

# CSE 176 Introduction to Machine Learning Lecture 10: Back Propagation

Some slides from O. Veksler, Y. Boykov, A. Ng, Y. LeCun, G. Hinton, A. Ranzato, R. Fergus

#### **Recap: Neural Unit**





#### Recap: Multi-layer perceptron



Example of Multi Layer Perceptron (MLP)



#### **Recap:** activation function





Which of the followings would you consider to be valid activation functions?

(a) 
$$f(x) = -min(2, x)$$
  
(b)  $f(x) = 0.9x + 1$   
(c)  $f(x) = \begin{cases} min(x, 0.1x) \text{ if } x \ge 0\\ min(x, 0.1x) \text{ if } x < 0 \end{cases}$   
(d)  $f(x) = \begin{cases} max(x, 0.1x) \text{ if } x \ge 0\\ min(x, 0.1x) \text{ if } x < 0 \end{cases}$ 



#### How to train neural networks?



Example of Multi Layer Perceptron (MLP)





# Optimization via Gradient Descent

Optimization of continuous differentiable functions

□ How to minimize a function of a single variable

$$f(x) = (x-5)^2$$

- Take derivative and set it to 0

$$\frac{d}{dx}f(x) = 0$$

- May find a closed form solution

$$\frac{d}{dx}f(x) = 2(x-5) = 0 \qquad \Rightarrow \qquad x=5$$





What is "slope" of  $L(x_1, x_2)$  at a given point  $\mathbf{x} = (x_1, x_2)$ ?



"heat-map" visualization of L



"heat-map" visualization of L



The most common optimization method for continuous differentiable (multi-variate) functions:

# gradient descent

take a step  $\mathbf{x}' = \mathbf{x} - \alpha \nabla L$ towards lower values of the function "heat-map" visualization of L



**direction of the steepest descent at point**  $\mathbf{x}=(x_1,x_2)$ 

#### Multi-variate functions

### **Gradient Descent**

#### Example: for a function of two variables



- direction of (negative) gradient at point  $\mathbf{x} = (x_1, x_2)$  is direction of the steepest descent towards lower values of function  $L_{5,12}$
- magnitude of gradient at  $\mathbf{x} = (x_1, x_2)$  gives the value of the slope

#### Multi-variate functions

### **Gradient Descent**

#### Example: for a function of two variables



#### Multi-variate functions

#### **Gradient Descent**

#### Example: for a function of two variables





#### How to Set Learning Rate $\alpha$ ?

$$\mathbf{x}' = \mathbf{x} - \alpha \, \nabla L$$

If α too small, too many iterations to converge



 If α too large, may overshoot the local minimum and possibly never even converge



### Variable Learning Rate

If desired, can change learning rate  $\alpha$  at each iteration

 $\begin{aligned} \mathbf{k} &= 1 \\ \mathbf{x}^{(1)} &= \text{any initial guess} \\ \text{choose } \alpha, \varepsilon \\ \text{while } \alpha ||\nabla \mathbf{L}(\mathbf{x}^{(k)})|| &> \varepsilon \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \alpha \nabla \mathbf{L}(\mathbf{x}^{(k)}) \\ \mathbf{k} &= \mathbf{k} + 1 \end{aligned}$ 

fixed α gradient descent variable  $\alpha$  gradient descent



### Learning Rate

 Monitor learning rate by looking at how fast the objective function decreases





# Derivative and Back Propagation

#### How to take derivatives w.r.t. weights?



Example of Multi Layer Perceptron (MLP)



# **Computing Derivatives: Small Example**

- Small network f(x,y,z) = (x+y)z
- Rewrite using
  - q = x + y
- f(x,y,z) = qz
- each node does one operation





# **Computing Derivatives: Small Example**

- Small network f(x,y,z) = (x+y)z
- Rewrite using
  - q = x + y
  - f(x,y,z) = qz
- Example of computing f(-2,5,-4)



# **Computing Derivatives: Small Example**

- Small network f(x,y,z) = (x+y)z
- Rewrite using  $\mathbf{q} = \mathbf{x} + \mathbf{y} \Rightarrow \mathbf{f}(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \mathbf{q}\mathbf{z}$
- Want  $\frac{\partial \mathbf{f}}{\partial \mathbf{x}^{'} \partial \mathbf{y}^{'} \partial \mathbf{z}}$

chain rule for 
$$f(y(x))$$
  
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial y} \frac{\partial y}{\partial x}$$

- Compute  $\frac{\partial \mathbf{f}}{\partial}$  from the end backwards
  - for each edge, with respect to the main variable at edge origin
  - using chain rule with respect to the variable at edge end, if needed



# Computing Derivatives: Chain of Chain Rule

• Compute  $\frac{\partial \mathbf{d}}{\partial \mathbf{d}}$  from the end backwards

direction of computation

- for each edge, with respect to the main variable at edge origin
- using chain rule with respect to the variable at edge end, if needed



# **Computing Derivatives Backwards**





- Have loss function **L**(**o**)
- Need derivatives for all

- Will compute derivatives from end to front, backwards
- On the way will also compute intermediate derivatives  $\frac{\partial \mathbf{L}}{\partial \mathbf{h}}$



- Simplified view at a network node
  - inputs **x**,**y** come in
  - node computes some function h(x,y)





- At each network node
  - inputs **x**,**y** come in
  - nodes computes activation function h(x,y)
- Have loss function L(·)



- Need  $\frac{\partial \mathbf{L}}{\partial \mathbf{x}}, \frac{\partial \mathbf{L}}{\partial \mathbf{y}}$
- Easy to compute local node derivatives  $\frac{\partial h}{\partial x}$ ,  $\frac{\partial h}{\partial y}$





- More complete view at a network node
  - inputs x,y come in, get multiplied by weight w and v
  - node computes function h(wx,vy)
  - node output **h** gets multiplied by **u**







• To be concrete, let **h**(**i**,**j**) = **i** + **j** 





- h(i,j) = i + j
- Break into more computational nodes
  - all computation happens inside nodes, not on edges





- Some of these partial derivatives are intermediate
  - their values will not be used for gradient descent



# **Computing Derivatives: Vector Notation**

• Inputs outputs are often vectors

$$\begin{array}{c|c} x & h(W^{1}x + b^{1}) \\ \end{array} & \begin{array}{c} h^{1} & h(W^{2}h^{1} + b^{2}) \\ \end{array} & \begin{array}{c} h^{2} & h(W^{3}h^{2} + b^{3}) \\ \end{array} & \begin{array}{c} o \\ \end{array} & \begin{array}{c} L(o) \\ \end{array} \end{array}$$

- h(a) is a function from  $\mathbb{R}^n$  to  $\mathbb{R}^m$
- Chain rule generalizes to vector functions



# **Computing Derivatives: Vector Notation**

- Let  $\mathbf{f}(\mathbf{x}): \mathbf{R}^n \rightarrow \mathbf{R}^m$ ,
  - $\mathbf{x}$  is  $\mathbf{n}$ -dimensional vector and output  $\mathbf{f}(\mathbf{x})$  is  $\mathbf{m}$ -dimensional vector
- Jacobian matrix
  - has **m** rows and **n** columns
  - has  $\frac{\partial \mathbf{f}_i}{\partial \mathbf{x}_j}$  in row **i**, column **j**



# **Computing Derivatives: Vector Notation**

- $f(x): \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $g(x): \mathbb{R}^k \rightarrow \mathbb{R}^n$
- $f(g(x)): \mathbb{R}^k \rightarrow \mathbb{R}^m$
- Chain rule for vector functions

 $\partial \mathbf{f} - \partial \mathbf{f} \partial \mathbf{g}$  $\partial \mathbf{x} \quad \partial \mathbf{g} \, \partial \mathbf{x}$ Jacobian matrices



## Vector Notation: Look at One Node

- h, x, y are vectors
- already computed Jacobian  $\frac{\partial \mathbf{L}}{\partial \mathbf{h}}$
- Need Jacobians  $\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}$



# Vector Notation: Look at One Node

- Can apply to matrices (and tensors) as well
- But first vectorize matrix (or tensor)
- Say W is 10 x 5, stretch into 50x1 vector
- Still denote Jacobian by  $\frac{\partial \mathbf{h}}{\partial \mathbf{W}}$



# Vector Notation: Look at One Node

- Easy to compute local node Jacobians  $\frac{\partial \mathbf{h}}{\partial \mathbf{x}}, \frac{\partial \mathbf{h}}{\partial \mathbf{W}}$
- But they can get very large (although sparse)
- Say **h** is 1000 x 1, **W** is 1000 x 500, then  $\frac{\partial \mathbf{h}}{\partial \mathbf{W}}$  is 1000 x 500,000





□Gradient Descent Optimization

- Chain rule of derivatives
- □Back propagation

