

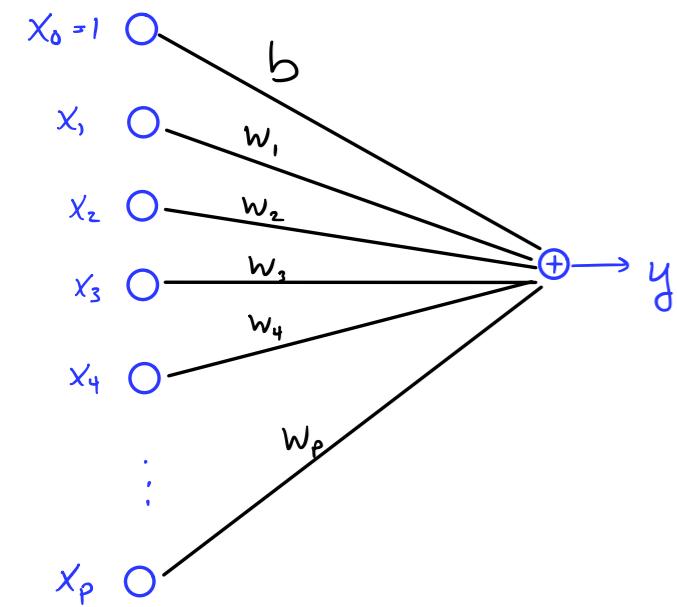
Neural Networks Part I

STAT 37710 / CMSC 35400
Rebecca Willett and Yuxin Chen

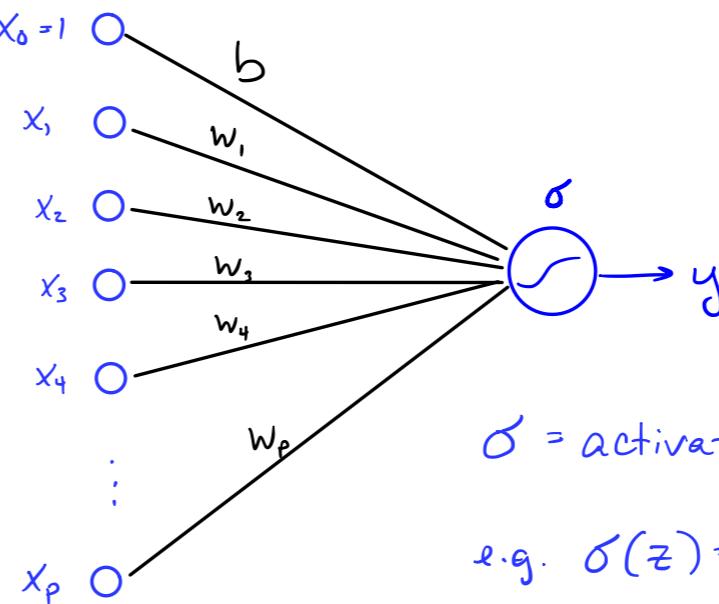
Introduction

A Simple Neural Network

$$y = \underline{x}^\top \underline{w} + b$$



$$y = \sigma(\underline{x}^\top \underline{w} + b)$$



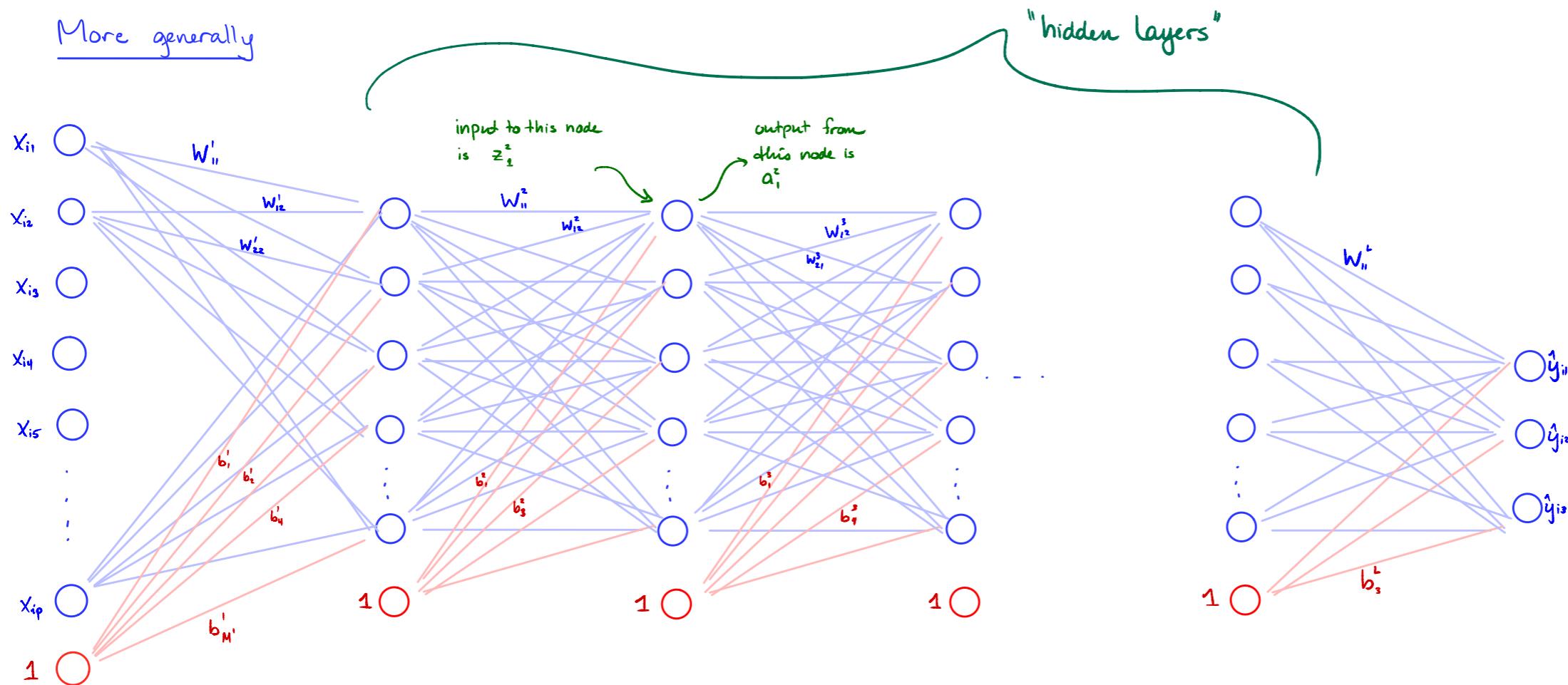
σ = activation function

e.g. $\sigma(z) = \max(0, z) = \text{ReLU}$

$$\sigma(z) = \frac{1}{1+e^{-z}} = \text{logistic}$$

$$\sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \tanh(z)$$

Multi-Layered Networks



$$\underline{x} \quad W, b'$$

$$= a^0 \quad z^1 = W^1 a^0 + b^1$$

$$a^1 = \sigma(z^1)$$

$$W^2, b^2$$

$$z^2 = W^2 a^1 + b^2$$

$$a^2 = \sigma(z^2)$$

$$W^3, b^3$$

$$z^3 = W^3 a^2 + b^3$$

$$a^3 = \sigma(z^3)$$

...

$$W^L, b^L$$

$$z^L = W^L a^{L-1} + b^L = \hat{y}$$

let $M^l = \# \text{ of blue nodes in layer } l$

$$W^l \sim M^l \times M^{l-1},$$

$$S^l \sim M^l \times 1,$$

$$a^l, z^l \sim M^l \times 1,$$

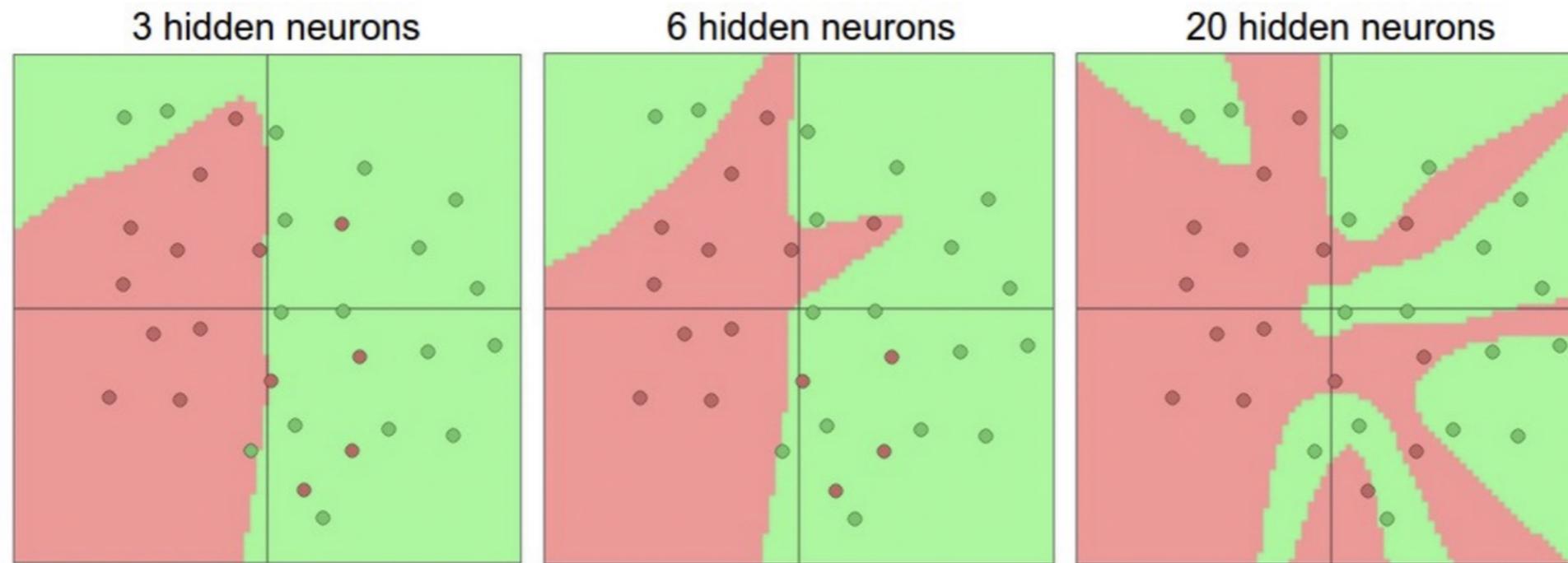
$$b^l \sim M^l \times 1$$

$$\hat{y}_{i,k} = W_k^L \sigma(W_{k-1}^L \sigma(\dots \sigma(W_1^L \sigma(W^0 x_i + b^0) + b^1) + \dots) + b^{L-1}) + b^L$$

Representational Power

- Neural network with at **least one hidden layer** is a universal approximator (can represent any function).

Proof in: Approximation by Superpositions of Sigmoidal Function, Cybenko, [paper](#)



- The capacity of the network increases with more hidden units and more hidden layers

Training

Let $\Theta : (W^1, \dots, W^L, b^1, \dots, b^L)$ denote all the weights in the neural network.

Our task in training is to use the training samples (x_i, y_i) , $i=1, \dots, n$, to select a $\hat{\Theta}$ that minimizes the training loss.

That is, let $f_\Theta(x)$ be the output of the network with weights Θ and input x :

$$f_\Theta(x) = W_k^L \sigma(W_{k-1}^{L-1} \cdot \sigma(W_{k-2}^2 \sigma(W^1 x + b^1) + b^2) + \dots) + b^L$$

We want to find $\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, f_\Theta(x_i))$, where

$L(\hat{y}, y)$ might be $\|\hat{y} - y\|^2$ or $\exp(-y\hat{y})$ for binary classification.

For K-class classification, map $y_i \rightarrow z_i \in \{0, 1\}^K$, where $z_{ik} = \begin{cases} 1 & \text{if } y_i = k \\ 0 & \text{otherwise} \end{cases} = [0 \ 0 \ \dots \ 1 \ \dots \ 0]^T$

Now we want to predict z_i for each x_i . kth position

Let \hat{y}_{ik} be kth output of network with input x_i , and say $P(\text{class}=k|x_i, \Theta) = \frac{\exp(\hat{y}_{ik})}{\sum_{j=1}^K \exp(\hat{y}_{ij})}$

Then $L_i(\Theta) = -\sum_{k=1}^K z_{ik} \log P(\text{class}=k|x_i, \Theta)$

"cross-entropy loss"

Stochastic Gradient Descent

Optimizing network weights

Write $L_i(\theta) := L(y_i, f_\theta(x_i))$

and $\hat{R}(\theta) := \frac{1}{n} \sum_{i=1}^n L_i(\theta)$.

Then $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \hat{R}(\theta)$

Our basic tool for this is

stochastic gradient descent.

Stochastic Gradient Descent:

@ iteration t , choose $i_t \in \{1, 2, \dots, n\}$

$$\theta^{(t+1)} = \theta^{(t)} - \tau \nabla L_{i_t}(\theta^{(t)})$$

- each iteration easier/faster to compute

- need more iterations

- choose i_t uniformly at random

Ex: $\hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \langle \underline{x}_i, \theta \rangle)^2 + \lambda \|\theta\|_2^2$

$$L_i(\theta) = (y_i - \langle \underline{x}_i, \theta \rangle)^2 + \lambda \|\theta\|_2^2$$

$$\text{check: } \frac{1}{n} \sum_{i=1}^n L_i(\theta) = \hat{R}(\theta)$$

$$\nabla L_i(\theta) = -2(y_i - \langle \underline{x}_i, \theta \rangle) \underline{x}_i + 2\lambda \theta$$

SGD: $\theta^{(t+1)} = \theta^{(t)} + 2\tau (y_{i_t} - \langle \underline{x}_{i_t}, \theta^{(t)} \rangle) \underline{x}_{i_t} - 2\tau \lambda \theta^{(t)}$

Then Gradient Descent =

$$\theta^{(t+1)} = \theta^{(t)} - \frac{\tau}{n} \sum_{i=1}^n \nabla L_i(\theta^{(t)})$$

SGD on simple network

for $\hat{y} = \sigma(\underline{x}^\top \underline{w} + b)$ and $\sigma(z) = \frac{1}{1+e^{-z}}$, how do we learn weights?

$$\text{loss } R(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\underbrace{\langle \underline{x}_i, \underline{w} \rangle}_{z_i} + b))^2$$



$$\hat{y}_i = \sigma(\underline{x}_i^\top \underline{w} + b) = \sigma(z_i)$$

$$z_i = \underline{x}_i^\top \underline{w} + b$$

$$\text{aside : } \frac{d\sigma}{dz} = \sigma'(z)$$

$$= \sigma(z)(1 - \sigma(z))$$

learn \underline{w}, b via Stochastic Gradient Descent!

@ iteration t

- choose i_t uniformly at random
- set $\underline{w}^{(t+1)} = \underline{w}^{(t)} - \gamma \nabla L_{i_t}(\underline{w}^{(t)})$

What is $\nabla L_{i_t}(\underline{w}^{(t)})$?

$$\frac{dL_i}{dw_j} \Big|_{\underline{w}^{(t)}} = \frac{dL_i}{d\hat{y}_i} \cdot \frac{d\hat{y}_i}{dz_i} \cdot \frac{dz_i}{dw_j} \Big|_{\underline{w}^{(t)}}$$

$$= 2(\hat{y}_i - y_i) \cdot \sigma'(z_i) \cdot x_{ij} \Big|_{\underline{w}^{(t)}}$$

$$= 2(\hat{y}_i - y_i) \cdot \sigma(z_i)(1 - \sigma(z_i)) x_{ij} \Big|_{\underline{w}^{(t)}}$$

$$= 2(\hat{y}_i - y_i) \underbrace{\hat{y}_i(1 - \hat{y}_i)}_{\text{scalar, independent of } j} x_{ij} \Rightarrow \nabla L_{i_t}(\underline{w}^{(t)}) = \mathcal{S}_i \underline{x}_i$$

scalar, independent of j

$$\text{call } \mathcal{S}_i = \mathcal{S}_i(\underline{w}^{(t)})$$

$$b^{(t+1)} = b^{(t)} - \gamma \nabla L_{i_t}(b^{(t)})$$

$$\text{What is } \nabla L_{i_t}(b^{(t)}) = \frac{dL_i}{db} \Big|_{b^{(t)}} ?$$

$$\frac{dL_i}{db} \Big|_{b^{(t)}} = \frac{dL_i}{d\hat{y}_i} \cdot \frac{d\hat{y}_i}{dz_i} \cdot \frac{dz_i}{db} \Big|_{b^{(t)}}$$

$$= \mathcal{S}_i \cdot 1$$

$$\text{SGD : } \underline{w}^{(t+1)} = \underline{w}^{(t)} - \gamma \mathcal{S}_{i_t} \underline{x}_{i_t}, b^{(t+1)} = b^{(t)} - \gamma \mathcal{S}_{i_t}$$

Backpropagation

To train deeper networks (i.e. networks with one or more hidden layers), we use

BACKPROPAGATION: an efficient way of computing gradients needed for SGD in a multilayer network

Basic backpropagation structure:

Loop until convergence
for each sample $i \in \{1, \dots, n\}$

1. given feature x_i , "propagate activity forward"
(i.e. calculate network output with input x_i and current estimate of network weights) \Rightarrow "FORWARD PASS"
2. Propagate gradients backward \Rightarrow "BACKWARD PASS"
3. Update each weight (via stochastic gradient descent)

Hidden activities

We don't have targets for a hidden unit, but we can compute how fast the error changes as we change its activity

- ▶ Instead of using desired activities to train the hidden units, use **error derivatives w.r.t. hidden activities**
- ▶ Each hidden activity can affect many output units and can therefore have many separate effects on the error. These effects must be combined
- ▶ We can compute error derivatives for all the hidden units efficiently
- ▶ Once we have the error derivatives for the hidden activities, its easy to get the error derivatives for the weights going into a hidden unit

This is just the chain rule!

General weight updates

Weight updates for more complex network

z^l is input to nodes in layer l , a^l is output of those nodes

$w_{j,k}^l$ = weight in layer l , applied to layer $l-1$ output (a_k^{l-1}), feeding into next layer (z_j^l)

To update W^l , need to compute $\nabla L_i(W^l)$

$$\frac{dL_i}{dW_{j,k}^l} = \frac{dL_i}{dz_j^l} \cdot \frac{dz_j^l}{dW_{j,k}^l} = S_j^l \cdot a_k^{l-1} \Rightarrow \nabla L_i(W^l) = S^l (a^{l-1})^\top$$

$$S_j^l := \frac{dL_i}{dz_j^l} = \begin{cases} \frac{dL_i}{da_j^l} \cdot \frac{da_j^l}{dz_j^l} = [(W^{l+1})^\top S^{l+1}]_j \cdot \sigma'(z_j^l) & \text{for } l < L \\ \frac{dL_i}{d\hat{y}_{ij}} \frac{d\hat{y}_{ij}}{dz_j^l} & \text{for } l = L \end{cases}$$

$$S^l = \left\{ \begin{array}{l} ((W^{l+1})^\top S^{l+1}) \odot \sigma'(z^l) \\ \frac{dL_i}{d\hat{y}_{ij}} \end{array} \right. \quad \begin{array}{l} \text{for } l < L \\ \text{for } l = L \end{array}$$

element-wise multiplication

$$\Rightarrow \nabla L_i(W^l) = S^l (a^{l-1})^\top$$

$$\nabla L_i(b^l) = S^l$$

$$\frac{dL_i}{da_j^l} = \sum_{k=1}^{M^{l+1}} \frac{dL_i}{dz_k^{l+1}} \frac{dz_k^{l+1}}{da_j^l} = \sum_{k=1}^{M^{l+1}} S_k^{l+1} W_{kj}^{l+1} = [(W^{l+1})^\top S^{l+1}]_j$$

$$\frac{da_j^l}{dz_j^l} = \sigma'(z_j^l)$$

$$\frac{dz_j^l}{dW_{j,k}^l} = a_k^{l-1} \quad \text{because } z_j^l = (W^l a^{l-1} + b^l)_j = \sum_k w_{jk}^l a_k^{l-1} + b_j^l$$

$$\frac{dL_i}{db_j^l} = \frac{dL_i}{dz_j^l} \cdot \frac{dz_j^l}{db_j^l} = S_j^l \cdot 1 \Rightarrow \nabla L_i(b^l) = S^l$$

Backpropagation Algorithm

Backpropagation Algorithm

for $t=1, 2, 3, \dots$

Select $i_t \sim \text{unif}(1, 2, \dots, n)$

forward pass:

$$\underline{a}^0 = \underline{x}_{i_t}$$

for $\ell = 1, 2, \dots, L$

$$\underline{z}^\ell = W^\ell \underline{a}^{\ell-1} + \underline{b}^\ell$$

$$\underline{a}^\ell = \sigma(\underline{z}^\ell)$$

end

backprop

$$\underline{s}^L = \nabla L_{i_t}(\hat{y}_{i_t})$$

for $\ell = L-1, L-2, \dots, 1$

$$\underline{s}^\ell = [(W^{\ell+1})^\top \underline{s}^{\ell+1}] \odot \sigma'(\underline{z}^\ell)$$

$$\nabla L_{i_t}(W^\ell) = \underline{s}^\ell (\underline{a}^{\ell-1})^\top$$

$$\nabla L_{i_t}(b^\ell) = \underline{s}^\ell$$

$$W^{\ell+1} = W^\ell - \tau \nabla L_{i_t}(W^\ell)$$

$$b^{\ell+1} = b^\ell - \tau \nabla L_{i_t}(b^\ell)$$

end

Initialization and Standardization

Initializing weights:

Usually initial weights are chosen to be random values near zero.

In this regime, if σ is a sigmoid (e.g. $\sigma(z) = \frac{1}{1+e^{-z}}$), then for each node, the corresponding z is approximately 0, which in turn suggests $\sigma(z) \approx \frac{1}{2} + \frac{z}{4}$

In other words, when all weights ≈ 0 , each neural network node behaves approximately LINEARLY.

Thus the model starts out nearly linear, and as training moves some weights away from 0, the model becomes more nonlinear where needed.

(Starting with weights exactly equal to zero leads to "vanishing gradients" (i.e. derivatives = 0), so training cannot move the weights and we never learn a good model.)

Scaling inputs (training data)

Rule of thumb: standardize all inputs to have zero mean and standard deviation = 1

⇒ makes it easier to choose a good range of initial weights and easier to regularize (discussed later)

With standard inputs, can take initial weights by drawing from $\text{Unif}([-0.7, +0.7])$

Minibatch and Step Size

How often to update weights:

- after computing gradients for single training sample (stochastic gradient descent)

$$\theta^{(t+1)} = \theta^{(t)} - \tau \nabla L_{i_t}(\theta^{(t)})$$

- after a full sweep through the training data:

$$\theta^{(t+1)} = \theta^{(t)} - \tau \nabla \hat{R}(\theta^{(t)}) = \theta^{(t)} - \frac{\tau}{n} \sum_{i=1}^n \nabla L_i(\theta^{(t)})$$

- after a mini-batch of training cases (e.g. a random sample of $m < n$ training samples)

How much to update

- τ is a fixed constant ("fixed learning rate")

- τ is changing over time, e.g. $\tau_t \propto \gamma_t$

- add "momentum": $\theta^{(t+1)} = \theta^{(t)} - \nu^{(t+1)}$

$$\nu^{(t+1)} = \gamma \nu^{(t)} + \tau \nabla \hat{R}(\theta^{(t)})$$

Learning Rate Effects

- Check how your loss behaves during training, to spot wrong hyperparameters, bugs, etc

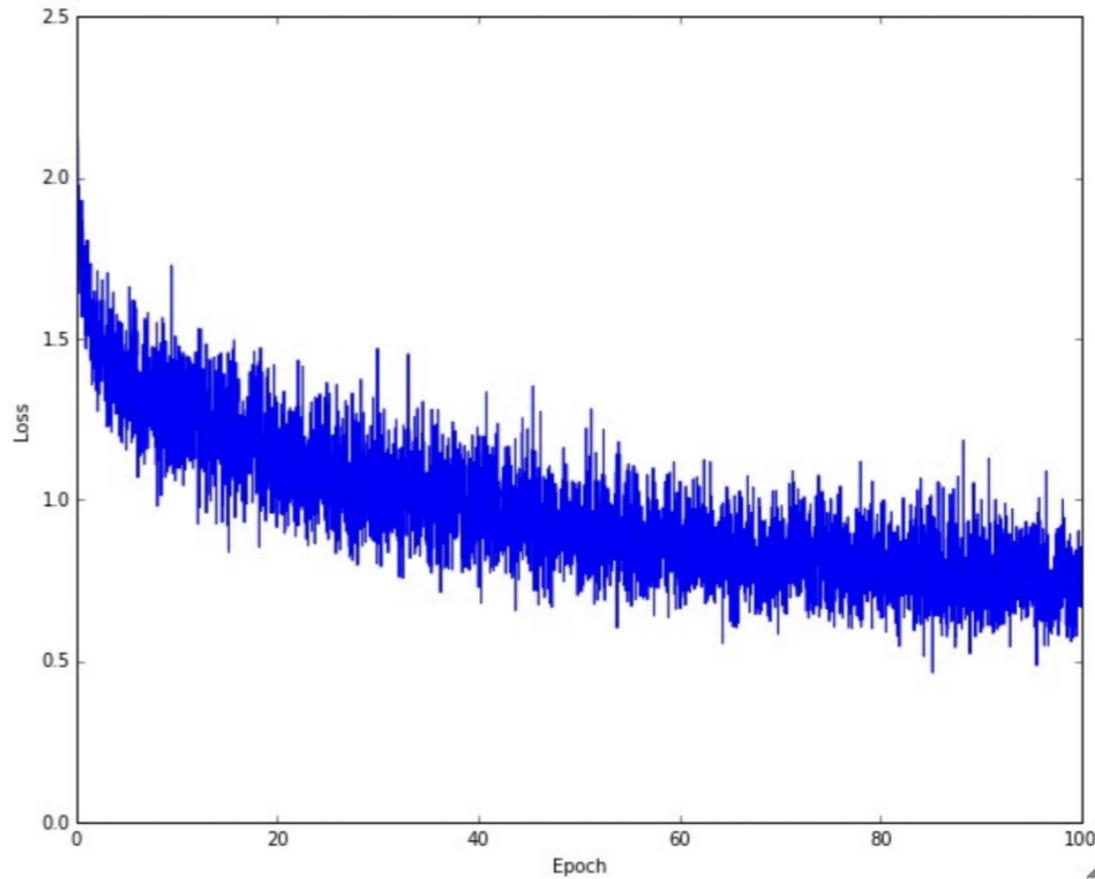
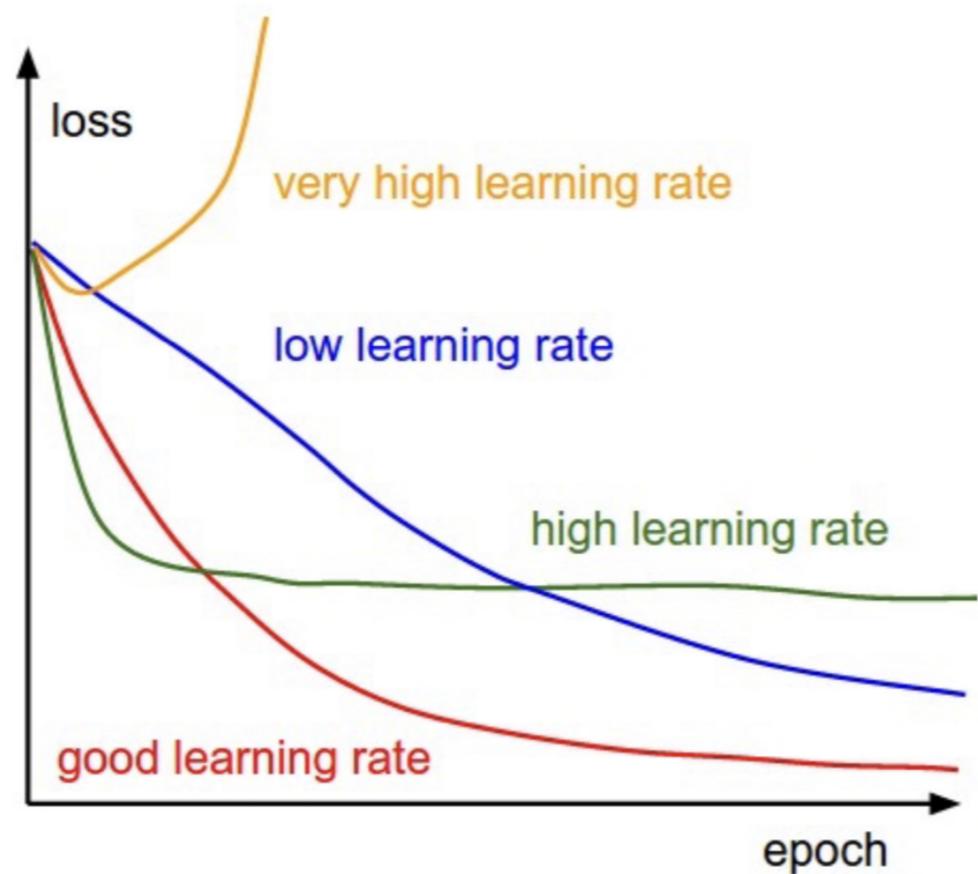


Figure : **Left:** Good vs bad parameter choices, **Right:** How a real loss might look like during training. What are the bumps caused by? How could we get a more smooth loss?