

## Recall

$$s = f(x; W) = Wx \quad \text{score function}$$

$$L_i = \sum_j y_i \max(0, s_j - s_{y_i} + 1) \quad \text{SVM loss}$$

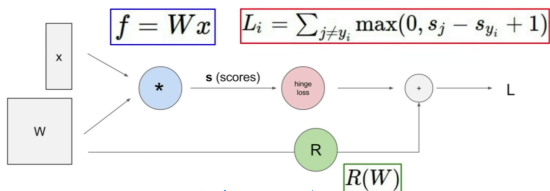
$$L = \frac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2 \quad \text{data loss + regularization}$$

want  $\nabla W L$

→ solve it with optimization  $\left( \begin{array}{l} \text{analytic gradient} \\ \text{numerical gradient} \end{array} \right)$

## analytic gradient (thru computational graphs)

### Computational graphs



- important to calculate gradient for neural network (complex function)

- e.g. find gradient

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

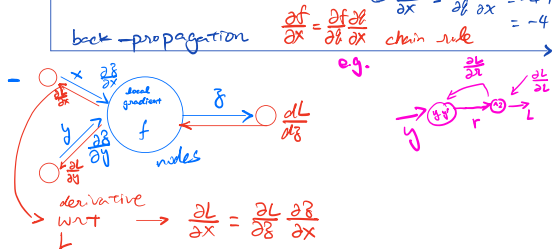
$$\textcircled{a} \quad x=-2, y=5, z=-4$$

$$f = x+y \quad \frac{\partial f}{\partial x} = 1 \quad \frac{\partial f}{\partial y} = 1$$

$$f = x \cdot y \quad \frac{\partial f}{\partial x} = y \quad \frac{\partial f}{\partial y} = x$$

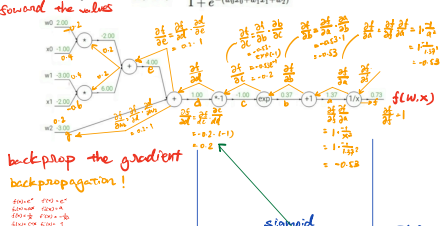
$$\text{find } \frac{df}{dx} \quad \frac{df}{dy} \quad \frac{df}{dz}$$

$$\text{back-propagation} \quad \frac{\partial f}{\partial x} = \frac{\partial f}{\partial z} \frac{\partial z}{\partial x} \quad \text{chain rule}$$



Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$



backprop the gradient backpropagation!

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

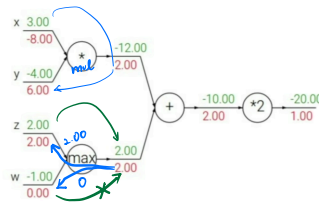
$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x))$$

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

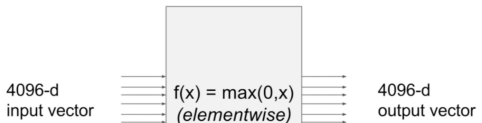
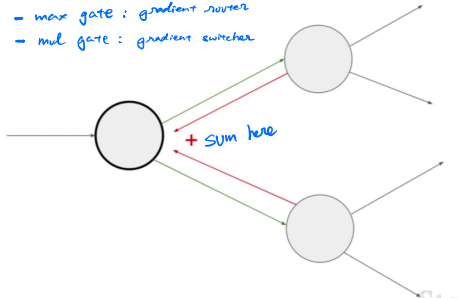
$$= (1 - 0.73) \cdot 0.73 \approx 0.2$$

→ to combine into a more complex node

## other nodes (more complex ones)



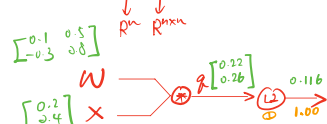
- max gate: gradient router  
- mul gate: gradient switcher



Q: what is the size of the Jacobian matrix?  
[4096 x 4096!]

in practice we process an entire minibatch (e.g. 100) of examples at one time:  
i.e. Jacobian would technically be a [409,600 x 409,600] matrix!  
(very sparse)

$$\text{e.g. } f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W_i \cdot x_i)^2$$



$$z = W \cdot x = \begin{bmatrix} W_{00}x_0 + \dots + W_{0n}x_n \\ W_{10}x_0 + \dots + W_{1n}x_n \\ \vdots \\ W_{m0}x_0 + \dots + W_{mn}x_n \end{bmatrix}$$

$$f(z) = \|z\|^2 = z_0^2 + \dots + z_n^2$$

$$\nabla_z f = 2z = \begin{bmatrix} 2z_0 \\ 2z_1 \\ \vdots \\ 2z_n \end{bmatrix} \in \mathbb{R}^2$$

$$\frac{\partial f}{\partial W_{ij}} = \sum_k \frac{\partial f}{\partial z_k} \frac{\partial z_k}{\partial W_{ij}} = \sum_k (2z_k) (x_j) = 2z_k x_j$$

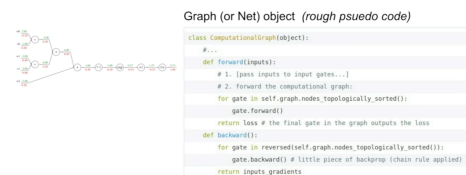
$$\nabla_W f = 2z \cdot x^T$$

$$\frac{\partial f}{\partial x_i} = W_{0i} x_0 + \dots + W_{ni} x_n = \sum_k W_{ki} x_k = W^T x$$

$$\nabla_x f = 2W^T z$$

Note that: the gradient with respect to a variable should be the same shape as the variable

Modularized implementation: forward / backward API

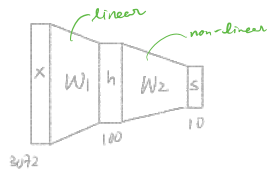


Graph (or Net) object (rough psuedo code)

```
class ComputationalGraphObject:
    def __init__(self):
        # 1. (zero inputs to input gates...)
        # 2. Forward the computational graph:
        for gate in self.graph.nodes, topologically, sorted():
            gate.forward()
        return loss at the final gate in the graph outputs the loss
    def backward(self):
        for gate in reversed(self.graph.nodes, topologically, sorted()):
            gate.backward() # little piece of backprop (chain rule applied)
        return inputs_gradients
```

## • Neural Network

- Before Linear score function :  $f = Wx$
- after 2-layer Neural Network:  $f = W_2 \max(0, W_1 x)$

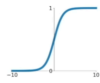


or even 3-layer ...

### Activation functions

#### Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



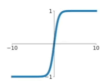
#### Leaky ReLU

$$\max(0.1x, x)$$



#### tanh

$$\tanh(x)$$

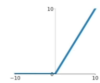


#### Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

#### ReLU

$$\max(0, x)$$



#### ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



### Neural networks: Architectures

