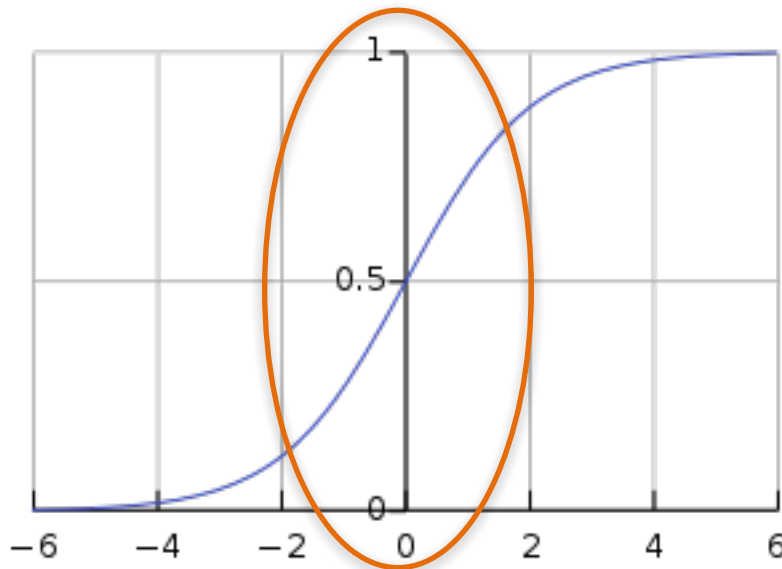


$$\frac{\partial L}{\partial \sigma}$$

Sigmoid

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

Forward



Active region
for gradient
descent

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



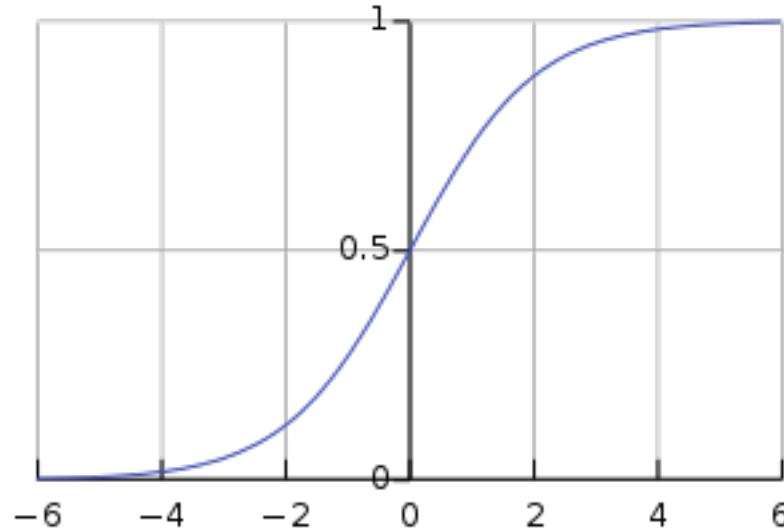
$$\frac{\partial \sigma}{\partial x}$$



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

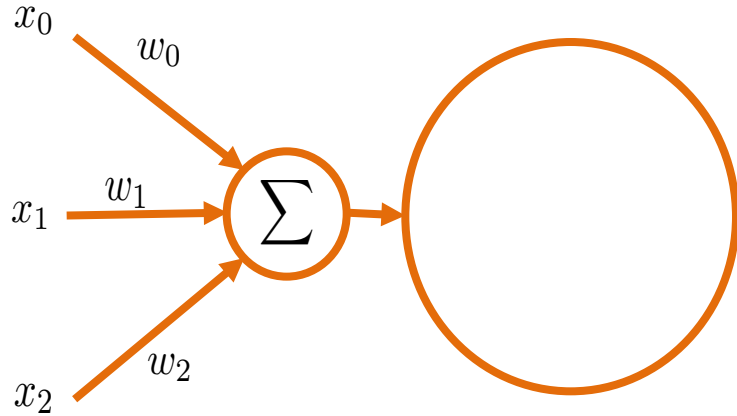
Sigmoid

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



Output is
always
positive

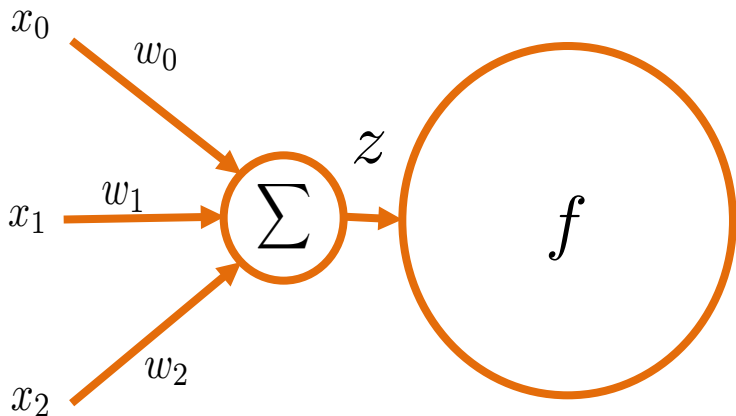
Problem of positive output



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

We want to compute the gradient wrt the weights

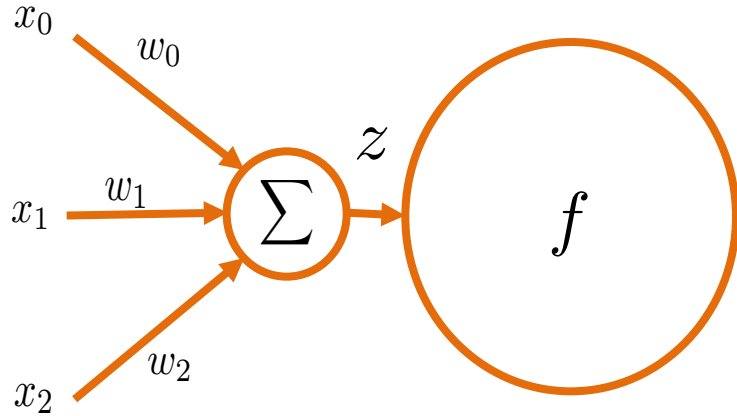
Problem of positive output



$$f \left(\underbrace{\sum_i w_i x_i + b}_{\frac{\partial z}{\partial w} = x_i > 0} \right)$$

We want to compute the gradient wrt the weights

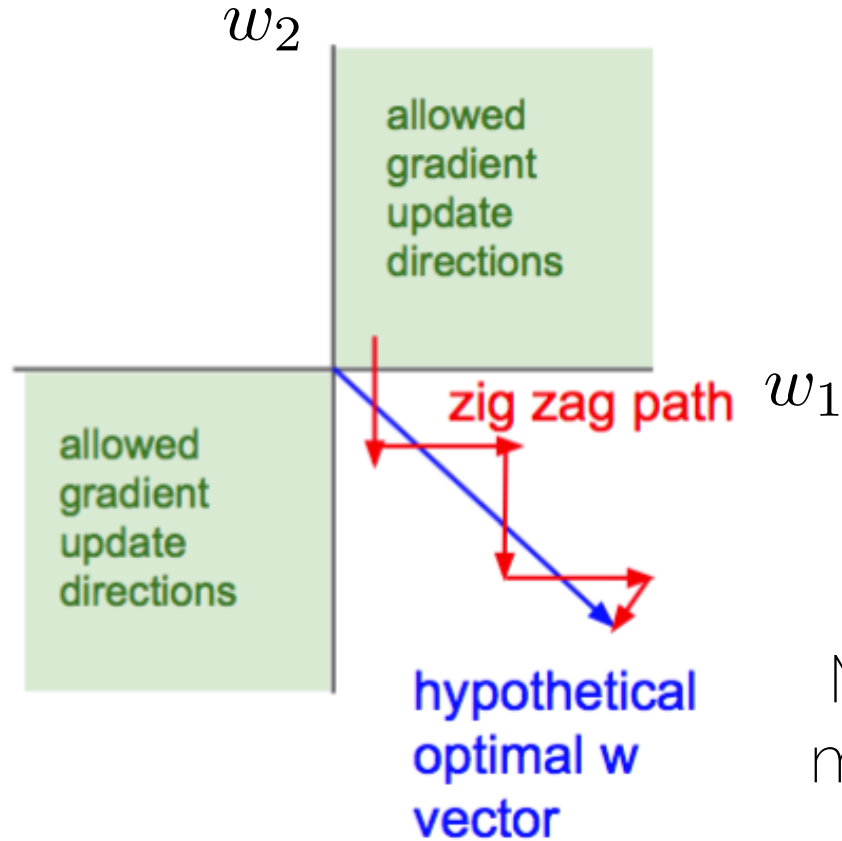
Problem of positive output



$$\frac{\partial f}{\partial z} \rightarrow f \left(\underbrace{\sum_i w_i x_i + b}_{\frac{\partial z}{\partial w} = x_i > 0} \right)$$

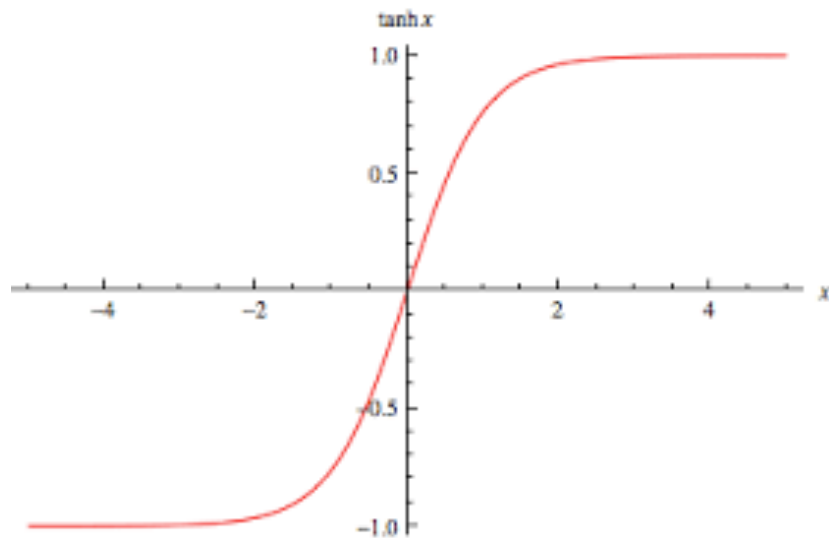
It is going to be either positive or negative for all weights

Problem of positive output



More on zero-mean data later

tanh



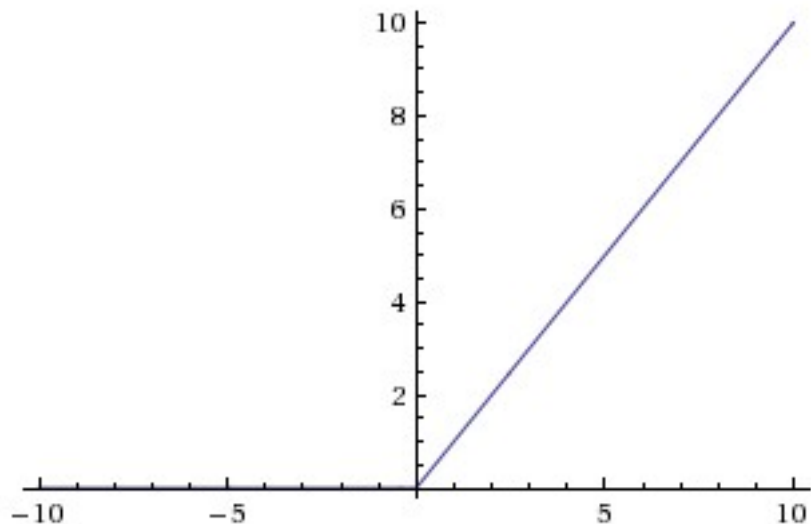
✗ Still saturates

✓ Zero-centered

✗ Still saturates

Rectified Linear Units (ReLU)

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



Large and
consistent
gradients ✓



Fast convergence



Does not saturate

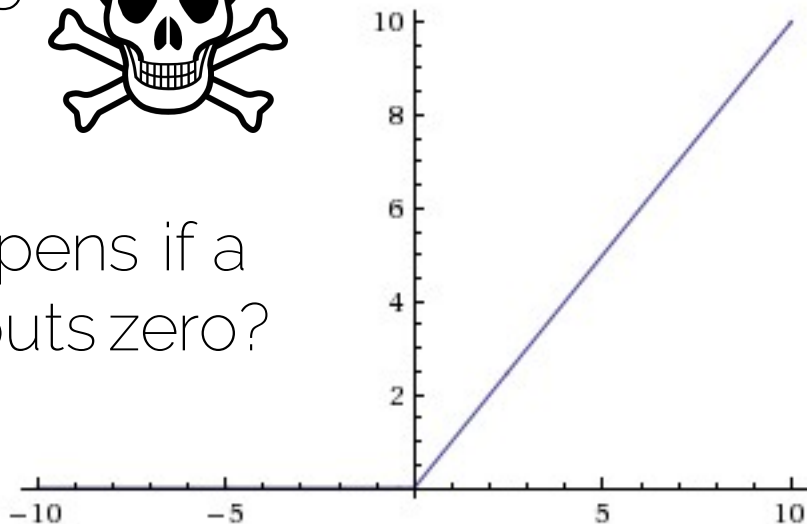
Rectified Linear Units (ReLU)



Dead ReLU



What happens if a ReLU outputs zero?



Large and consistent gradients



Fast convergence



Does not saturate

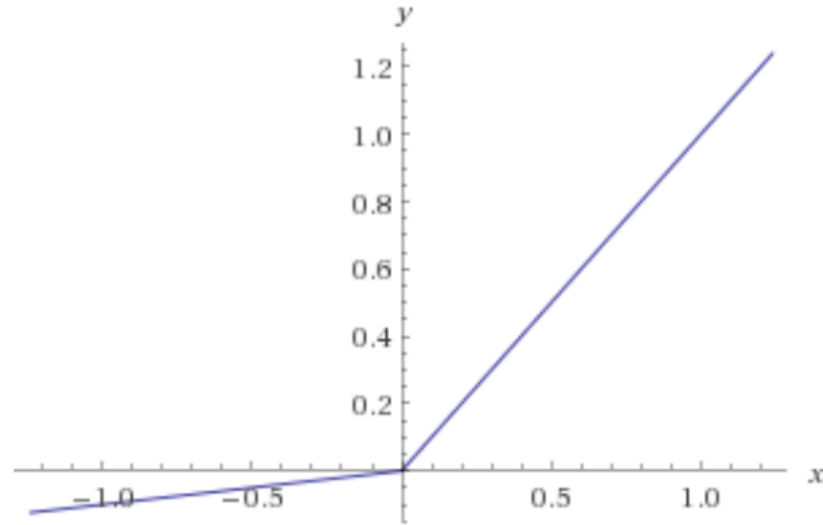
Rectified Linear Units (ReLU)

- Initializing ReLU neurons with slightly positive biases (0.1) makes it likely that they stay active for most inputs

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

Leaky ReLU

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



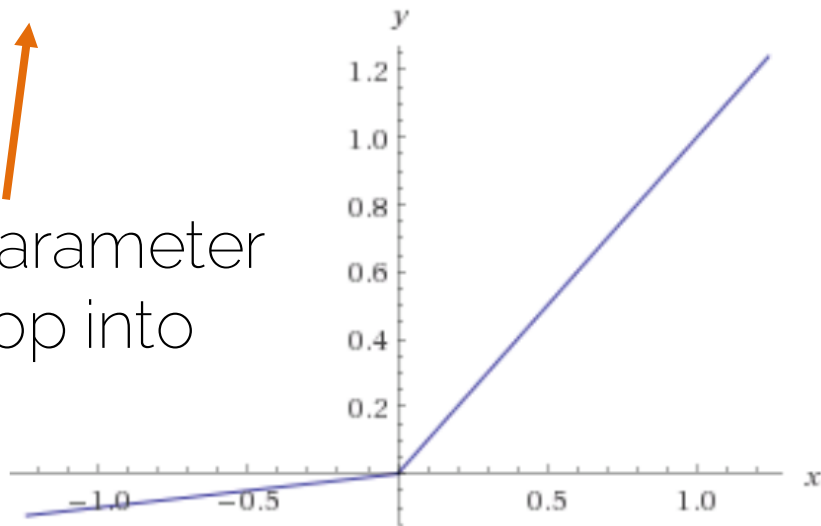
Does not die

Parametric ReLU

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

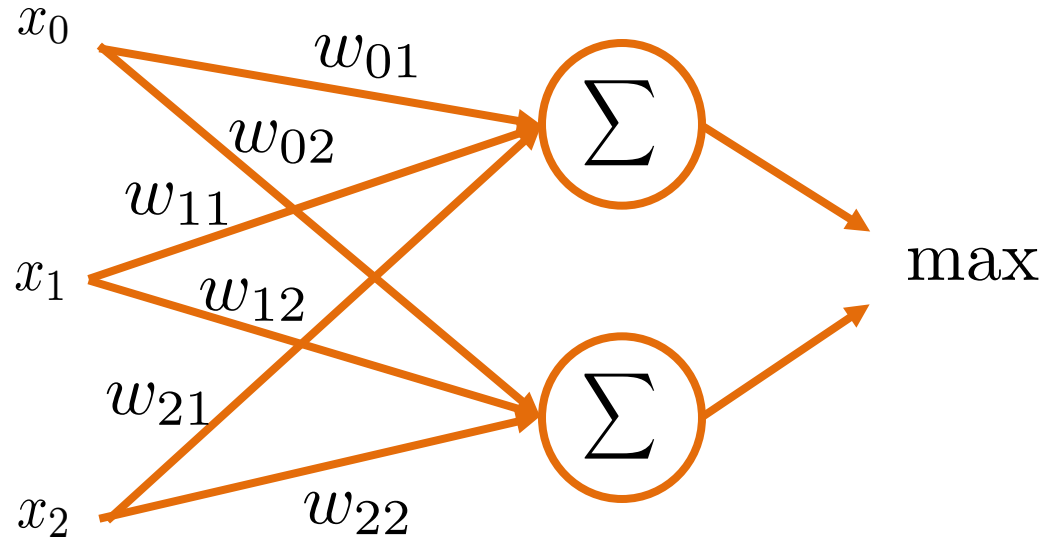


One more parameter
to backprop into

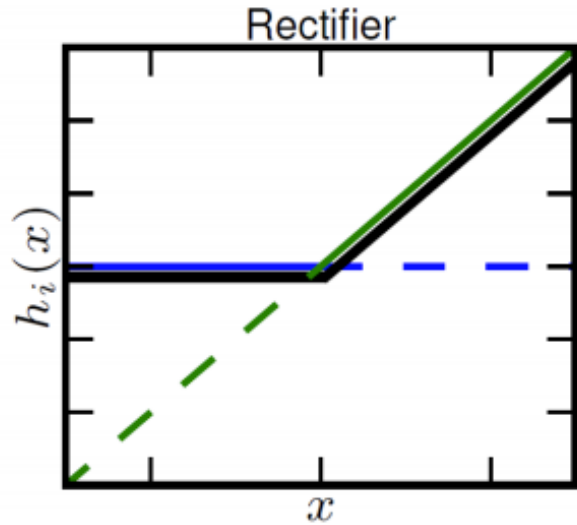


Does not die

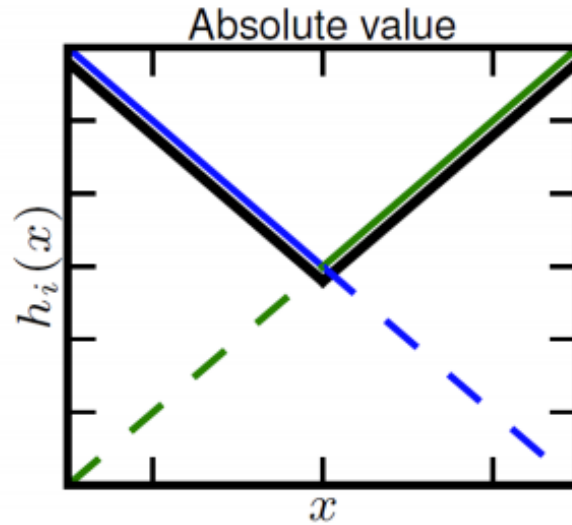
Maxout units



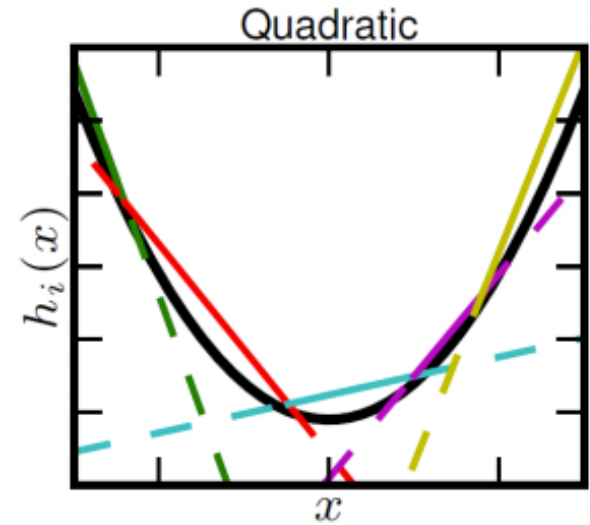
Maxout units



$k=2$



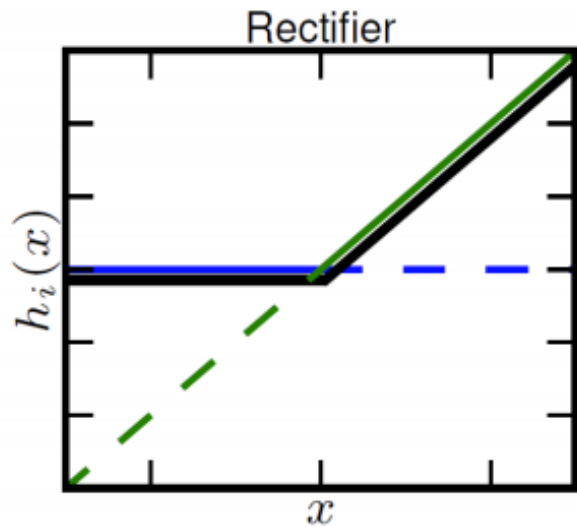
$k=2$



$k=5$

Piecewise linear approximation of
a convex function with N pieces

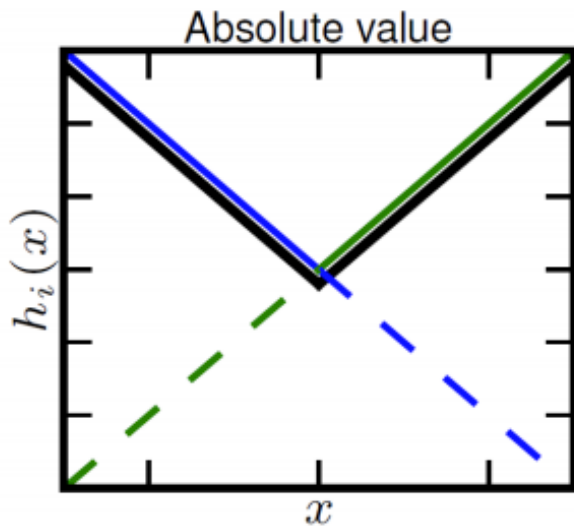
Maxout units



$k=2$



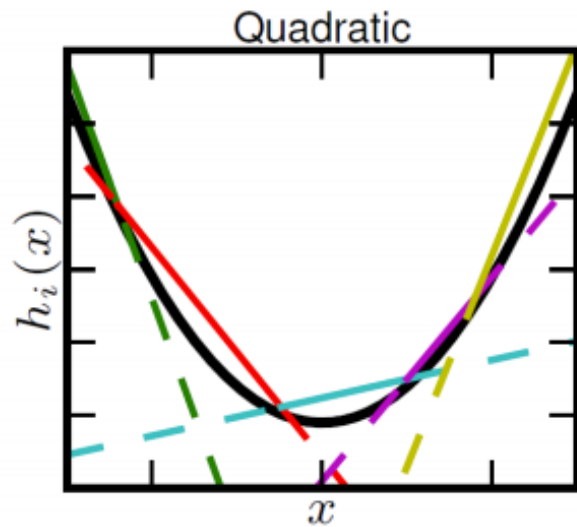
Generalization
of ReLUs



$k=2$



Linear
regimes



$k=5$



Does not
die

Does not
saturate



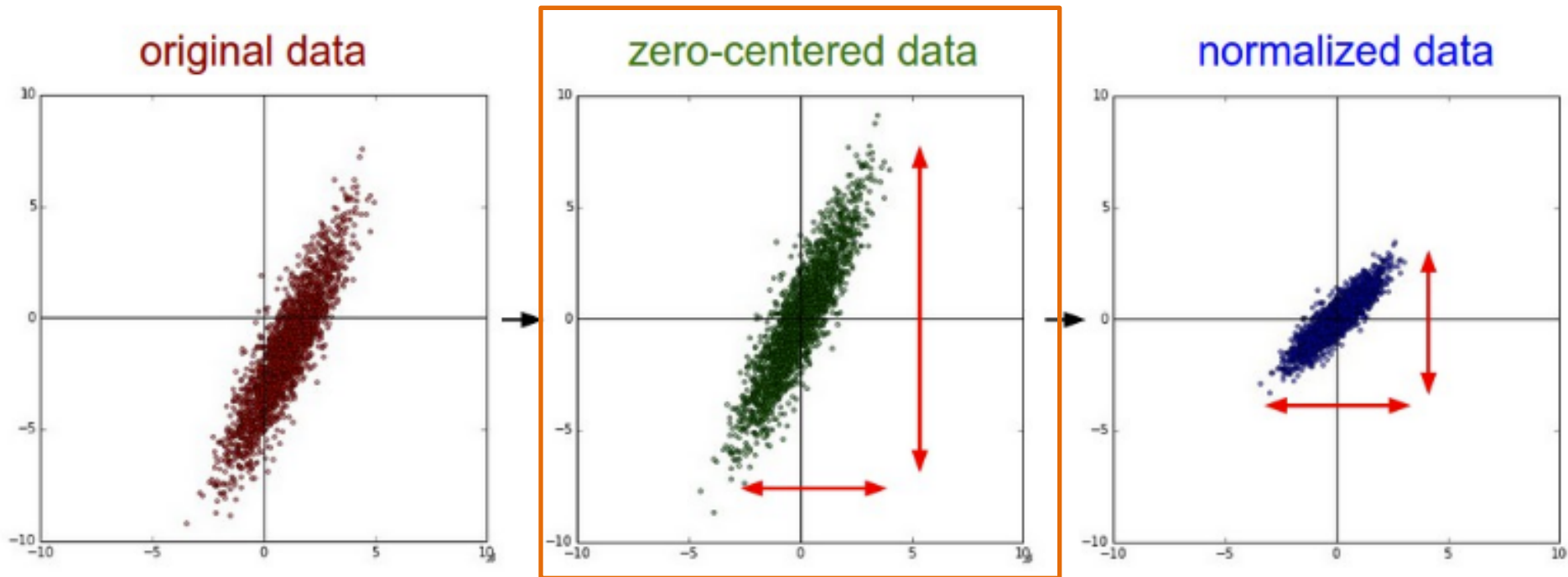
Increase of the number of parameters

Quick guide

- Sigmoid is not really used
- ReLU is the standard choice
- Second choice are the variants of ReLu or Maxout
- Recurrent nets will require tanh or similar

A quick word on data pre-processing

Data pre-processing

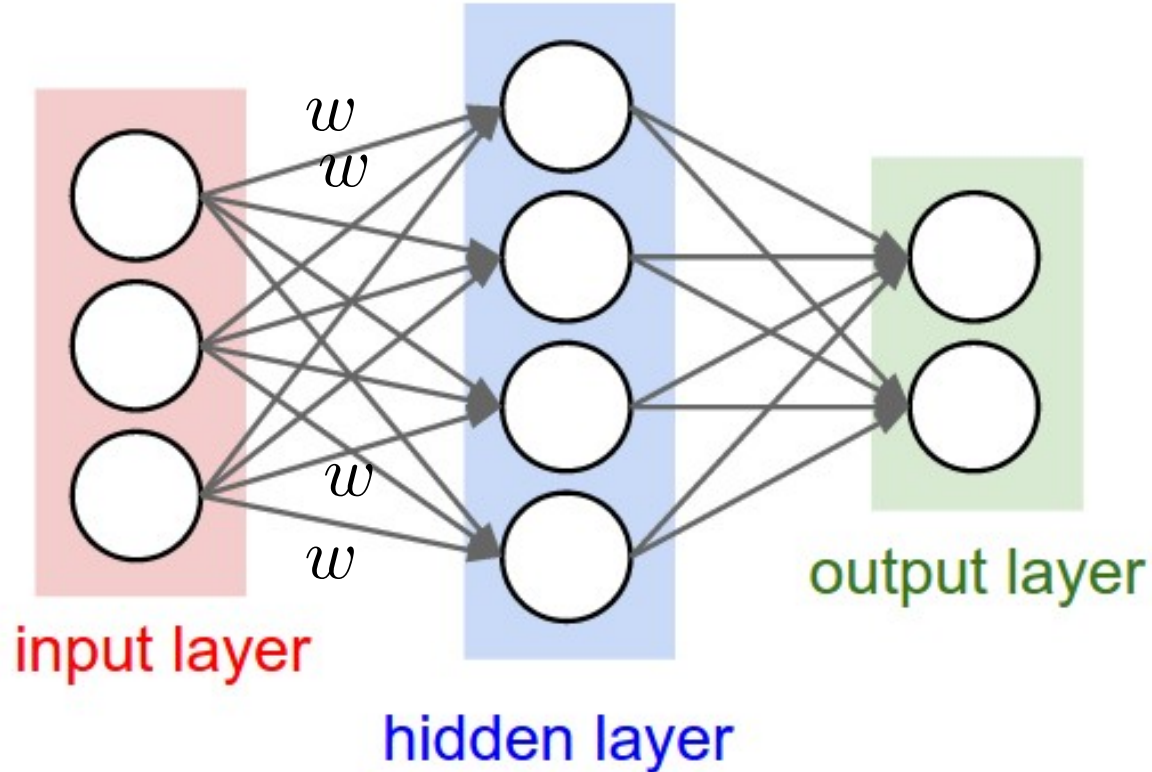


For images subtract the mean image (AlexNet) or per-channel mean (VGG-Net)

Weight initialization

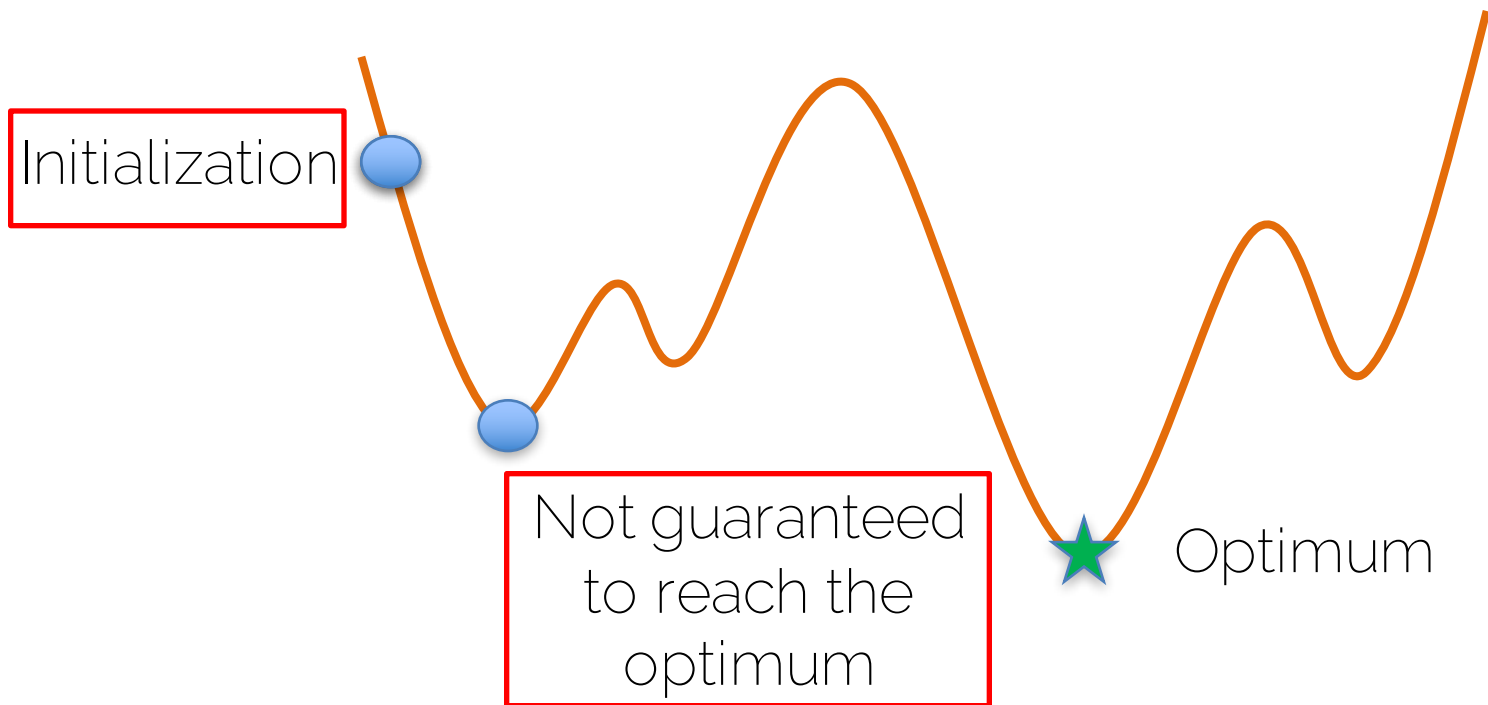
How do I start?

Forward



Initialization is extremely important

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



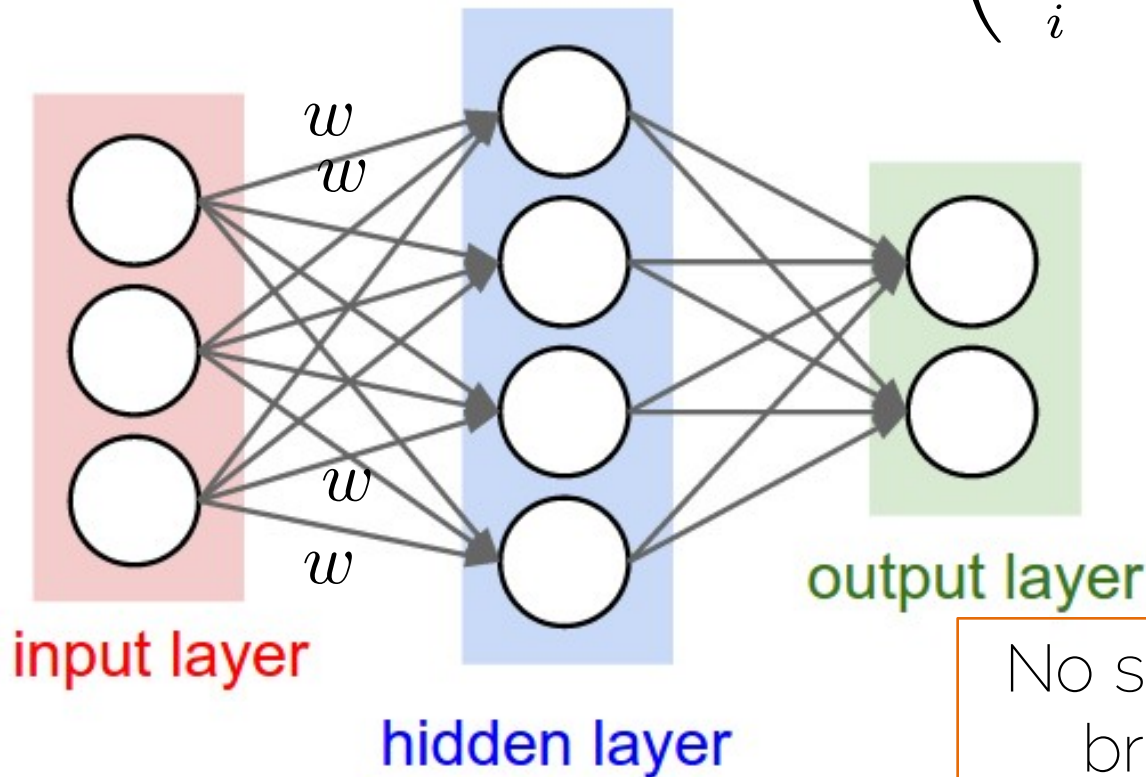
How do I start?

Forward



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

$w = 0$
What happens to the gradients?

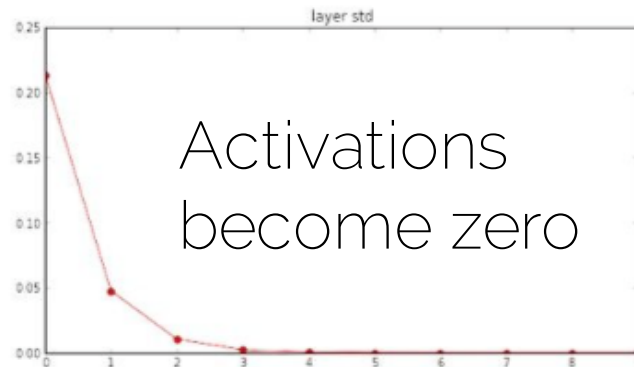
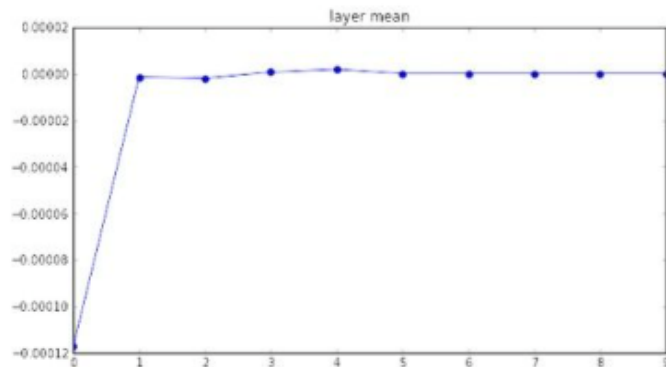


No symmetry breaking

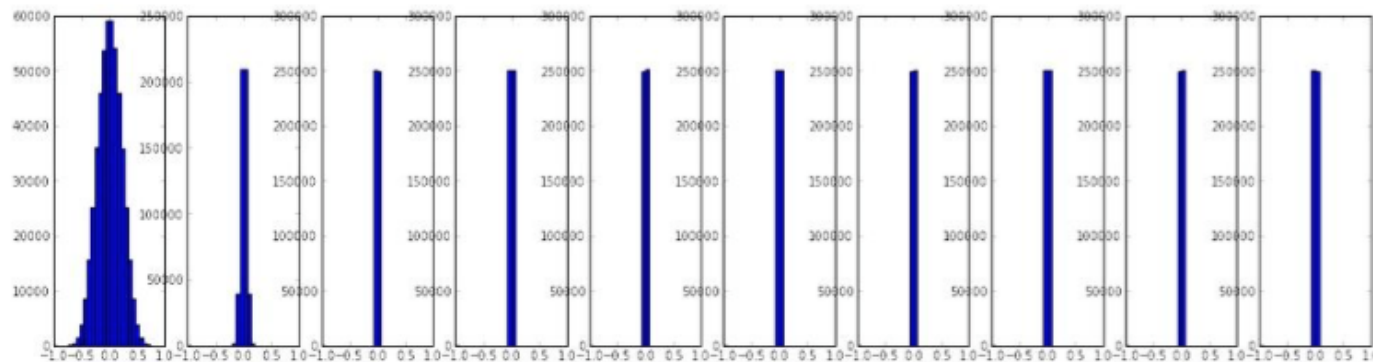
Small random numbers

- 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

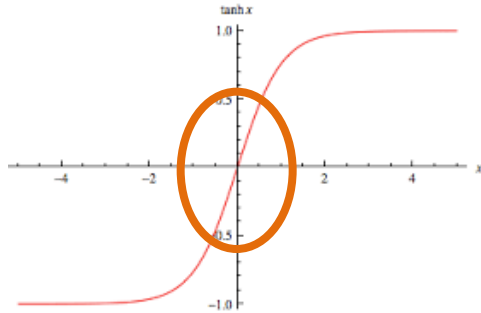
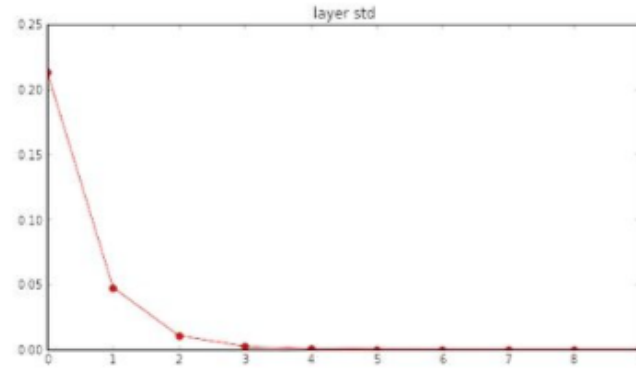
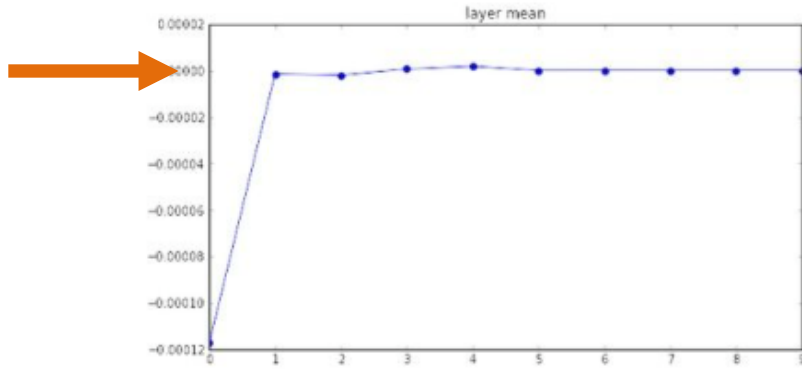


Last layer

Forward



Small random numbers

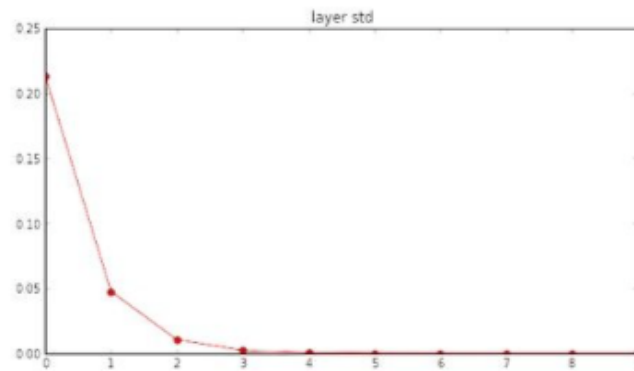
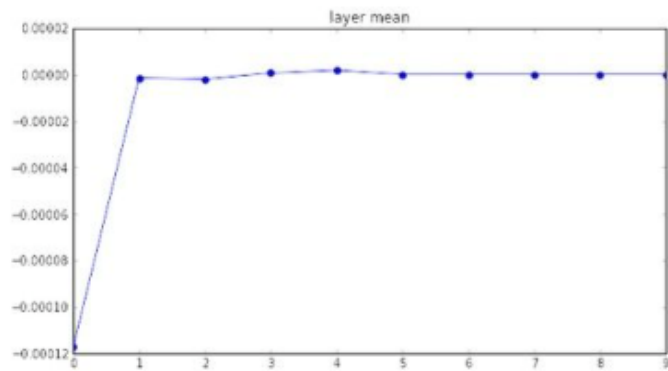


$$f \left(\sum_i^{\text{small}} w_i x_i + b \right)$$

Forward



Small random numbers



Gradients
vanish

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

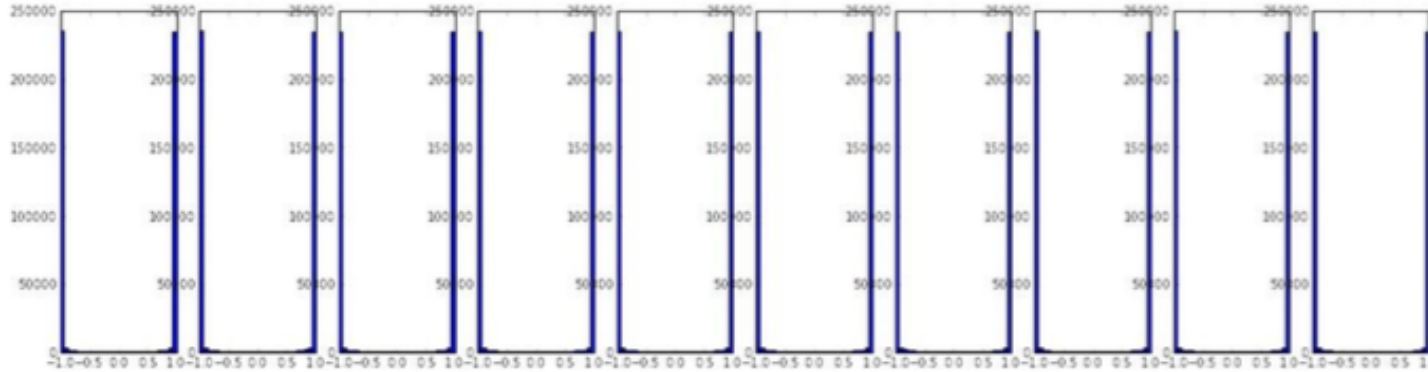
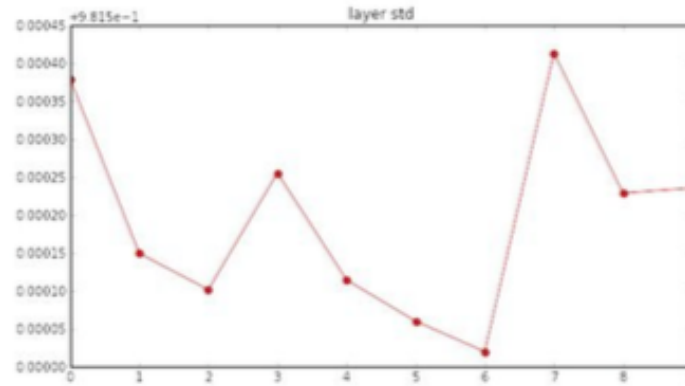
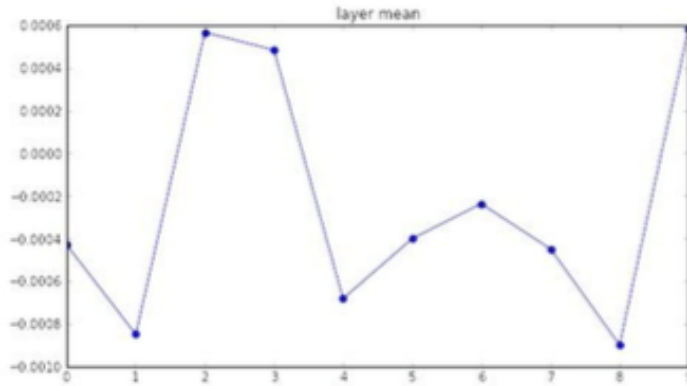
Backward



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

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?

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

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

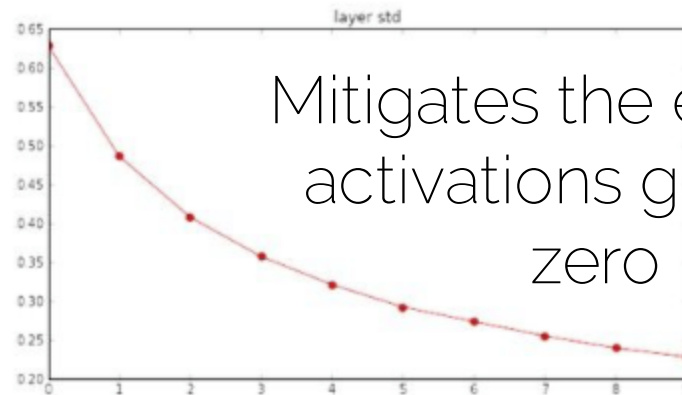
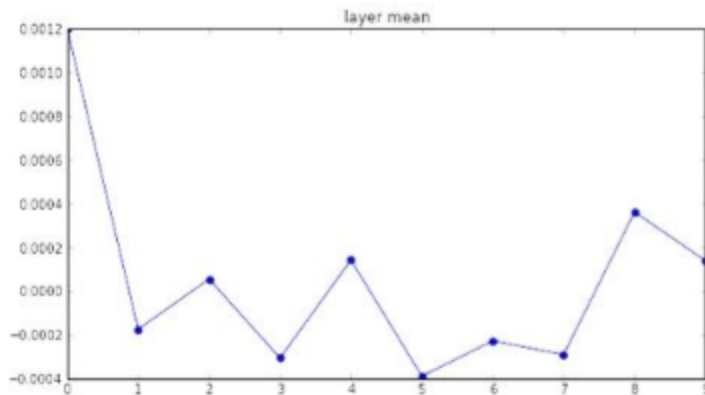
Xavier initialization

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

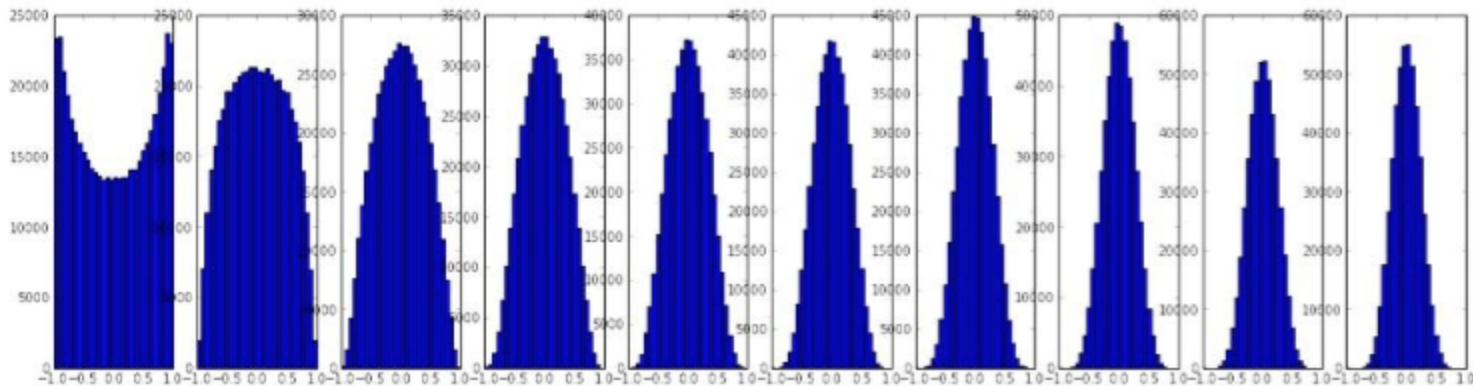
$$\frac{(n \text{Var}(w)) \text{Var}(x)}{1}$$

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

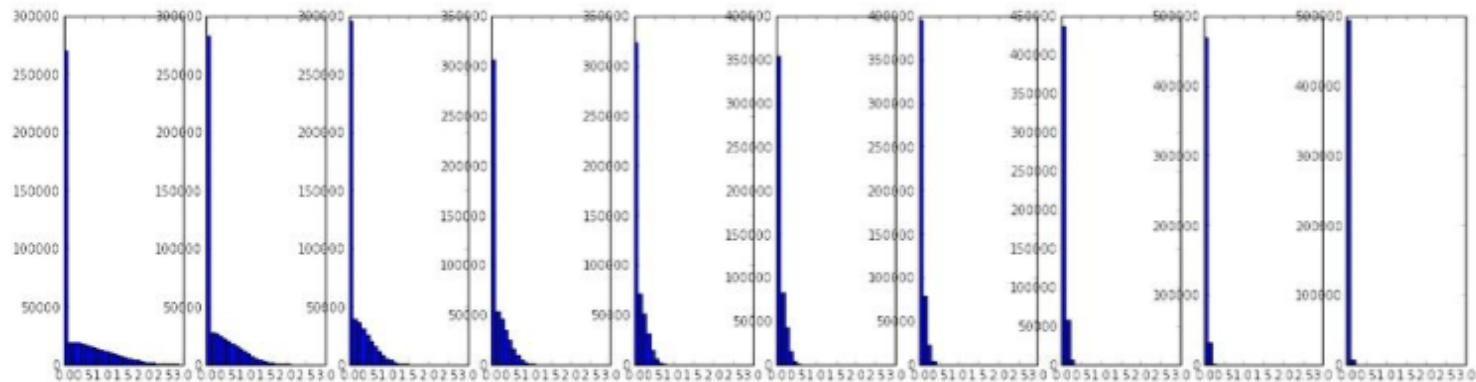
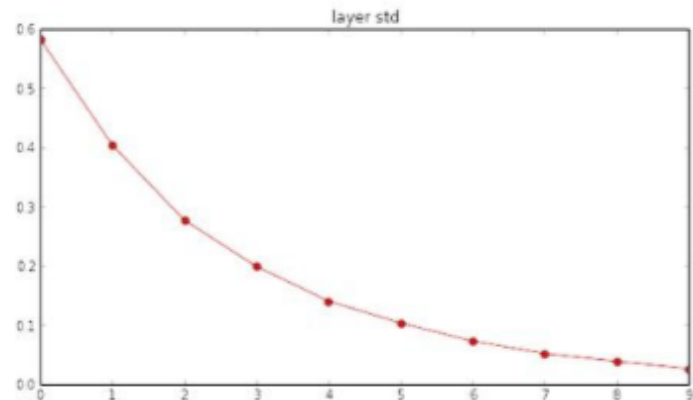
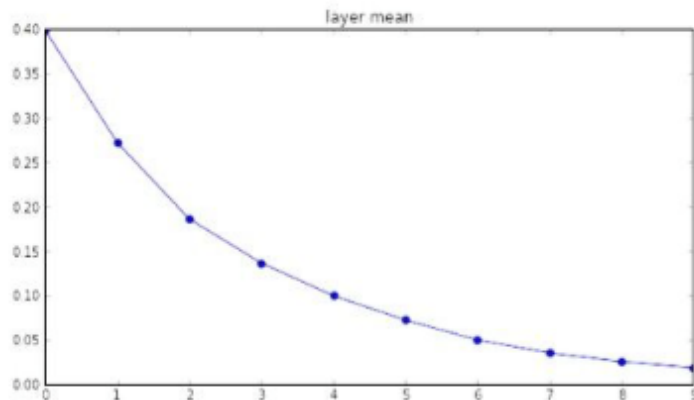
Xavier initialization



Mitigates the effect of
activations going to
zero

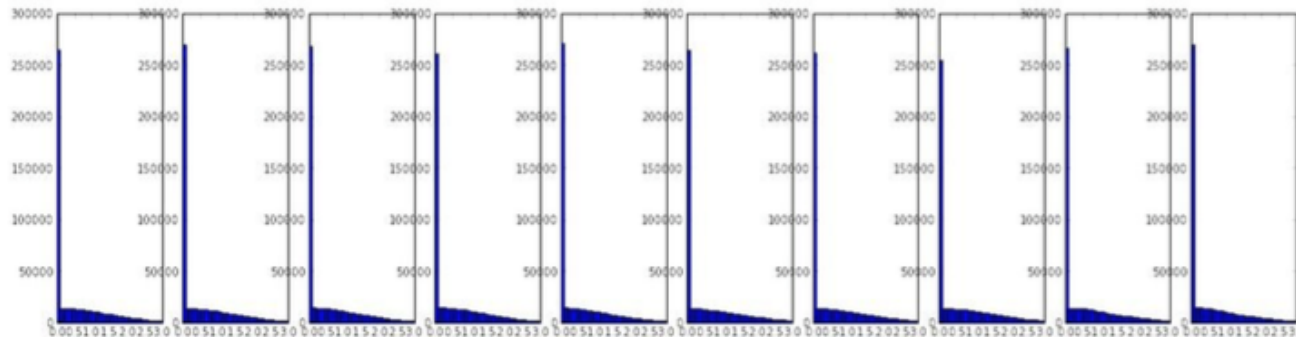
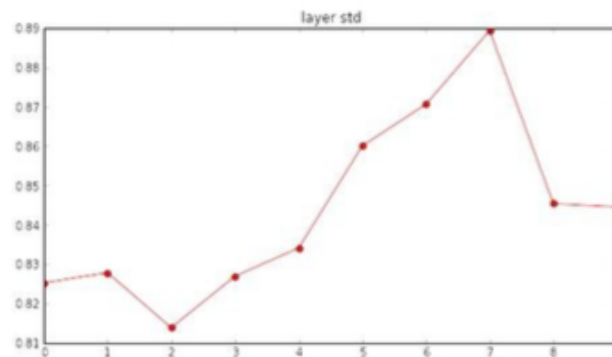
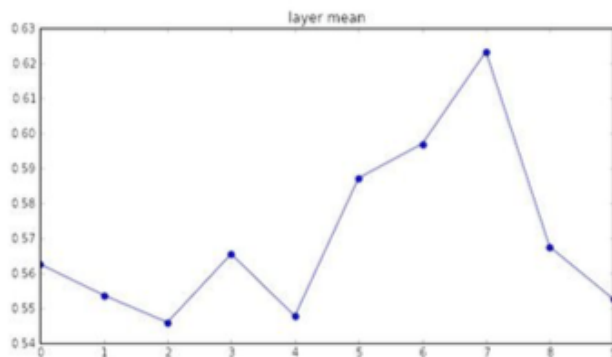


Xavier initialization with ReLU



ReLU kills half of the data

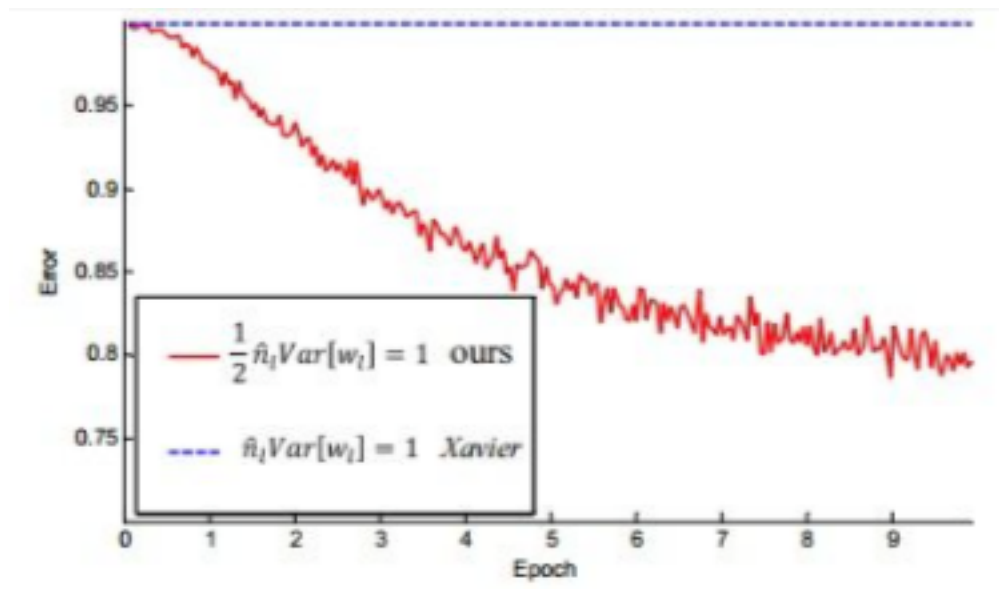
$$\text{Var}(w) = \frac{2}{n}$$



ReLU kills half of the data

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

It makes a huge difference!



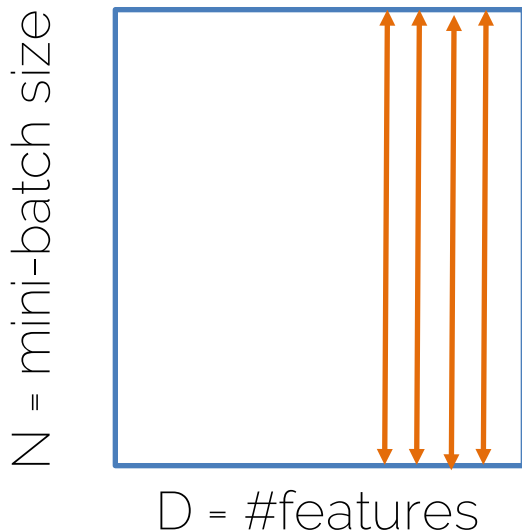
Tips and tricks

- Use ReLU and Xavier/2 initialization

Batch normalization

Batch normalization

- Wish: unit Gaussian activations
- Solution: let's do it

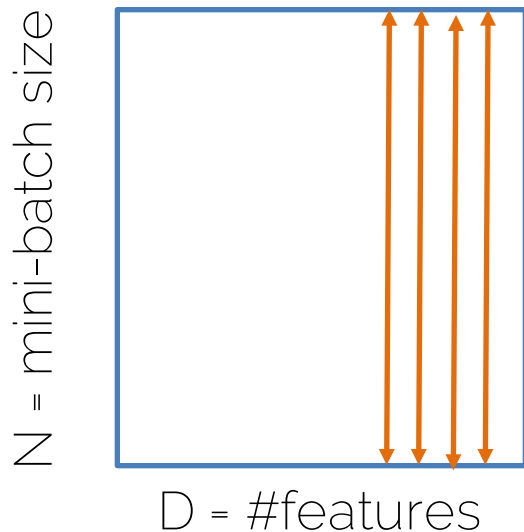


dimension

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

Batch normalization

- In each dimension of the features, you have a unit gaussian



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

An orange arrow points from the text "dimension" in the list item above to the term $x^{(k)}$ in the numerator of the equation.

Batch normalization

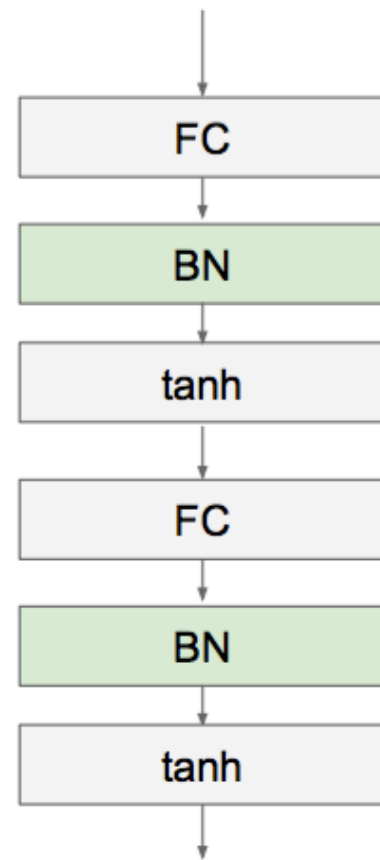
- In each dimension of the features, you have a unit Gaussian
- **Is it ok to treat dimensions separately?** Shown empirically that even if features are not decorrelated, convergence is still faster with this method

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

Differentiable function so
we can backprop
through it....

Batch normalization

- 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?



Batch normalization

- Normalize

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

- 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

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

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

BN for Exercise 2

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots x_m\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

Input: Network N with trainable parameters Θ ;
subset of activations $\{x^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference, $N_{\text{BN}}^{\text{inf}}$

1: $N_{\text{BN}}^{\text{tr}} \leftarrow N$ // Training BN network

2: **for** $k = 1 \dots K$ **do**

3: Add transformation $y^{(k)} = \text{BN}_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$ to $N_{\text{BN}}^{\text{tr}}$ (Alg. 1)

4: Modify each layer in $N_{\text{BN}}^{\text{tr}}$ with input $x^{(k)}$ to take $y^{(k)}$ instead

5: **end for**

6: Train $N_{\text{BN}}^{\text{tr}}$ to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$

7: $N_{\text{BN}}^{\text{inf}} \leftarrow N_{\text{BN}}^{\text{tr}}$ // Inference BN network with frozen parameters

8: **for** $k = 1 \dots K$ **do**

9: // For clarity, $x \equiv x^{(k)}, \gamma \equiv \gamma^{(k)}, \mu_{\mathcal{B}} \equiv \mu_{\mathcal{B}}^{(k)}$, etc.

10: Process multiple training mini-batches \mathcal{B} , each of size m , and average over them:

$$E[x] \leftarrow E_{\mathcal{B}}[\mu_{\mathcal{B}}]$$

$$\text{Var}[x] \leftarrow \frac{m}{m-1} E_{\mathcal{B}}[\sigma_{\mathcal{B}}^2]$$

11: In $N_{\text{BN}}^{\text{inf}}$, replace the transform $y = \text{BN}_{\gamma, \beta}(x)$ with

$$y = \frac{\gamma}{\sqrt{\text{Var}[x] + \epsilon}} \cdot x + \left(\beta - \frac{\gamma E[x]}{\sqrt{\text{Var}[x] + \epsilon}} \right)$$

12: **end for**

Algorithm 2: Training a Batch-Normalized Network

Administrative Things

- Next Thursday May 25th: No Lecture!
- Thursday June 1st :
 - More on Neural Networks 😊
- Tomorrow: Solution 1st exercise, presentation 2nd

