

#### Announcements:

• HW #2 is due Monday Jan 29, uploaded to Gradescope. To submit your Jupyter Notebook, print the notebook to a pdf with your solutions and plots filled in. You must also submit your .py files as pdfs.



## **Tonmoy Monsoor**



- Graduate student researcher at Big data and Complex Networks Group
- Fourth year Ph.D. student advised by Professor Vwani Roychowdhury
- Research interests
  - Reinforcement learning
  - Distributed optimization
  - Pattern extraction in dynamic networks
- Favorite classes at UCLA
  - Convex optimization, ECE 236B
  - Neural signal processing, ECE 243A



Softmax:

$$\arg\min_{\theta} \frac{1}{m} \sum_{i=1}^{m} \left( \log \sum_{j=1}^{c} e^{a_j(\mathbf{x})} - a_{y^{(i)}}(\mathbf{x}^{(i)}) \right)$$

Parameters?

$$\theta = \{ w_{1} b \}$$



110)

21

20

θ

 $\Box$ 

Prof J.C. Kao, UCLA ECE

- Our goal in machine learning is to optimize an objective function, f(x). (Without loss of generality, we'll consider minimizing f(x). This is equivalent to maximizing -f(x).)
- From basic calculus, we recall that the derivative of a function,  $\frac{df(x)}{dx}$  tells • the slope of f(x) at point x. • For small enough  $\epsilon$ ,  $f(x + \epsilon) \approx f(x) + \epsilon f'(x)$ . • This tells us how to reduce (or increase)  $f(\cdot)$  for small enough steps. • Percell that  $f(x) = \hat{f}(x) + \hat{f}(x)$ . us the slope of f(x) at point x.

  - Recall that when f'(x) = 0, we are at a stationary point or critical point. This may be a local or global minimum, a local or global maximum, or a saddle point of the function.
- In this class we will consider cases where we would like to maximize fw.r.t. vectors and matrices, e.g.,  $f(\mathbf{x})$  and  $f(\mathbf{X})$ .
- Further, often  $f(\cdot)$  contains a nonlinearity or non-differentiable function. In these cases, we can't simply set  $f'(\cdot) = 0$ , because this does not admit a closed-form solution.
- However, we can iteratively approach an critical point via gradient descent.

1(0)



To do so, we use the technique of gradient descent.





# Terminology

- A global minimum is the point,  $\mathbf{x}_g$ , that achieves the absolute lowest value of  $f(\mathbf{x})$ . i.e.,  $f(\mathbf{x}) \ge f(\mathbf{x}_g)$  for all  $\mathbf{x}$ .
- A local minimum is a point, x<sub>ℓ</sub>, that is a critical point of f(x) and is lower than its neighboring points. However, f(x<sub>ℓ</sub>) > f(x<sub>g</sub>).
- Analogous definitions hold for the **global maximum** and **local maximum**.
- A saddle point are critical point of f(x) that are not local maxima or minima. Concretely, neighboring points are both greater than and less than f(x).

saddle point



## Gradient

Recall the gradient,  $\nabla_{\mathbf{x}} f(\mathbf{x})$ , is a vector whose *i*th element is the partial derivative of  $f(\mathbf{x})$  w.r.t.  $x_i$ , the *i*th element of  $\mathbf{x}$ . Concretely, for  $\mathbf{x} \in \mathbb{R}^n$ ,

$$\Delta \mathbf{X} = \begin{pmatrix} \delta \mathbf{x}_{1} \\ \delta \mathbf{x}_{2} \\ \vdots \\ \delta \mathbf{x}_{N} \end{pmatrix} \qquad \nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_{1}} \\ \frac{\partial f(\mathbf{x})}{\partial x_{2}} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial (x_{n})} \end{bmatrix} \qquad \mathbf{X} + \Delta \mathbf{X}^{\mathsf{T}} \nabla_{\mathbf{x}} f(\mathbf{x})$$

• The gradient tells us how a small change in  $\Delta \mathbf{x}$  affects  $f(\mathbf{x})$  through  $f \rightarrow \mathcal{L}$   $\mathbf{x} \rightarrow \boldsymbol{\theta}$ • The gradient tells us how a small change in  $\Delta \mathbf{x}$  affects  $f(\mathbf{x})$  through  $f(\mathbf{x} + \mathbf{u}) \approx f(\mathbf{x}) + \mathbf{u}^{\mathsf{T}} \nabla f(\mathbf{x})$  $f(\mathbf{x} + \Delta \mathbf{x}) \approx f(\mathbf{x}) + \Delta \mathbf{x}^{\mathsf{T}} \nabla_{\mathbf{x}} f(\mathbf{x})$ 

• The directional derivative of  $f(\mathbf{x})$  in the direction of the unit vector  $\mathbf{u}$  is given by:

$$\mathbf{u}^T \nabla_{\mathbf{x}} f(\mathbf{x}) \qquad \qquad \| \mathbf{u} \| \in \mathbf{I}$$

• The directional derivative tells us the slope of f in the direction  $\mathbf{u}$ .



## Arriving at gradient descent

 To minimize f(x), we want to find the direction in which f(x) decreases the fastest. To do so, we find the direction u which minimizes the directional derivative.

where  $\theta$  is the angle between the vectors  $\mathbf{u}$  and  $\nabla_{\mathbf{x}} f(\mathbf{x})$ .

- This quantity is minimized for u pointing in the opposite direction of the gradient, so that cos(θ) = −1.
   u = -∇<sub>x</sub> (-lx)
- Hence, we arrive at gradient descent. To update x so as to minimize f(x), we repeatedly calculate:

$$\mathbf{x} := \mathbf{x} - \epsilon \nabla_x f(\mathbf{x})$$

•  $\epsilon$  is typically called the *learning rate*. It can change over iterations. Setting the value of  $\epsilon$  appropriately is an important part of deep learning.



### Example:

Animations thanks to: http://louistiao.me/notes/visualizing-and-animating-optimization-algorithms-with-matplotlib/





### Example:

Animations thanks to: http://louistiao.me/notes/visualizing-and-animating-optimization-algorithms-with-matplotlib/





```
def gd(func, x0, eps=1e-4, tol=1e-3):
   last diff = np.Inf
   x = x0
   path = [np.copy(x0)]
   costs = [func(x0)[0]]
   grads = []
   i = 1
   hit_max = False Softmax. loss_and-grad
   while last diff > tol:
       cost, q = func(x) # returns the cost and the gradient
       x -= eps*g # gradient step
       last diff = np.linalg.norm(x - path[-1]) # stopping criterion
       i += 1
       if i > max_iters:
           hit max = True
           break
       path.append(np.copy(x))
       costs.append(cost)
       grads.append(g)
   return path, costs, grads, hit max
```



http://seas.ucla.edu/~kao/opt\_anim/1gd.mp4 http://seas.ucla.edu/~kao/opt\_anim/2gd.mp4









Why not always use smaller learning rates?





### Interpreting the cost

The cost function can be very informative as to how to adjust your step sizes for gradient descent.





# Why not use a numerical gradient?

f is like & x is like O

$$\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$



How does this example differ from what we will really encounter?

- In this example, we know the function f() exactly, and thus at every point in space, we can calculate the gradient at that point exactly.
- In optimization, we differentiate the cost function f() with respect to the parameters.
  - The gradient of f() w.r.t. parameters is a function of the training data!  $\chi^{(c)}, \chi^{(c)}$
  - Hence, we can think of each data point as providing a noisy estimate of the gradient at that point.

$$\nabla_{\theta} \mathbb{E}_{(x^{(i)},y^{(i)},D)} \left[ \log P(\theta) \right] = \mathbb{E} \left[ \nabla_{\theta} \log P(\theta, x^{(i)}, y^{(i)}) \right]$$
$$= \frac{1}{m} \sum_{i=1}^{m} \nabla_{\theta} \log P(\theta, x^{(i)}, y^{(i)})$$



k = 64

However, it's expensive to have to calculate the gradient by using *every example* in the training set.

To this end, we may want to get a noisier estimate of the gradient with fewer examples.

# Batch vs minibatch (cont)

Calculating the gradient exactly is expensive, because it requires evaluating the model on all m examples in the dataset. This leads to an important distinction.

#### 50,000

- Batch algorithm: uses all *m* examples in the training set to calculate the gradient.
- Minibatch algorithm: approximates the gradient by calculating it using k training examples, where m > k > 1.
- Stochastic algorithm: approximates the gradient by calculating it over one example. m = 1

It is typical in deep learning to use minibatch gradient descent. Note that some may also use minibatch and stochastic gradient descent interchangeably.



To get a more robust estimate of the gradient, we would use as many data samples as possible.

and its gradient is:

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{m} \sum_{i=1}^{m} \log p_{\text{model}}(\mathbf{x}^{(i)}, y^{(i)}; \theta) \\ &= \frac{1}{m} \sum_{i=1}^{m} \nabla_{\theta} \log p_{\text{model}}(\mathbf{x}^{(i)}, y^{(i)}; \theta) \\ &\approx \mathbb{E} \left[ \nabla_{\theta} \log p_{\text{model}}(\mathbf{x}^{(i)}, y^{(i)}; \theta) \right] \end{aligned}$$



You'll do this in the HW. More on this later in the optimization lecture...

- And a lot more to be said about optimization.
- First order vs second order methods
- Momentum
- Adaptive gradients.
- ... all of these will become quite important when we get to neural networks. We'll cover these in an optimization lecture.



In this lecture, we'll introduce the neural network architecture, parameters, and its inspiration from biological neurons.



Reading:

Deep Learning, 6 (intro), 6.1, 6.2, 6.3, 6.4











Neurons are the main signaling units of the nervous system.

Neurons have four regions:

1) Cell body (soma) – metabolic center, with nucleus, etc. 2) Dendrites – tree like structure for receiving **input** signals. 2.5. Axon hillock: integrative part 3) Axon - single, long, tubular structure for sending output signals.

4) Presynaptic terminals – sites of communication to next  $f(z) = \begin{cases} z > threshold, 1 \\ z \leq threshold, 0 \end{cases}$ neurons.

Axons (the output) convey signals to other neurons:

- Conveys electrical signals long distances (0.1mm
- Conveys action potentials (~100 mV, ~1 ms pulses).
- O Action potentials initiate at the axon hillock.
- Propagate w/o distortion or failure at 1-100 m/s.



#### Neurons are diverse (unlike in neural networks)

Figure 2-4 Neurons can be classified as unipolar, bipolar, or multipolar according to the number of processes that originate from the cell body.

A. Unipolar cells have a single process, with different segments serving as receptive surfaces or releasing terminals. Unipolar cells are characteristic of the invertebrate nervous system.

**B.** Bipolar cells have two processes that are functionally specialized: the dendrite carries information to the cell, and the axon transmits information to other cells.

**C.** Certain neurons that carry sensory information, such as information about touch or stretch, to the spinal cord belong to a subclass of bipolar cells designated as pseudo-unipolar. As such cells develop, the two processes of the embryonic bipolar cell become fused and emerge from the cell body as a single process. This outgrowth then splits into two processes, *both* of which function as axons, one going to peripheral skin or muscle, the other going to the central spinal cord.

D. Multipolar cells have an axon and many dendrites. They are the most common type of neuron in the mammalian nervous system. Three examples illustrate the large diversity of these cells. Spinal motor neurons (left) innervate skeletal muscle fibers. Pyramidal cells (middle) have a roughly triangular cell body; dendrites emerge from both the apex (the apical dendrite) and the base (the basal dendrites). Pyramidal cells are found in the hippocampus and throughout the cerebral cortex. Purkinje cells of the cerebellum (right) are characterized by the rich and extensive dendritic tree in one plane. Such a structure permits enormous synaptic input. (Adapted from Ramón y Cajal 1933.)



D Three types of multipolar cells



erebellum



Neurons fundamentally communicate through all-or-nothing spikes (not analog values!):





... And the spikes are probabilistic.





The spikes reflect an underlying rate.



This rate is what the neural networks are "encoding."



How does the artificial neuron compare to the real neuron?





How does the artificial neuron compare to the real neuron?

# The artificial neuron (cont.)

- The incoming signals, a vector  $\mathbf{x} \in \mathbb{R}^N$ , reflects the output of N neurons that are connected to the current artificial neuron.
- The incoming signals, x, are pointwise multiplied by a vector, w ∈ ℝ<sup>N</sup>. That is, we calculate w<sub>i</sub>x<sub>i</sub> for i = 1,..., N. This computation reflects dendritic processing.
- The "dendritic-processed" signals are then summed, i.e., we calculate 
   ∑<sub>i</sub> w<sub>i</sub>x<sub>i</sub> + b. This computation reflects integration at the axon hillock (the first "Node of Ranvier") where action potentials are generated if the integrated signal is large enough.
- The output of the artificial neuron is then a nonlinearly transformation of the integrated signal, i.e., f(∑<sub>i</sub> w<sub>i</sub>x<sub>i</sub> + b). Rather than reflecting whether an action potential was generated or not (which is a noisy process), this nonlinear output is typically treated as the *rate* of the neuron. The higher the rate, the more likely the neuron is to fire action potentials.



### How does the artificial neuron compare to the real neuron?





## Caution when comparing to biology

These computing analogies are not precise, with large approximations. Limitations in the analogy include:

- Synaptic transmission is probabilistic, nonlinear, and dynamic.
- Dendritic integration is probabilistic and may be nonlinear.
- Dendritic computation has associated spatiotemporal decay.
- Integration is subject to biological constraints; for example, ion channels (which change the voltage of the cell) undergo refractory periods when they do not open until hyperpolarization.
- Different neurons may have different action potential thresholds depending on the density of sodium-gated ion channels.
- Feedforward and convolutional neural networks have no recurrent connections.
- Many different cell types.
- Neurons have specific dynamics that can be modulated by e.g., calcium concentration.
- And so many more...



## Caution when comparing to biology

On the prior list, several of these bullet points constitute entire research areas. E.g., several labs work specifically on studying the details of synaptic transmission.

Big picture: though neural networks are inspired by biology, they approximate biological computation at a fairly crude level. These networks ought not be thought of as models of the brain, although recent work (including my research group's work) has used them as a means to propose mechanistic insight into neurophysiological computation.



## Nomenclature

Some naming conventions.

- We call the first layer of a neural network the "input layer." We typically represent this with the variable **x**.
- We call the last layer the "output layer." We typically represent this with the variable z. (Note: why not y to match our prior nomenclature for the supervised outputs? Because the output of the network may be a processed version of z, e.g., softmax(z).)
- We call the intermediate layers the "hidden layers." We typically represent this with the variable  ${f h}$ .
- When we specify that a network has N layers, this does not include the input layer.







## **Neural networks**



This network has 6 neurons (not counting the input). It has (3x4) + (4x2) = 20 weights, and 4+2 = 6 biases for a total of 26 learnable parameters.



## Neural network architecture 2

An example 3-layer network is shown below.



Here,  $h_{ij}$  denotes the *j*th element of  $h_i$ . There are many considerations in architecture design, which we will later discuss.



## **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)



## **Neural networks**



This network has 10 neurons (not counting the input). It has (3x4) + (4x4) + (4x2) = 36 weights, and 4+4+2=10 biases for a total of 46 learnable parameters.



Prof J.C. Kao, UCLA ECE



f(x) = x

## Neural networks

The above figure suggests the following equation for a neural network.

• Layer 1:  $\mathbf{h}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$ • Layer 2:  $\mathbf{h}_2 = \mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2$ •  $\vdots$ • Layer N:  $\mathbf{z} = \mathbf{W}_N \mathbf{h}_{N-1} + \mathbf{b}_N$ • Layer N:  $\mathbf{z} = \mathbf{W}_N \mathbf{h}_{N-1} + \mathbf{b}_N$ 

Any composition of linear functions can be reduced to a single linear function. Here, z = Wx + b, where

$$\mathbf{W} = \mathbf{W}_N \cdots \mathbf{W}_2 \mathbf{W}_1$$

and

$$\mathbf{b} = \mathbf{b}_N + \mathbf{W}_N \mathbf{b}_{N-1} + \dots + \mathbf{W}_N \cdots \mathbf{W}_3 \mathbf{b}_2 + \mathbf{W}_N \cdots \mathbf{W}_2 \mathbf{b}_1$$



## **Neural networks**

The above figure suggests the following equation for a neural network.

- Layer 1:  $\mathbf{h}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$
- Layer 2:  $h_2 = W_2 h_1 + b_2$
- •
- Layer N:  $\mathbf{z} = \mathbf{W}_N \mathbf{h}_{N-1} + \mathbf{b}_N$

Any composition of linear functions can be reduced to a single linear function. Here, z = Wx + b, where

$$\mathbf{W} = \mathbf{W}_N \cdots \mathbf{W}_2 \mathbf{W}_1$$

and

$$\mathbf{b} = \mathbf{b}_N + \mathbf{W}_N \mathbf{b}_{N-1} + \dots + \mathbf{W}_N \cdots \mathbf{W}_3 \mathbf{b}_2 + \mathbf{W}_N \cdots \mathbf{W}_2 \mathbf{b}_1$$



$$f(x) = ax + b$$

- This may be useful in some contexts. For example, when  $\dim(\mathbf{h}) \ll \dim(\mathbf{x})$ , this corresponds to finding a low-rank representation of the inputs.
- However, a system with greater complexity may require a higher capacity model.

Autoencodor





## Introducing nonlinearity

To increase the network capacity, we can make it nonlinear. We do this by introducing a nonlinearity,  $f(\cdot)$ , at the output of each artificial neuron.

- Layer 1:  $h_1 = f(W_1x + b_1)$
- Layer 2:  $h_2 = f(W_2h_1 + b_2)$
- Laver N:  $\mathbf{z} = \mathbf{W}_N \mathbf{h}_{N-1} + \mathbf{b}_N$

A few notes:

- These equations describe a *feedforward neural network*, also called a *multilayer perceptron*.
- $f(\cdot)$  is typically called an *activation function* and is applied elementwise on its input. "unit"
- The activation function does not typically act on the output layer, z, as these are meant to be interpreted as scores. Instead, separate "output activations" are used to process z. While these output activations may be the same as the activation function, they are typically different. For example, it may comprise a softmax or SVM classifier.



# A perspective on feature learning

One area of machine learning is very interested in finding *features* of the data that are then good for use as the input data to a classifier (like a SVM). Why might this be important?



The intermediate layers of the neural network (i.e.,  $h_1$ ,  $h_2$ , etc.) are features that the later layers then use for decoding. If the performance of the neural network is well, these features are good features.

Importantly, these features don't have to be handcrafted.



Consider a system that produces training data that follows the  $xor(\cdot)$  function. The xor function accepts a 2-dimensional vector  $\mathbf{x}$  with components  $x_1$  and  $x_2$ and returns 1 if  $x_1 \neq x_2$ . Concretely,

$x_1$	$x_2$	$\operatorname{xor}(\mathbf{x})$
0	0	0 0
0	1	1 🗙
1	0	1 🗙
1	1	0 0

$$J(\theta) = \frac{1}{2} \sum_{\mathbf{x}} (g(\mathbf{x}) - y(\mathbf{x}))^2$$

(Note, we wouldn't know xor(x), but we would have samples of corresponding inputs and outputs from training data. Hence, it may be better to simply replace xor(x) with y(x) representing training examples.)



# Example: XOR

Consider first a linear approximation of xor, via  $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ . Then,

$$\frac{\partial J(\mathbf{w}, b)}{\partial \mathbf{w}} = \sum_{\mathbf{x}} (\mathbf{w}^T \mathbf{x} + b - y(\mathbf{x})) \mathbf{x}$$
$$\frac{\partial J(\mathbf{w}, b)}{\partial b} = \sum_{\mathbf{x}} (\mathbf{w}^T \mathbf{x} + b - y(\mathbf{x}))$$

Equating these to 0, we arrive at:

$$(w_1+b-1)\begin{bmatrix}1\\0\end{bmatrix}+(w_2+b-1)\begin{bmatrix}0\\1\end{bmatrix}+(w_1+w_2+b)\begin{bmatrix}1\\1\end{bmatrix} = \begin{bmatrix}0\\0\end{bmatrix}$$

These two equations can be simplified as:

$$(w_1 + b - 1) + (w_1 + w_2 + b) = 0$$
  
(w\_2 + b - 1) + (w\_1 + w\_2 + b) = 0

These equations are symmetric, implying  $w_1 = w_2 = w$ . This means:

$$3w + 2b - 1 = 0 \implies b = \frac{1 - 3w}{2}$$

Prof J.C. Kao, UCLA ECE



# Example: XOR

Now let's consider using a two-layer neural network, with the following equation:

$$g(\mathbf{x}) = \mathbf{w}^T \max(0, \mathbf{W}^T \mathbf{x} + \mathbf{c}) + b$$

We haven't yet discussed how to optimize these parameters, but the point here is to show that by introducing a simple nonlinearity like  $f(x) = \max(0, x)$ , we can now solve the  $\operatorname{xor}(\cdot)$  problem. Consider the solution:

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$
$$\mathbf{c} = \begin{bmatrix} 0, -1 \end{bmatrix}^{T}$$
$$\mathbf{w} = \begin{bmatrix} 1, -2 \end{bmatrix}^{T}$$



There are a variety of activation functions. We'll discuss some more commonly encountered ones.



"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)