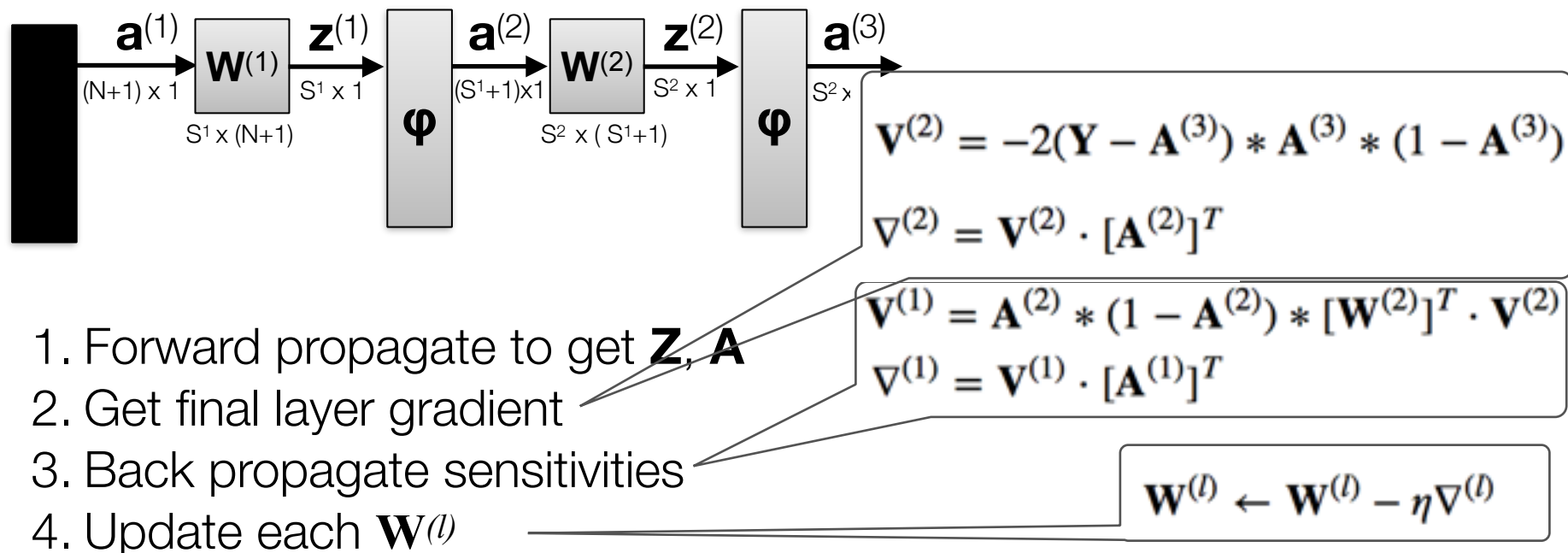


Lecture Notes for **Machine Learning in Python**

Professor Eric Larson
Town Hall + MLP History

Changing the Objective Function



• Self Test:

True or False: If we change the cost function, $J(\mathbf{W})$, we only need to update the final layer sensitivity calculation, $\mathbf{V}^{(2)}$, of the back propagation steps. The remainder of the algorithm is unchanged.

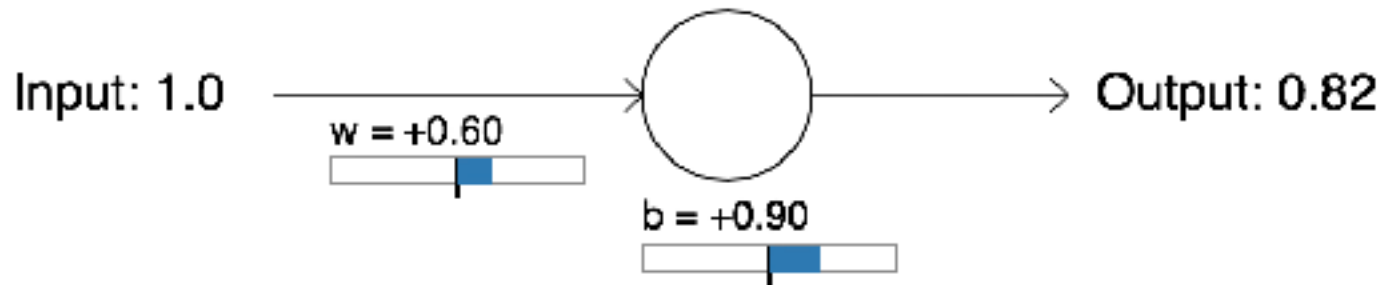
- A. True
- B. False

Practical Implementation of Architectures

- MSE

$$J(\mathbf{W}) = \sum_k^M (\mathbf{y}^{(k)} - \mathbf{a}^{(L)})^2$$

least squares objective,
tends to slow training initially



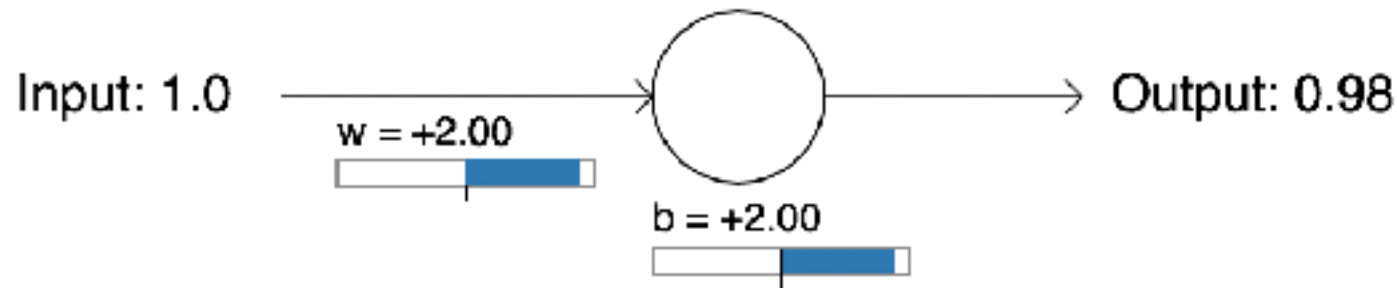
Run

Practical Implementation of Architectures

- MSE

$$J(\mathbf{W}) = \sum_k^M (\mathbf{y}^{(k)} - \mathbf{a}^{(L)})^2$$

least squares objective,
tends to slow training initially



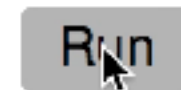
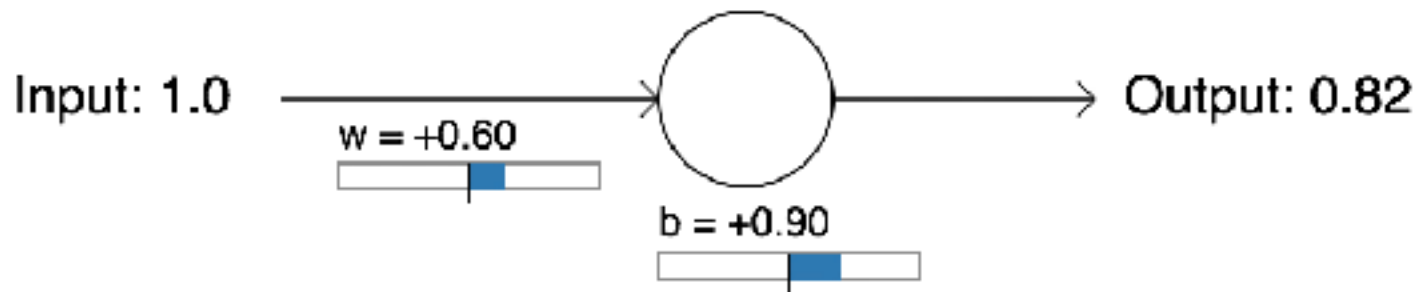
Run

Practical Implementation of Architectures

- Negative of MLE: **Binary Cross entropy**

$$J(\mathbf{W}) = - [\mathbf{y}^{(i)} \ln([\mathbf{a}^{(L+1)}]^{(i)}) + (1 - \mathbf{y}^{(i)}) \ln(1 - [\mathbf{a}^{(L+1)}]^{(i)})]$$

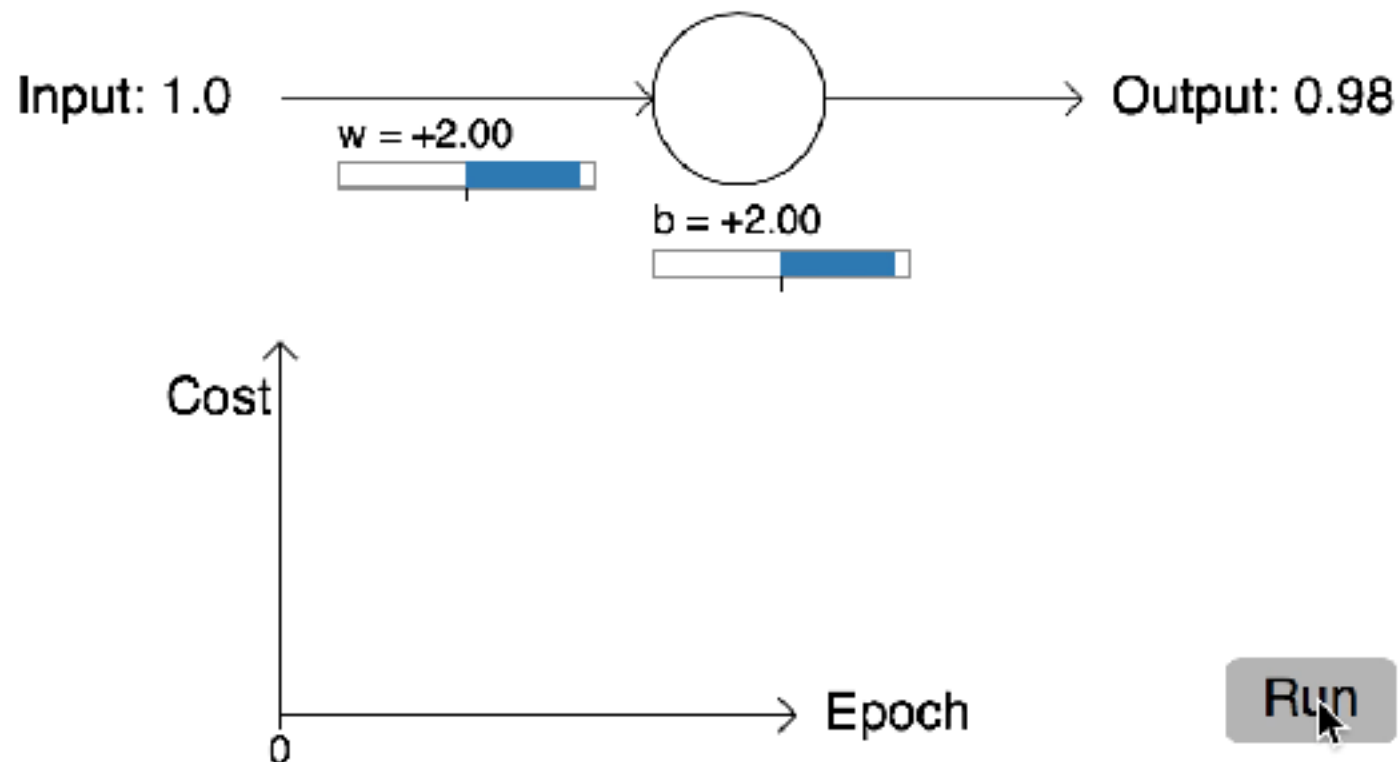
speeds up initial training



Practical Implementation of Architectures

- Negative of MLE: **Binary Cross entropy**

$$J(\mathbf{W}) = - \left[\mathbf{y}^{(i)} \ln([\mathbf{a}^{(L+1)}]^{(i)}) + (1 - \mathbf{y}^{(i)}) \ln(1 - [\mathbf{a}^{(L+1)}]^{(i)}) \right] \quad \text{speeds up initial training}$$



Practical Implementation of Architectures

- Negative of MLE: **Binary Cross entropy**

$$J(\mathbf{W}) = - \left[\mathbf{y}^{(i)} \ln([\mathbf{a}^{(L+1)}]^{(i)}) + (1 - \mathbf{y}^{(i)}) \ln(1 - [\mathbf{a}^{(L+1)}]^{(i)}) \right] \quad \text{speeds up initial training}$$

$$\left[\frac{\partial J(\mathbf{W})}{\mathbf{z}^{(L)}} \right]^{(i)}$$

$$\mathbf{V}^{(2)} = -2(\mathbf{Y} - \mathbf{A}^{(3)}) * \mathbf{A}^{(3)} * (1 - \mathbf{A}^{(3)}) \quad \text{old update}$$

Practical Implementation of Architectures

- Back to our old friend: **Cross entropy**

$$J(\mathbf{W}) = - \left[\mathbf{y}^{(i)} \ln([\mathbf{a}^{(L+1)}]^{(i)}) + (1 - \mathbf{y}^{(i)}) \ln(1 - [\mathbf{a}^{(L+1)}]^{(i)}) \right] \quad \begin{array}{l} \text{speeds up} \\ \text{initial training} \end{array}$$

$$\left[\frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(L)}} \right]^{(i)} = ([\mathbf{a}^{(L+1)}]^{(i)} - \mathbf{y}^{(i)})$$

$$\left[\frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(2)}} \right]^{(i)} = ([\mathbf{a}^{(3)}]^{(i)} - \mathbf{y}^{(i)})$$

$$\mathbf{V}^{(2)} = \mathbf{A}^{(3)} - \mathbf{Y}$$

new update

```
# vectorized backpropagation
V2 = (A3 - Y_enc) # <- this is only line t
V1 = A2 * (1 - A2) * (W2.T @ V2)

grad2 = V2 @ A2.T
grad1 = V1[1:,:] @ A1.T
```

$$\mathbf{V}^{(2)} = -2(\mathbf{Y} - \mathbf{A}^{(3)}) * \mathbf{A}^{(3)} * (1 - \mathbf{A}^{(3)}) \quad \text{old update}$$

bp-5

cross entropy



Practical Implementation of Architectures

Gradient when using cosine annealing with warm restarts learning rate scheduler

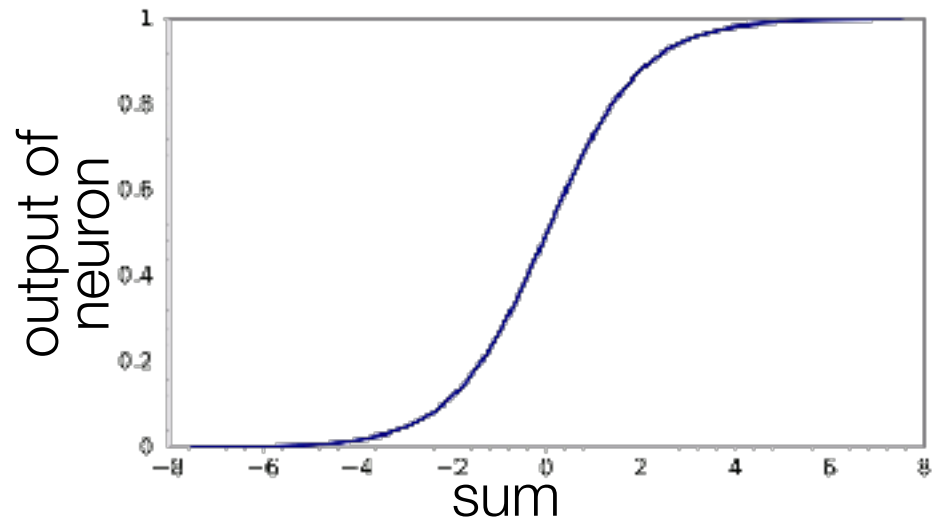


Formative Self Test

- for adding Gaussian distributions, variances add together

$\mathbf{a}^{(L+1)} = \phi(\mathbf{W}^{(L)}\mathbf{a}^{(L)})$ assume each element of \mathbf{a} is Gaussian

- If you initialized the weights, \mathbf{W} , with too large variance, you would expect the output of the neuron, $\mathbf{a}^{(L+1)}$, to be:
 - A. saturated to “1”
 - B. saturated to “0”
 - C. could either be saturated to “0” or “1”
 - D. would not be saturated

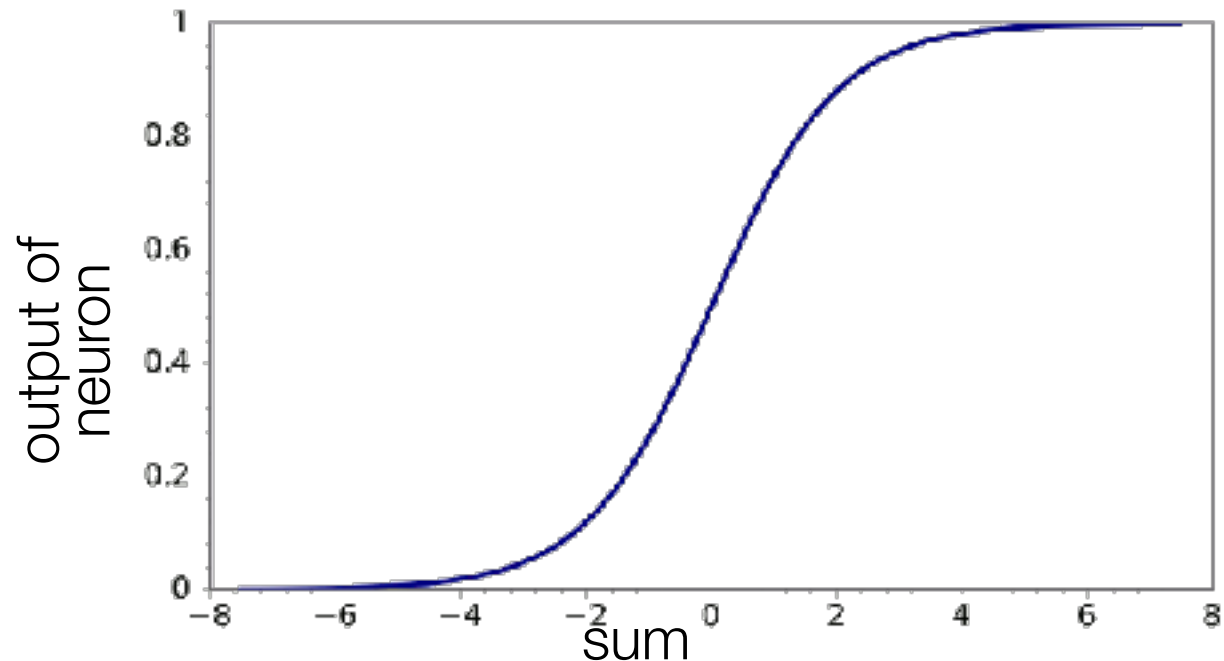


Formative Self Test

- for adding Gaussian distributions, variances add together

$\mathbf{a}^{(L+1)} = \phi(\mathbf{W}^{(L)}\mathbf{a}^{(L)})$ assume each element of \mathbf{a} is Gaussian

- What is the derivative of a saturated sigmoid neuron?
 - A. zero
 - B. one
 - C. $a * (1-a)$
 - D. it depends



Practical Implementation of Architectures

- **Weight initialization**

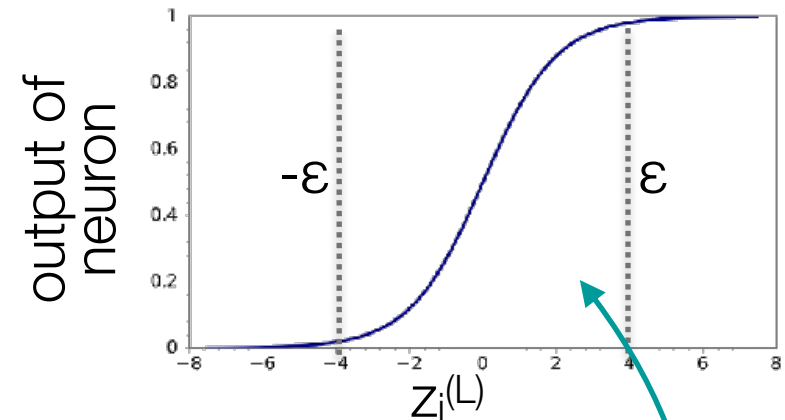
- try not to **saturate** your neurons right away!

$$\mathbf{a}^{(L+1)} = \phi(\mathbf{z}^{(L)})$$

$$\mathbf{z}^{(L)} = \mathbf{W}^{(L)} \mathbf{a}^{(L)}$$



each row is summed before sigmoid



want each $z^{(L)}$ to be between $-\epsilon < \Sigma < \epsilon$ for no saturation

solution: squash initial weights magnitude

- one choice: each element of **W** selected from a Gaussian with **zero mean** and **specific standard deviation**

$$w_{ij}^{(L)} \leftarrow \mathcal{N}\left(0, \sqrt{\frac{1}{n^{(L)}}}\right)$$

For a sigmoid, want $-\epsilon < z_i^{(L)} < \epsilon$

$\epsilon=4$

More Weight Initialization

Understanding the difficulty of training deep feedforward neural networks

Xavier Glorot

JMLR 2010

Yoshua Bengio

DIRO, Université de Montréal, Montréal, Québec, Canada

Goal: We should not saturate **feedforward** or **back propagated** variance

Relate variance of current layer to variance in z , so $\sigma(z_i^{(L)})$ isn't saturated

try not to saturate z $z_i^{(L)} = \sum_j^{n^{(L)}} w_{ij} a_j^{(L)}$ *break down feed forward by each multiply*

$$\text{Var}[z_i^{(L)}] = \sum_j^{n^{(L)}} \underbrace{E[w_{ij}]^2 \text{Var}[a_j^{(L)}] + \text{Var}[w_{ij}] E[a_j^{(L)}]^2}_{0, \text{ if uncorrelated}} + \underbrace{\text{Var}[w_{ij}] \text{Var}[a_j^{(L)}]}_{\approx 1}$$

Want to keep $\text{Var}[\cdot] \sim 1$ *assume i.i.d. expand variance calc*

$$\text{Var}[z_i^{(L)}] = 4 = n^{(L)} \text{Var}[w_{ij}] \text{Var}[a_j^{(L)}]$$

Similar for back prop.

$$\text{Var}[v_i^{(L)}] = n^{(L+1)} \text{Var}[w_{ij}] \text{Var}[v_j^{(L+1)}]$$

$$w_{ij}^{(L)} \approx \mathcal{N}\left(0, 4 \cdot \sqrt{\frac{1}{n^{(L)}}}\right) \quad \left| \quad w_{ij}^{(L)} \approx \mathcal{N}\left(0, 4 \cdot \sqrt{\frac{1}{n^{(L+1)}}}\right) \quad \left| \quad w_{ij}^{(L)} \approx \mathcal{N}\left(0, 4 \cdot \sqrt{\frac{2}{n^{(L)} + n^{(L+1)}}}\right)$$

forward from data backward from sensitivity compromise

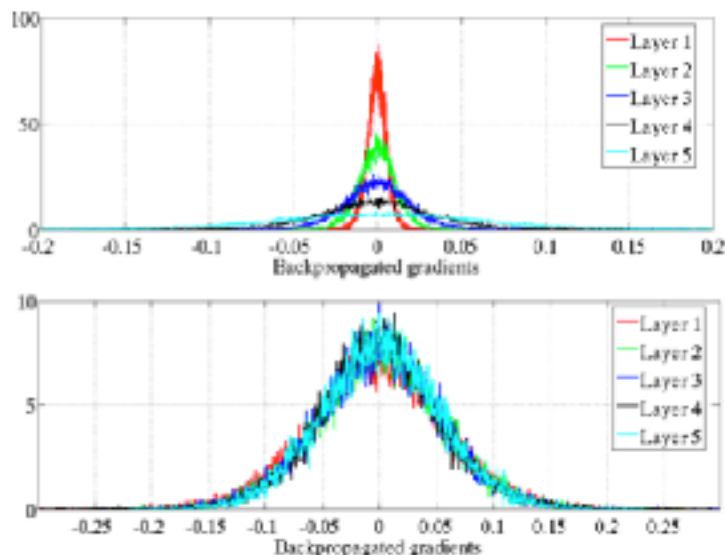
More Weight Initialization

Understanding the difficulty of training deep feedforward neural networks

Xavier Glorot

DIRO, Université de Montréal, Montréal, Québec, Canada

Yoshua Bengio



Starting gradient histograms
per layer
standard normalization

Starting gradient histograms
per layer
Glorot normalization

Figure 7: *Back-propagated gradients normalized histograms with hyperbolic tangent activation, with standard (top) vs normalized (bottom) initialization. Top: 0-peak decreases for higher layers.*

Glorot and He Initialization

We have solved this assuming the activation output is in the range -4 to 4 (for a sigmoid) and assuming that x is distributed Gaussian

This range, epsilon, is different depending on the activation and assuming Gaussian or Uniform

Uniform

Gaussian

Tanh

$$w_{ij}^{(L)} = \sqrt{\frac{6}{n^{(L)} + n^{(L+1)}}}$$

$$w_{ij}^{(L)} = \sqrt{\frac{2}{n^{(L)} + n^{(L+1)}}}$$

Sigmoid

$$w_{ij}^{(L)} = 4\sqrt{\frac{6}{n^{(L)} + n^{(L+1)}}}$$

$$w_{ij}^{(L)} = 4\sqrt{\frac{2}{n^{(L)} + n^{(L+1)}}}$$

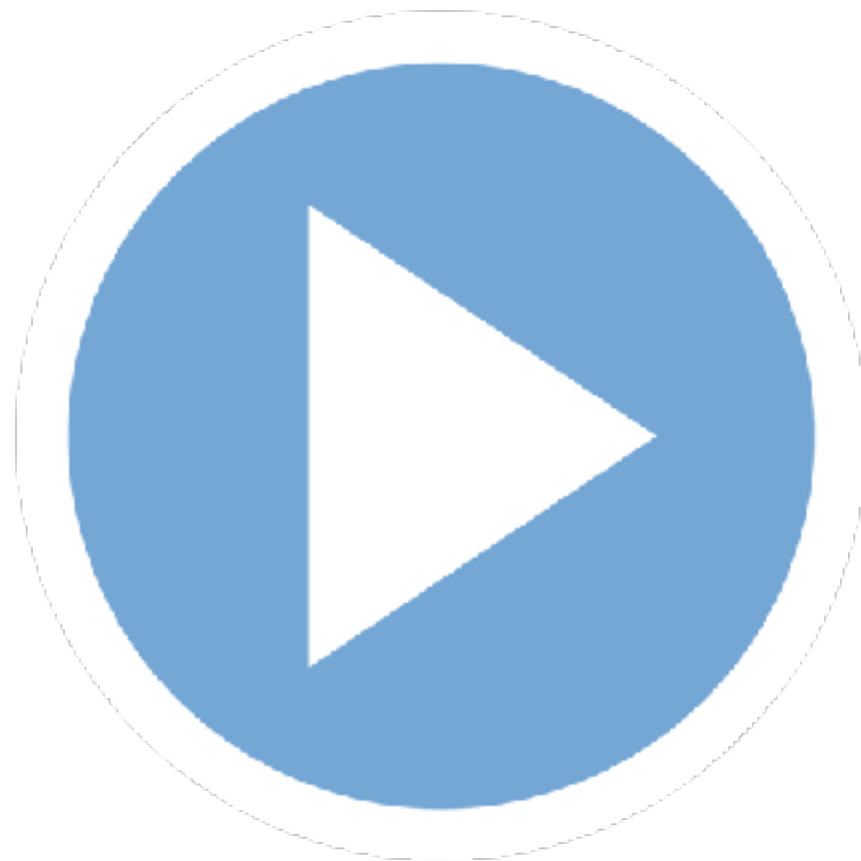
ReLU

$$w_{ij}^{(L)} = \sqrt{2}\sqrt{\frac{6}{n^{(L)} + n^{(L+1)}}}$$

$$w_{ij}^{(L)} = \sqrt{2}\sqrt{\frac{2}{n^{(L)} + n^{(L+1)}}}$$

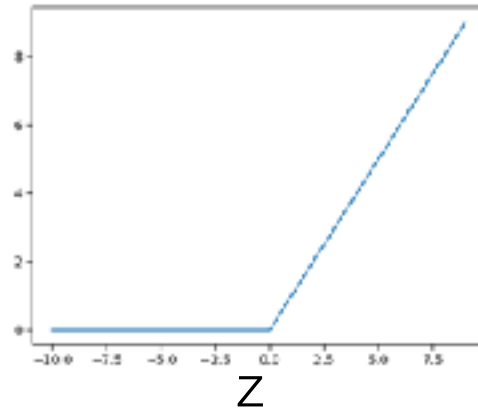
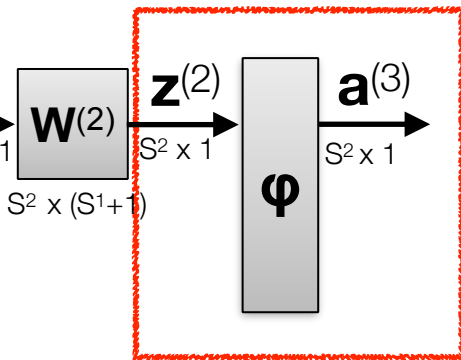
Summarized by Glorot and He

Smarter Weight Initialization



Practical Implementation of Architectures

- A new nonlinearity: **rectified linear units**



$$\phi(z) = \begin{cases} z, & \text{if } z > 0 \\ 0, & \text{else} \end{cases}$$

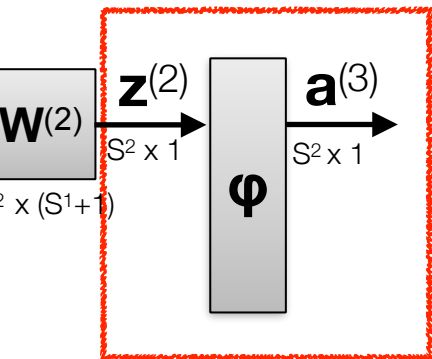
it has the advantage of **large gradients** and **extremely simple** derivative

$$\nabla \phi(z) = \begin{cases} 1, & \text{if } z > 0 \\ 0, & \text{else} \end{cases}$$



ReLU Nonlinearities
Important for deep networks

Other Activation Functions



- Sigmoid Weighted Linear Unit **SiLU**
 - also called Swish
- Mixing of sigmoid, σ , and ReLU

$$\varphi(z) = z \cdot \sigma(z)$$

$$\frac{\partial \varphi(z)}{\partial z} = \varphi(z) + \sigma(z)[1 - \varphi(z)]$$

$$= a^{(l+1)} + \sigma(z^{(l)}) \cdot [1 - a^{(l+1)}]$$

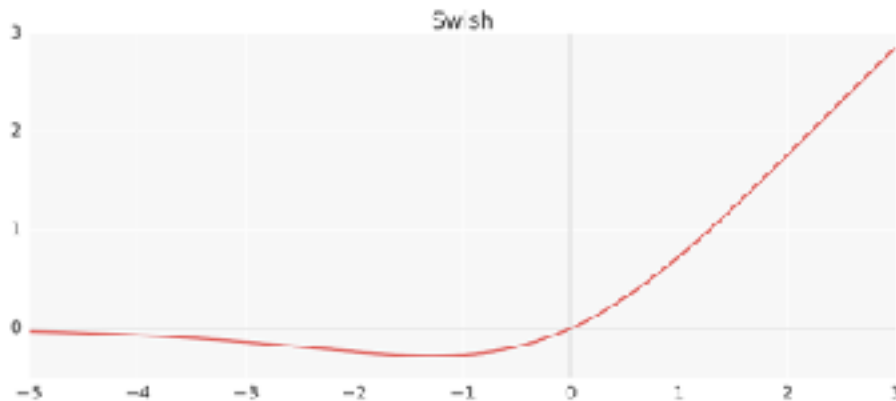


Figure 1: The Swish activation function.

Ramachandran P, Zoph B, Le QV. Swish: a Self-Gated Activation Function. arXiv preprint arXiv:1710.05941. 2017 Oct 16

Elfwing, Stefan, Eiji Uchibe, and Kenji Doya. "Sigmoid-weighted linear units for neural network function approximation in reinforcement learning." Neural Networks (2018).

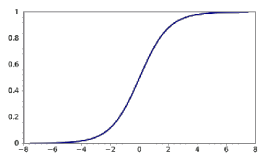
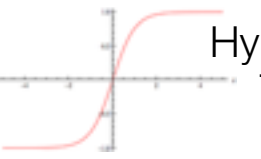
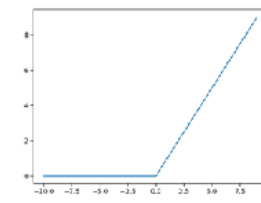
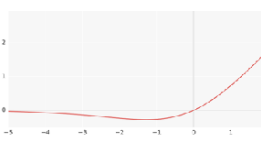
Derivative Calculation:

$$= \sigma(x) + x \cdot \sigma(x)(1 - \sigma(x))$$

$$= \sigma(x) + x \cdot \sigma(x) - x \cdot \sigma(x)^2$$

$$= x \cdot \sigma(x) + \sigma(x)(1 - x \cdot \sigma(x))$$

Activations Summary

	Definition	Derivative	Weight Init (Uniform Bounds)
 <p>Sigmoid</p>	$\phi(z) = \frac{1}{1 + e^{-z}}$	$\nabla \phi(z) = a(1 - a)$	$w_{ij}^{(L)} = 4\sqrt{\frac{6}{n^{(L)} + n^{(L+1)}}}$
 <p>Hyperbolic Tangent</p>	$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$	$\nabla \phi(z) = \frac{4}{(e^z + e^{-z})^2}$	$w_{ij}^{(L)} = \sqrt{\frac{6}{n^{(L)} + n^{(L+1)}}}$
 <p>ReLU</p>	$\phi(z) = \begin{cases} z, & \text{if } z > 0 \\ 0, & \text{else} \end{cases}$	$\nabla \phi(z) = \begin{cases} 1, & \text{if } z > 0 \\ 0, & \text{else} \end{cases}$	$w_{ij}^{(L)} = \sqrt{2}\sqrt{\frac{6}{n^{(L)} + n^{(L+1)}}}$
 <p>SiLU</p>	$\phi(z) = \frac{z}{1 + e^{-z}}$	$\nabla \phi(z) = a + \frac{(1 - a)}{1 + e^{-z}}$	$w_{ij}^{(L)} = \sqrt{2}\sqrt{\frac{6}{n^{(L)} + n^{(L+1)}}}$

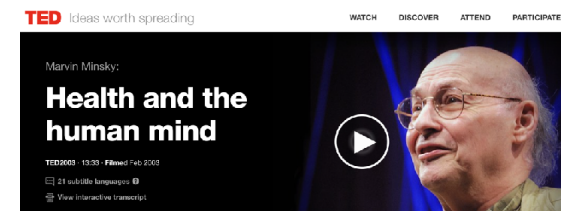
Practical Details

- Neural networks can separate any data through multiple layers. The true realization of Rosenblatt:

"Given an elementary α -perceptron, a stimulus world W , and any classification $C(W)$ for which a solution exists; let all stimuli in W occur in any sequence, provided that each stimulus must reoccur in finite time; then beginning from an arbitrary initial state, an error correction procedure will always yield a solution to $C(W)$ in finite time..."



- **Universality:** No matter what function we want to compute, we know that there is a neural network which can do the job.



- One nonlinear hidden layer with an output layer can perfectly train any problem with enough data, but might just be memorizing...
 - ... it might be better to have even more layers for decreased computation and generalizability

End of Session

- Next Time: Final Flipped Module!
- Then: Deep Learning in Keras

Back Up Slides