

#### Announcements:

- HW #2 is due tonight, uploaded to Gradescope. Please budget time for submission. Please be sure to print out Jupyter Notebooks and .py files for your submission.
- We will upload HW #3 tonight. Since we're going a bit slower than prior iterations of the course, we decided to make it due on Friday, Feb 9, 2024 (instead of Monday, Feb 5, 2024).

• \* \* \* UPDATE \* \* \* Midtern will be held in class on Feb. 21, 2024. (senting will be tight.)



#### **Neural networks**



First layer: $\mathbf{h}_1 = f(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$ Fully Connected (FC)First layer: $\mathbf{h}_1 = f(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$ Neural NetworksSecond layer: $\mathbf{h}_2 = f(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2)$ Multilayer Perceptron (MLP)Third (output layer): $\mathbf{z} = \mathbf{W}_3 \mathbf{h}_2 + \mathbf{b}_3$ Multilayer Perceptron (MLP)



"One recurring theme throughout neural network design is that the gradient of the cost function must be large and predictable enough to serve as a good guide for the learning algorithm." (Goodfellow et al., p. 173)



## Sigmoid unit



Its derivative is:

$$\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$$

"One recurring theme throughout neural network design is that the gradient of the cost function must be large and predictable enough to serve as a good guide for the learning algorithm." (Goodfellow et al., p. 173)

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \sigma(w)}{\partial w} \frac{\partial \mathcal{L}}{\partial \sigma(w)}$$
  
a very small #







Hyperbolic tangent,  $tanh(x) = 2\sigma(x) - 1$ 

The hyperbolic tangent is a zero-centered sigmoid-looking activation.





Its derivative is:

$$\frac{d\tanh(x)}{dx} = 1 - \tanh^2(x)$$



# Hyperbolic tangent, $tanh(x) = 2\sigma(x) - 1$

Pros:

- Around x = 0, the unit behaves linearly.
- It is differentiable everywhere.
- It is zero-centered. It's both pos & neg. So it isn't constrained to Cons: Zig-Zag.
  - Like the sigmoid unit, when a unit saturates, i.e., its values grow larger or smaller, the unit saturates and no additional learning occurs.





# **Rectified linear unit,** $\mathbf{ReLU}(x) = \max(0, x)$





Its derivative is:

$$\frac{d\text{ReLU}(x)}{dx} = \left\{ \begin{array}{cc} 1 & x > 0 \\ 0 & x < 0 \end{array} \right.$$

This function is not differentiable at x = 0. However, we can define its subgradient by setting the derivative to be between [0, 1] at x = 0.



# **ReLU unit** z = Wx + bRectified linear unit, $\operatorname{ReLU}(x) = \max(0, x)$ relu(2) = $\begin{pmatrix} 0 \\ 4 \\ 0 \end{pmatrix}$ Pros: -9.5

- In practice, learning with the ReLU unit converges faster than sigmoid and AlexNet, relu was 6x faster than tanh. tanh.
- When a unit is active, it behaves as a linear unit.
- The derivative at all points, except x = 0, is 0 or 1. When x > 0, the gradients are large, and not scaled by second order effects.
- There is no saturation if x > 0.

Cons:

- $\operatorname{ReLU}(x)$ , like sigmoid, is non-negative. SGD with ReLU() therefore zig-zags.
- $\operatorname{ReLU}(x)$  is not differentiable at x = 0. However, in practice, this is not a large issue. A heuristic when evaluating  $\frac{d\text{ReLU}(x)}{dx}\Big|_{x=0}$  is to return the left derivative (0) or the right derivative (1); this is reasonable given digital computation is subject to numerical error.
- Learning does not happen for examples that have zero activation. This can be fixed by e.g., using a leaky ReLU or maxout unit.



Softplus unit,  $\zeta(x) = \log(1 + \exp(x))$ 

One may consider using the softplus function,  $\zeta(x) = \log(1 + e^x)$ , in place of  $\operatorname{ReLU}(x)$ . Intuitively, this ought to work well as it resembles  $\operatorname{ReLU}(x)$  and is differentiable everywhere. However, empirically, it performs worse than  $\operatorname{ReLU}(x)$ .



Its derivative is:

$$\frac{d\zeta(x)}{dx} = \sigma(x)$$

Prof J.C. Kao, UCLA ECE



"One might expect it to have an advantage over the [ReLU] due to being differentiable everywhere or due to saturating less completely, but empirically it does not." (Goodfellow et al., p. 191)

	Neuron	MNIST	CIFAR10	NISTP	NORB
	With unsupervised pre-training				
	Rectifier	1.20%	49.96%	32.86%	16.46%
	Tanh	1.16%	50.79%	35.89%	17.66%
	Softplus	1.17%	49.52%	33.27%	19.19%
	Without unsupervised pre-training				
Л	Rectifier	1.43%	50.86%	32.64%	16.40%
	Tanh	1.57%	52.62%	36.46%	19.29%
/ >	Softplus	1.77%	53.20%	35.48%	17.68%
relu				Glore	ot et al., 2011a
softplus					



# Leaky rectified linear unit, $f(x) = \max(\alpha x, x)$

f = lambda x: x \* (x>0) + 0.1\*x \* (x<0)



The leaky ReLU avoids the stopping of learning when x < 0.  $\alpha$  may be treated as a selected hyperparameter, or it may be a parameter to be optimized in learning, in which case it is typically called the "PReLU" for parametrized rectified linear unit.





f = lambda x: x \* (x>0) + 0.2\*(np.exp(x) - 1) \* (x<0)



The exponential linear unit avoids the stopping of learning when x < 0. A con of this activation function is that it requires computation of the exponential, which is more expensive.





### Which LU do I use?



Figure 4: Comparison of ReLUs, LReLUs, and SReLUs on CIFAR-100. Panels (a-c) show the Clevert et al., 2015



## Maxout unit

A generalization of the ReLU and PReLU units is the maxout unit, where:

$$maxout(\mathbf{x}) = max(\mathbf{w}_1^T\mathbf{x} + b_1, \mathbf{w}_2^T\mathbf{x} + b_2)$$

This can be generalized to more than two components. If  $\mathbf{w}_1 = \mathbf{0}$  and  $b_1 = 0$ , this is the rectified linear unit.





$$ext{GELU}(x) = xP(X \leq x) = x\Phi(x) = x \cdot rac{1}{2} \Big[ 1 + ext{erf}(x/\sqrt{2}) \Big]$$



Figure 1: The GELU ( $\mu = 0, \sigma = 1$ ), ReLU, and ELU ( $\alpha = 1$ ).



#### In practice...

In practice...

- The ReLU unit is very popular.
- The sigmoid unit is almost never used; tanh is preferred.
- It may be worth trying out leaky ReLU / PReLU / ELU / maxout for your application.

• This is an active area of research.

swish, GELU





- Layer 1:  $\mathbf{h}_1 = f(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)$
- Layer 2:  $h_2 = f(W_2h_1 + b_2)$
- Layer N:  $\mathbf{z} = \mathbf{W}_N \mathbf{h}_{N-1} + \mathbf{b}_N$

What output activation do we use?

lurge negative 2 -> class 0 ; y(i) = 0

#### What outputs and cost functions?

There are several options to process the output scores, z, to ultimately arrive at a cost function. The choice of output units interacts with what cost function to use.  $\sigma(z) = \Pr\{\chi^{(i)} \text{ be longe to class 1}\}$ 

*Example:* Consider a neural network that produces a single score, z, for binary classification. As the output unit, we choose the sigmoid nonlinearity, so that  $\hat{y} = \sigma(z)$ . On a given example,  $y^{(i)}$  is either 0 or 1, and  $\hat{y}^{(i)} = \sigma(z^{(i)})$  can be interpreted as the algorithm's probability the output is in class 1. Is it better to use mean-square error or cross-entropy (i.e., corresponding to maximum-likelihood estimation) as the cost function? For n examples:

MSE = 
$$\frac{1}{2} \sum_{i=1}^{n} \left( y^{(i)} - \sigma(z^{(i)}) \right)^2$$
  $y^{(i')} = 1$   
CE =  $-\sum_{i=1}^{n} \left[ y^{(i)} \log \sigma(z^{(i)}) + (1 - y^{(i)}) \log(1 - \sigma(z^{(i)})) \right]$ 



Prof J.C. Kao, UCLA ECE



# What outputs and cost functions? (cont.)

*Example (cont):* Consider just one example, where  $y^{(i)} = 1$ . For this example,

$$\frac{\partial \text{MSE}}{\partial z^{(i)}} = -2(y^{(i)} - \sigma(z^{(i)}))(\sigma(z^{(i)})(1 - \sigma(z^{(i)})))$$

This derivative looks like the following:



When z is very negative, indicating a large MSE, the gradient saturates to zero, and no learning occurs.



#### What outputs and cost functions? (cont.)

Example (cont): Now let's consider the cross-entropy. For one example, where  $y^{(i)} = 1$ ,

$$\frac{\partial \text{CE}}{\partial z^{(i)}} = \sigma(z^{(i)}) - 1$$

This derivative looks like the following:



Notice that when z is very negative, learning will occur, and it will only "stall" when z gets close to the right answer.



• Linear output units:  $\hat{\mathbf{y}} = \mathbf{z}$ .

These output units typically specify the conditional mean of a Gaussian distribution, i.e.,

 $p(\mathbf{y}|\mathbf{z}) = \mathcal{N}(\mathbf{z}, \mathbf{I})$ 

and in this case, MLE estimation is equivalent to minimizing squared error.



• Sigmoid outputs:  $\hat{\mathbf{y}} = \sigma(\mathbf{z})$ .

These outputs are typically used in binary classification to approximate a Bernoulli distribution.

Question: Why not use the following output?

 $\Pr(y=1|\mathbf{x}) = \max\left(0,\min\left(1,\mathbf{z}\right)\right)$ 



• Softmax output:  $\hat{\mathbf{y}}_i = \operatorname{softmax}_i(\mathbf{z})$ .

The softmax is the generalization of the sigmoid output to multiple classes.

A softmax output activation is fairly common.



Now that we've defined all the elements of the neural network, the question now becomes: how do we learn its parameters?

The short answer is that we use versions of gradient descent.

forward However, neural networks architectures have units that are several layers from the output. In these scenarios, how do we arrive at the gradient?





In this lecture, we'll introduce backpropagation as a technique to calculate the gradient of the loss function with respect to parameters in a neural network.





Reading:

Deep Learning, 6.5-6.6





Now that we've defined all the elements of the neural network, the question now becomes: how do we learn its parameters?

The short answer is that we use versions of gradient descent.

However, neural networks architectures have units that are several layers from the output. In these scenarios, how do we arrive at the gradient?





Intuitively, backpropagation is the application of the chain rule for derivatives.

## Motivation for backpropagation

To do gradient descent, we need to calculate the gradient of the objective with respect to the parameters,  $\theta$ . However, in a neural network, some parameters are not directly connected to the output. How do we calculate then the gradient of the objective with respect to these parameters? Backpropagation answers this question, and is a general application of the chain rule for derivatives.



#### Nomenclature

Forward propagation:

- Forward propagation is the act of calculating the values of the hidden and output units of the neural network given an input.
- It involves taking input x, propagating it to through each hidden unit sequentially, until you arrive at the output y. From forward propagation, we can also calculate the cost function J(θ).
- In this manner, the forward propagated signals are the activations.
- With the input as the "start" and the output as the "end," information propagates in a forward manner.

Backpropagation (colloquially called backprop):

- As its name suggests, now information is passed backward through the network, from the cost function and outputs to the inputs.
- The signal that is backpropagated are the gradients.
- It enables the calculation of gradients at every stage going back to the input layer.



## Why do we need backpropagation?

- Backpropagation is computationally efficient because we will find that most of the terms used in backpropagation can be cached from the forward pass.
- Informally (at least for me) sometimes taking analytical multivariate gradients can be challenging. Backpropagation breaks down these multivariate gradients into easier steps.
- A few further notes on backpropagation.
  - Backpropagation is not the learning algorithm. It's the method of computing gradients.
  - Backpropagation is not specific to multilayer neural networks, but a general way to compute derivatives of functions.



$$f(x, y, z) = (x + y)z$$

$$\frac{\partial f}{\partial z} = x + y \qquad \frac{\partial f}{\partial x} = \frac{\partial f}{\partial y} = z$$





$$f(x, y, z) = (x + y)z$$

#### Forward propagation:

























**Forward pass:** 









**Forward pass:** 



#### **Backward pass:**

$$\frac{\partial f}{\partial x} \frac{\partial \mathcal{L}}{\partial f} \xrightarrow{f} \frac{\partial \mathcal{L}}{\partial f}$$



**Forward pass:** 







## Idea: computational graphs apply to gradients



The basic intuition of backpropagation is that we break up the calculation of the gradient into **small and simple steps**. Each of these nodes in the graph is a straightforward gradient calculation, where we multiply an **input** (the "upstream derivative") with a **local gradient** (an application of the chain rule).

Composing all of these gradients together returns the overall gradient.







Add gate: distributes the gradient





Add gate: distributes the gradient





Add gate: distributes the gradient Mult gate: switches the gradient





Add gate: distributes the gradient Mult gate: switches the gradient

relu (x) = max(x, o)

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

$$\frac{\partial f}{\partial x} = \begin{cases} 1, & \text{if } x > y \\ 0, & \text{else} \end{cases}$$

$$= \mathbf{I} \{x > y\}$$





Add gate: distributes the gradient Mult gate: switches the gradient Max gate: routes the gradient





$$f(\mathbf{w}, \mathbf{x}) = \frac{1}{1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2))}$$



$$f(\mathbf{w}, \mathbf{x}) = \frac{1}{1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2))}$$





$$f(\mathbf{w}, \mathbf{x}) = \frac{1}{1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2))}$$









$$f(\mathbf{w}, \mathbf{x}) = \frac{1}{1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2))}$$



$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial f}{\partial z} \frac{\partial \mathcal{L}}{\partial f} \qquad \qquad \frac{\partial f}{\partial z} = -\frac{1}{\frac{z^2}{z^2}}$$



## A more involved scalar example



 $f(x) = e^{x}$   $\frac{\partial f}{\partial f} = e^{x}$ 



#### A more involved scalar example







$$f(\mathbf{w}, \mathbf{x}) = \frac{1}{1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2))}$$





$$f(\mathbf{w}, \mathbf{x}) = \frac{1}{1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2))}$$





$$f(\mathbf{w}, \mathbf{x}) = \frac{1}{1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2))}$$





With backpropagation, as long as you can break the computation into components where you know the local gradients, you can take the gradient of anything.



### What happens when two gradient paths converge?

