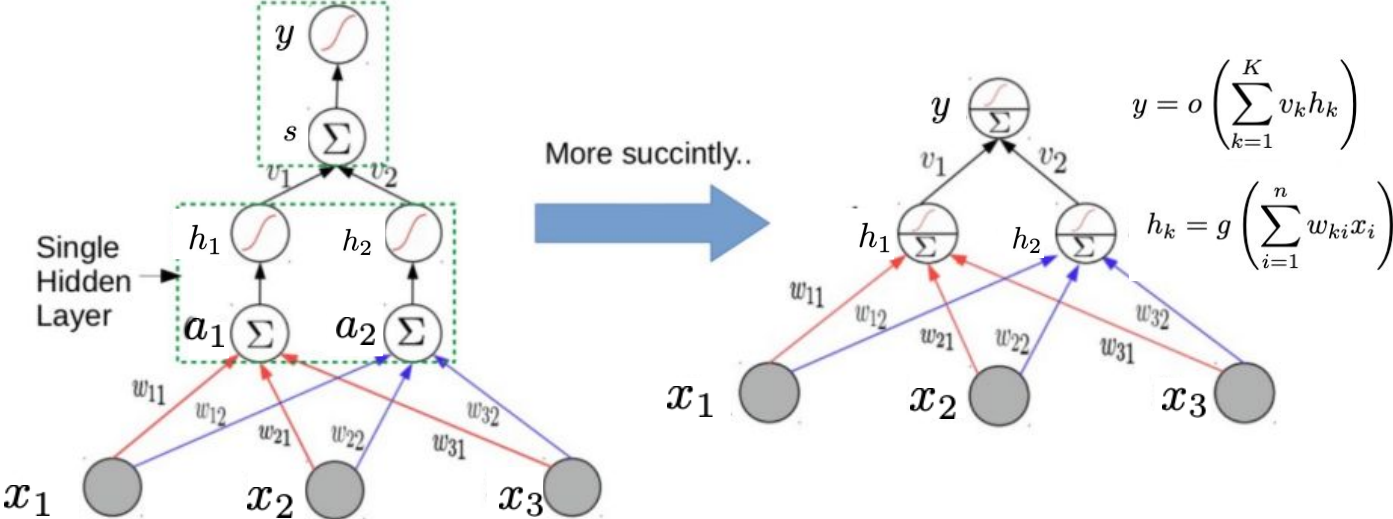


CS4641 Spring 2025

Neural Networks: Backpropagation

Bo Dai
School of CSE, Georgia Tech
bodai@cc.gatech.edu

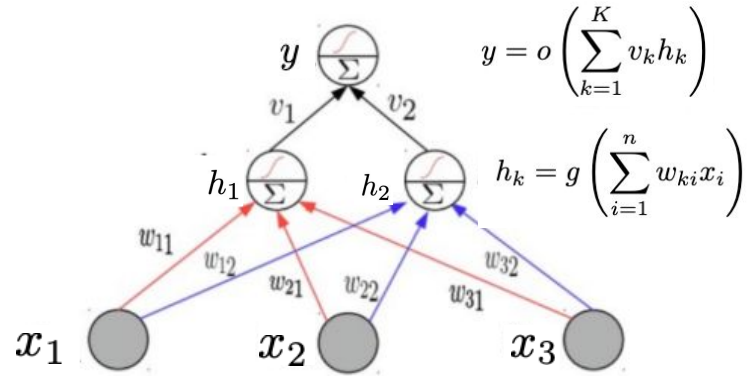
Neural Network Revisit



Vector Formation

$$W = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix}$$

$$\mathbf{x} = [x_1, x_2, x_3]^\top$$

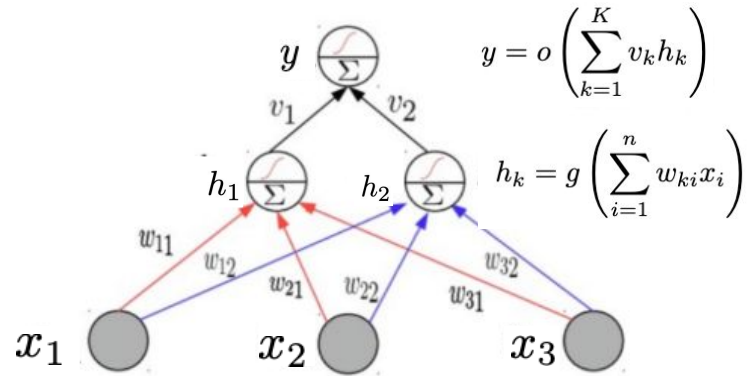


Vector Formation

$$h = [h_1, h_2]^\top = g(Wx)$$

$$W = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix}$$

$$x = [x_1, x_2, x_3]^\top$$



Vector Formation

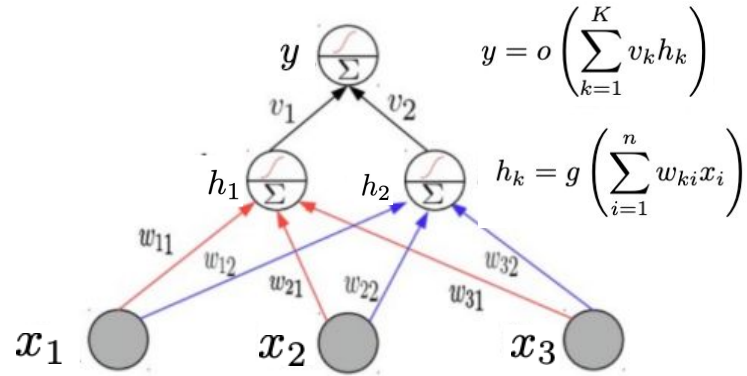
$$y = o(Vh)$$

$$V = [v_1, v_2]$$

$$h = [h_1, h_2]^T = g(Wx)$$

$$W = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix}$$

$$x = [x_1, x_2, x_3]^T$$



Vector Formation

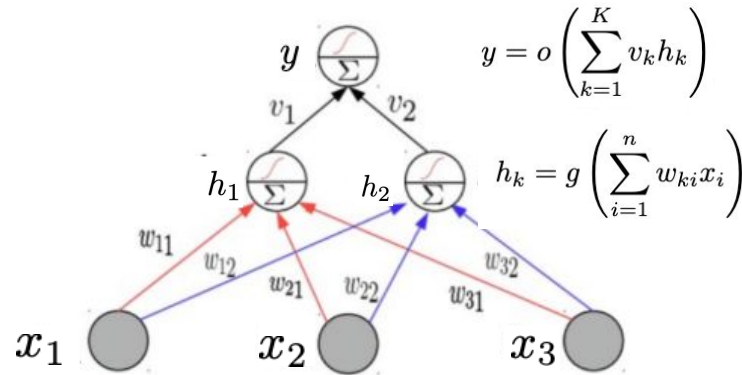
$$y = o(\mathbf{V}g(\mathbf{W}\mathbf{x}))$$

$$\mathbf{V} = [v_1, v_2]$$

$$\mathbf{h} = [h_1, h_2]^\top = g(\mathbf{W}\mathbf{x})$$

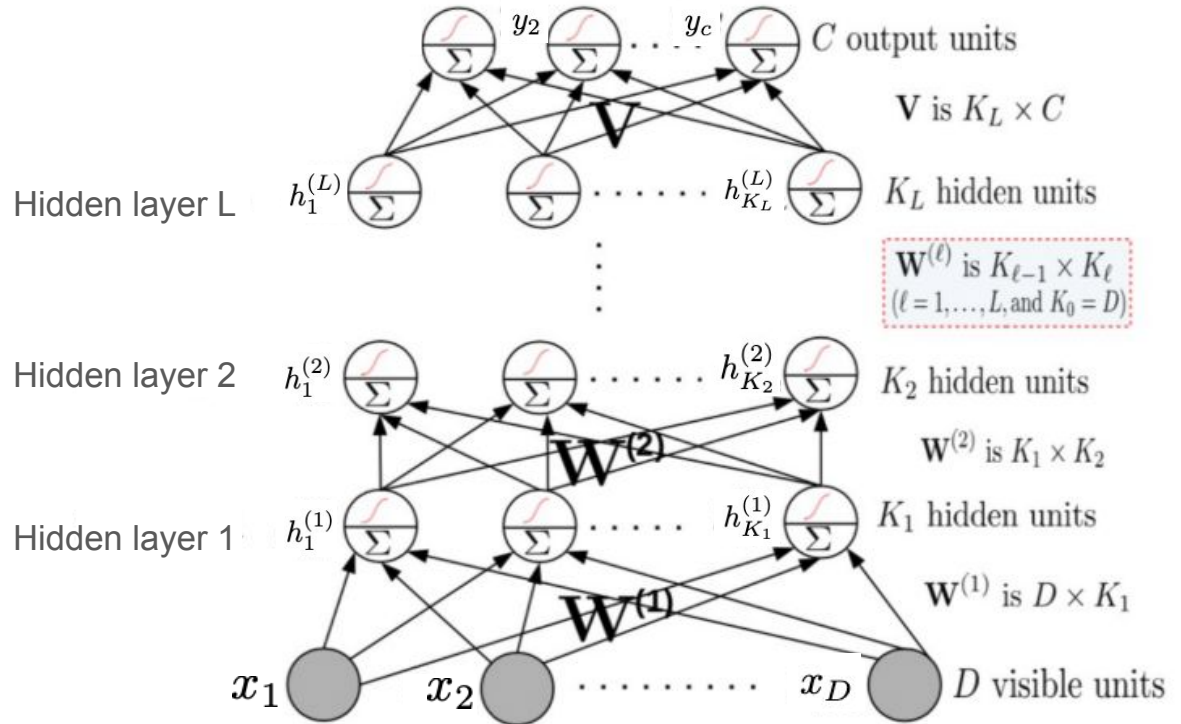
$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix}$$

$$\mathbf{x} = [x_1, x_2, x_3]^\top$$

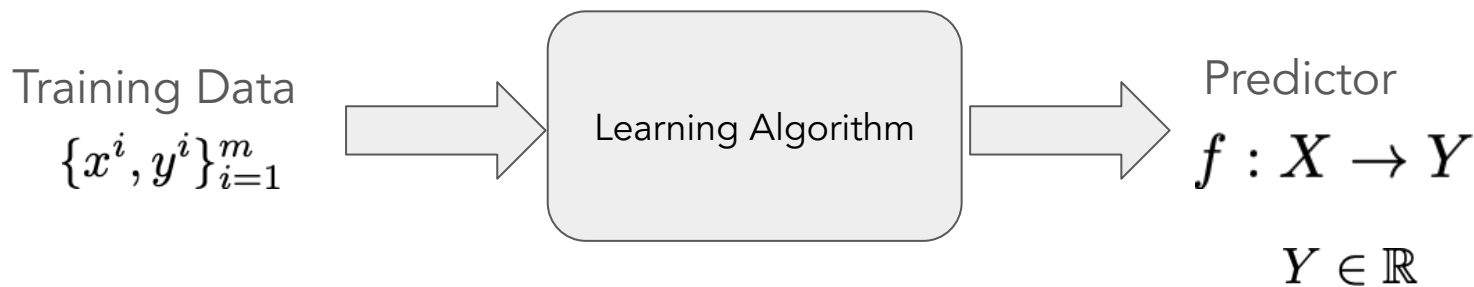


Multi-Layer Perception

$$y = f_L(\mathbf{W}_L f_{L-1}(\mathbf{W}_{L-1} \dots f_1(\mathbf{W}_1 \mathbf{x})))$$



Regression Algorithms

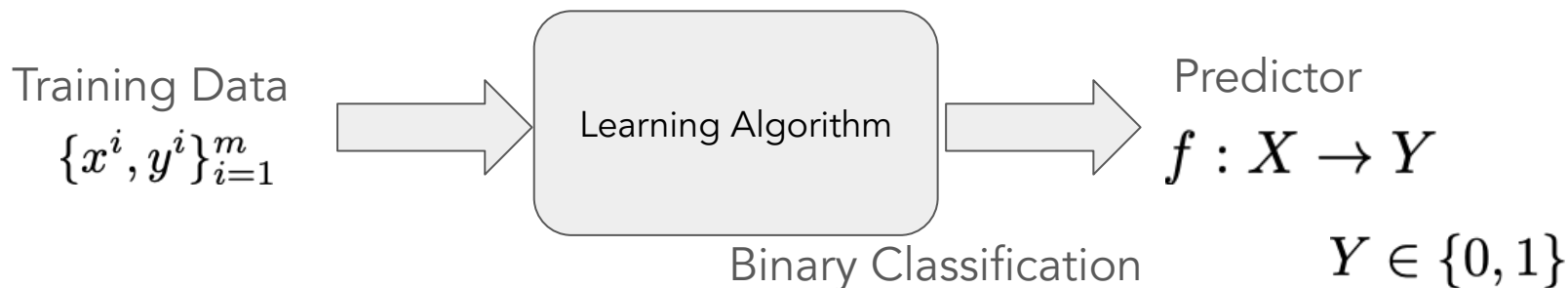


Linear Regression Pipeline

1. Build probabilistic models:
Gaussian Distribution + Neural Network
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) GD

$$y = o(\mathbf{V} g(\mathbf{W} x))$$

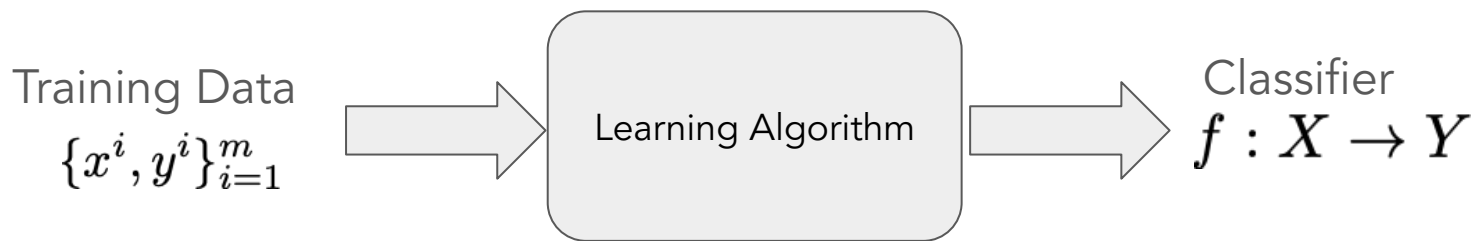
Binary Classification Algorithms



Binary Logistic Regression Pipeline

1. Build probabilistic models:
Bernoulli Distribution + Neural Network $y = o(\mathbf{V}g(\mathbf{W}x))$
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) Gradient Descent

Multiclass Logistic Regression Algorithms



Multiclass Classification $Y \in \{0, 1, \dots, k\}$
Multiclass Logistic Regression Pipeline

1. Build probabilistic models:
Categorical Distribution + Neural Network $y = o(\mathbf{V}g(\mathbf{W}x))$
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) Gradient Descent

Select Optimizer

$$L(\theta) = \sum_{i=1}^m \ell(x^i, y^i, \theta) + \lambda \Omega(\theta)$$

$$\theta = [V, W]$$

$$\ell(x^i, y^i, \theta) = (o(Vg(Wx^i)) - y^i)^2$$

Select Optimizer

$$L(\theta) = \sum_{i=1}^m \ell(x^i, y^i, \theta) + \lambda \Omega(\theta)$$
$$\theta = [V, W]$$

$$\ell(x^i, y^i, \theta) = (o(Vg(Wx^i)) - y^i)^2$$

$$\ell(x^i, y^i, \theta) = -y^i \log \sigma(o(Vg(Wx^i)))$$
$$-(1 - y^i) \log(1 - \sigma(o(Vg(Wx^i))))$$

Select Optimizer

$$L(\theta) = \sum_{i=1}^m \ell(x^i, y^i, \theta) + \lambda \Omega(\theta)$$
$$\theta = [V, W]$$

$$\ell(x^i, y^i, \theta) = (o(Vg(Wx^i)) - y^i)^2$$

$$\ell(x^i, y^i, \theta) = -y^i \log \sigma(o(Vg(Wx^i)))$$
$$-(1 - y^i) \log(1 - \sigma(o(Vg(Wx^i))))$$

$$\ell(x^i, y^i, \theta) = - \sum_{j=1}^k y^i \log \frac{\exp(o(V_j g(Wx^i)))}{\sum_{c=1}^k \exp(o(V_c g(Wx^i)))}$$

Select Optimizer

$$L(\theta) = \sum_{i=1}^m \ell(x^i, y^i, \theta) + \lambda \Omega(\theta)$$

$$\theta = [V, W]$$

$$\ell(x^i, y^i, \theta) = (o(Vg(Wx^i)) - y^i)^2$$

$$\ell(x^i, y^i, \theta) = -y^i \log \sigma(o(Vg(Wx^i))) \\ - (1 - y^i) \log(1 - \sigma(o(Vg(Wx^i))))$$

$$\ell(x^i, y^i, \theta) = - \sum_{j=1}^k y^i \log \frac{\exp(o(V_j g(Wx^i)))}{\sum_{c=1}^k \exp(o(V_c g(Wx^i)))}$$

- (Stochastic) Gradient Descent

(Stochastic) Gradient Descent

- Initialize parameter θ^0

- Sample $\{x^i, y^i\}_{i=1}^B$

- Do
$$\theta^{t+1} \leftarrow \theta^t - \eta \sum_{i=1}^B \nabla_{\theta} \ell(x^i, y^i, \theta^t) - (\lambda \nabla \Omega(\theta^t))$$

Chain Rule

- A **composite function** is the combination of two functions: a function that takes as input the output of another function



E.g, $f(\theta) = 2\theta + 1$, $g(\theta) = \theta^4$, $h(\theta) = (2\theta + 1)^4$

Let's call $u = f(\theta)$ the output of the inner function $\rightarrow h(\theta) = g(u)$

$$h' = \frac{dh}{d\theta} = \frac{dh}{du} \frac{du}{d\theta}$$

$$\frac{dh}{du} = 4(2\theta + 1)^3$$

$$\frac{du}{d\theta} = 2$$

$$h' = 8(2\theta + 1)^3$$

Derivative of outer part of $h(\theta)$

Derivative of inner part of $h(\theta)$

Chain Rule

$$h(\theta) = g(f(\theta))$$

$$u = f(\theta)$$

where $\theta \in \mathbb{R}^m$, $u \in \mathbb{R}^n$

$$h : \mathbb{R}^m \rightarrow \mathbb{R}$$

$$f : \mathbb{R}^m \rightarrow \mathbb{R}^n$$

$$g : \mathbb{R}^n \rightarrow \mathbb{R}$$

Vector function

The inner vector function f maps m inputs to n outputs, while the outer function g receives n inputs to produce one output, h .

The chain rule allows to compute the variation (i.e., the partial derivative) of the function w.r.t. each component of the multivariate input \rightarrow **Gradient vector** of $h(\theta)$

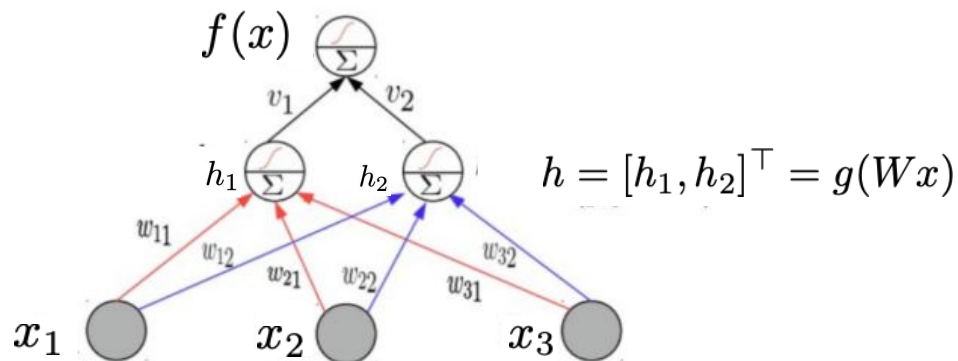
$$\frac{\partial h}{\partial \theta_i} = \frac{\partial h}{\partial u_1} \frac{\partial u_1}{\partial \theta_i} + \frac{\partial h}{\partial u_2} \frac{\partial u_2}{\partial \theta_i} + \dots + \frac{\partial h}{\partial u_n} \frac{\partial u_n}{\partial \theta_i} = \sum_{j=1}^n \frac{\partial h}{\partial u_j} \frac{\partial u_j}{\partial \theta_i}$$

$$i = 1, \dots, m$$

$$\nabla h(\theta) = \left(\frac{\partial h}{\partial \theta_1}, \frac{\partial h}{\partial \theta_2}, \dots, \frac{\partial h}{\partial \theta_m} \right)^T$$

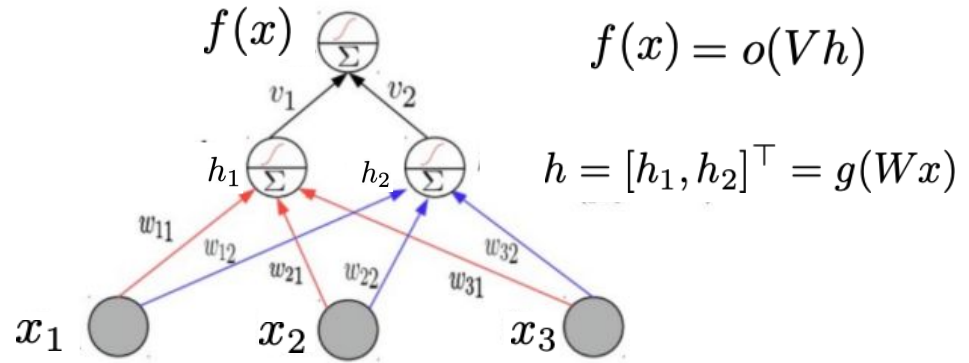
Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$



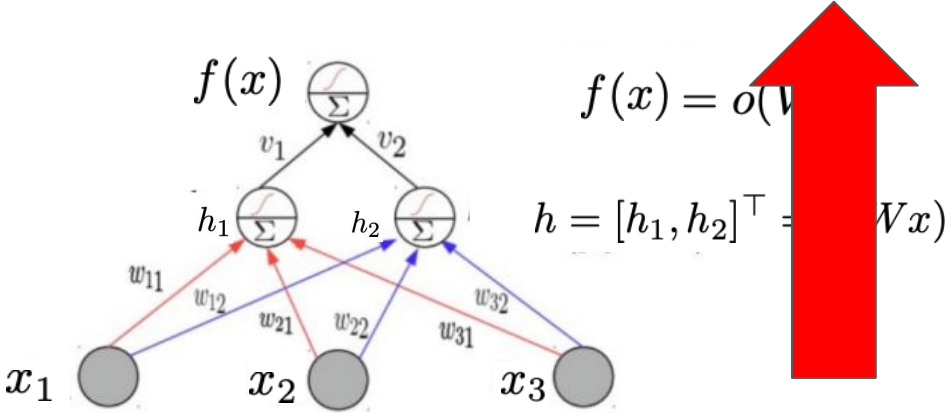
Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$



Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

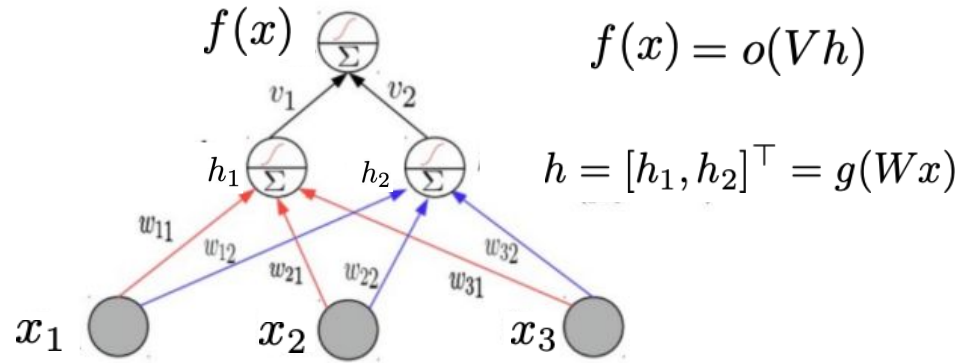


Forward Pass

Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\nabla_{\theta} \ell(x^i, y^i, \theta)$$

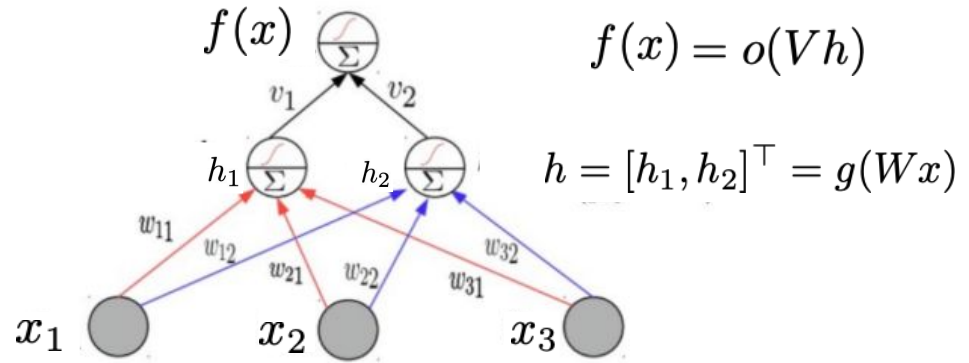


Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\nabla_{\theta} \ell(x^i, y^i, \theta) = \left[\frac{\partial \ell(x^i, y^i, \theta)}{\partial V}, \frac{\partial \ell(x^i, y^i, \theta)}{\partial W} \right]$$

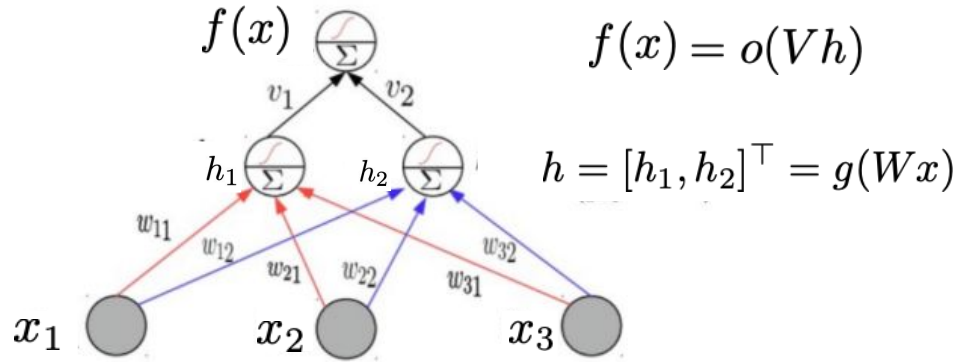
$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial V} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial f}{\partial V}$$



Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

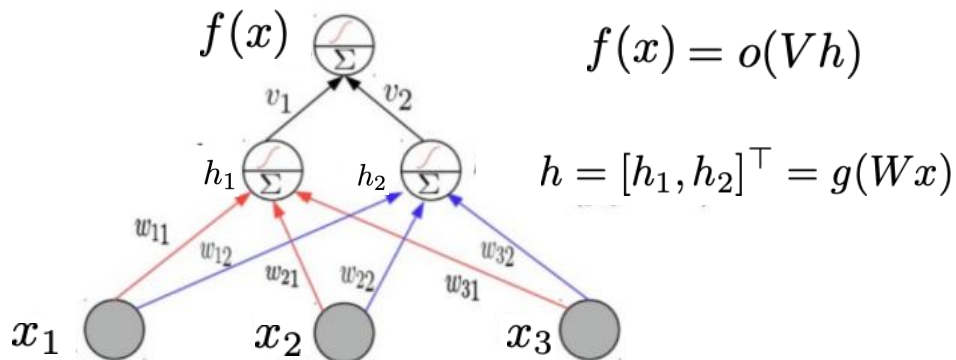
$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial V} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial f}{\partial V}$$
$$\frac{\partial f(x)}{\partial V} = \frac{\partial o(Vh)}{\partial V}$$



Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\begin{aligned} \frac{\partial \ell(x^i, y^i, \theta)}{\partial V} &= \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial f}{\partial V} \\ &= \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial o(Vh)}{\partial V} \end{aligned}$$

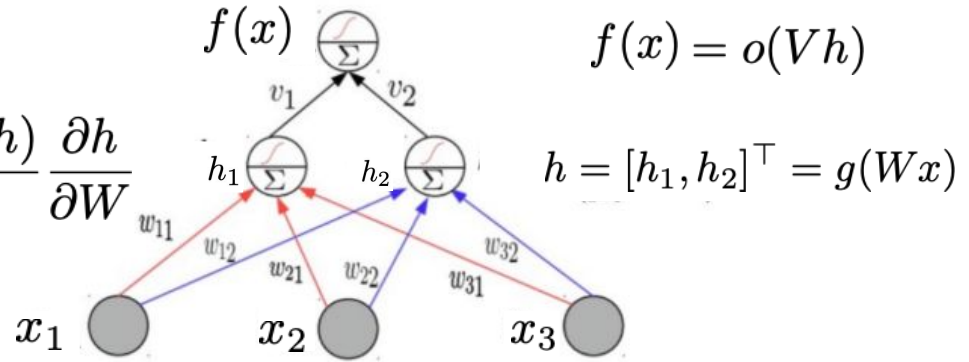


Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial \mathbf{W}} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial f}{\partial \mathbf{W}}$$

$$\frac{\partial f}{\partial \mathbf{W}} = \frac{\partial o(\mathbf{V}h)}{\partial h} \frac{\partial h}{\partial \mathbf{W}}$$



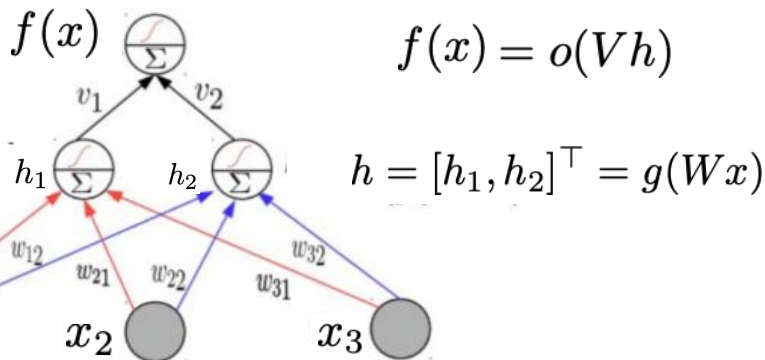
Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial W} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial f}{\partial W}$$

$$\frac{\partial f}{\partial W} = \frac{\partial o(Vh)}{\partial h} \frac{\partial h}{\partial W}$$

$$\frac{\partial h}{\partial W} = \frac{\partial g(Wx)}{\partial W} x_1$$



Backpropagation: Chain Rule on Neural Network

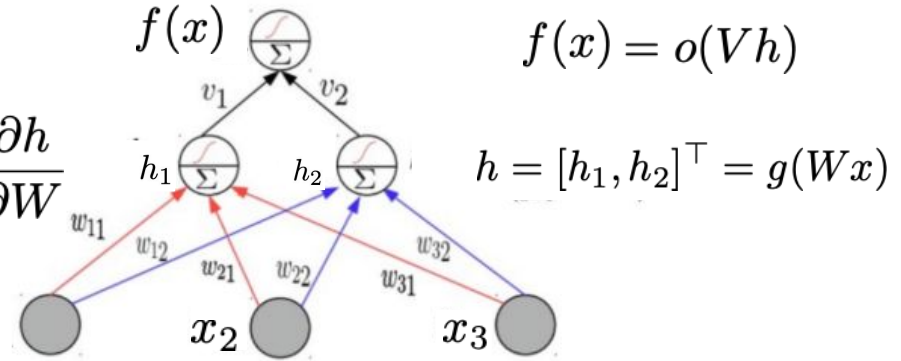
$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial \mathbf{W}} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial f}{\partial \mathbf{W}}$$

$$\frac{\partial f}{\partial \mathbf{W}} = \frac{\partial o(\mathbf{V}h)}{\partial h} \frac{\partial h}{\partial \mathbf{W}}$$

$$\frac{\partial h}{\partial \mathbf{W}} = \frac{\partial g(\mathbf{W}x)}{\partial \mathbf{W}} x_1$$

$$= \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial o(\mathbf{V}h)}{\partial h} \frac{\partial h}{\partial \mathbf{W}}$$

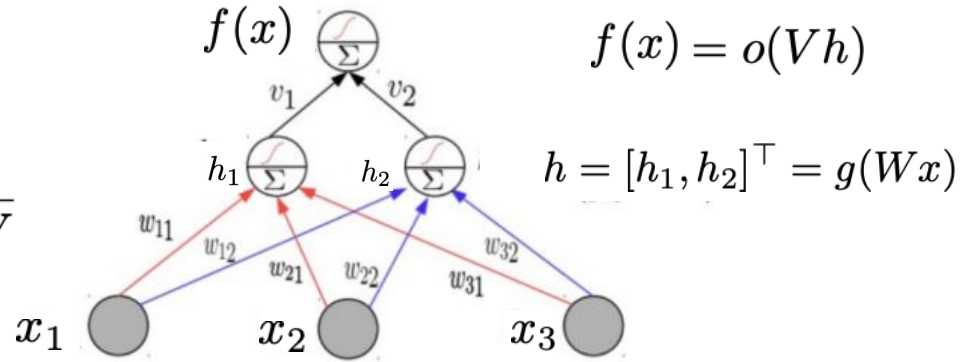


Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial V} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial o(Vh)}{\partial V}$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial W} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial o(Vh)}{\partial h} \frac{\partial h}{\partial W}$$

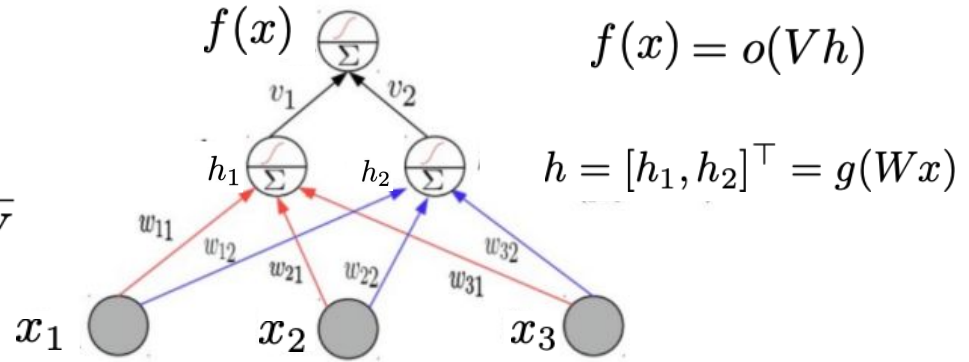


Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial V} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial o(Vh)}{\partial V}$$

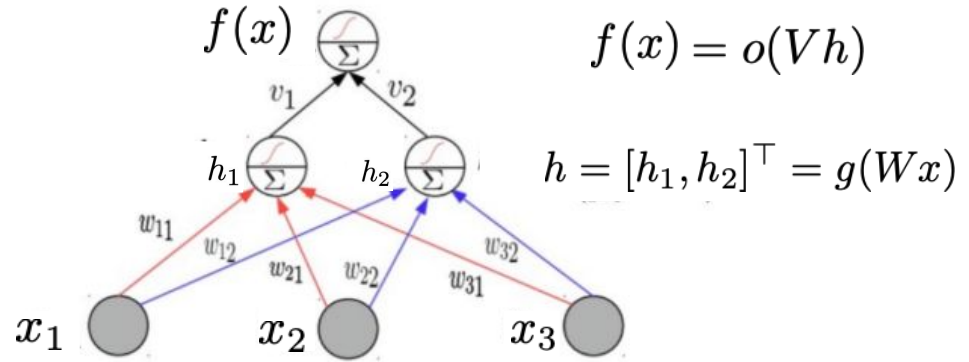
$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial W} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial o(Vh)}{\partial h} \frac{\partial h}{\partial W}$$



Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

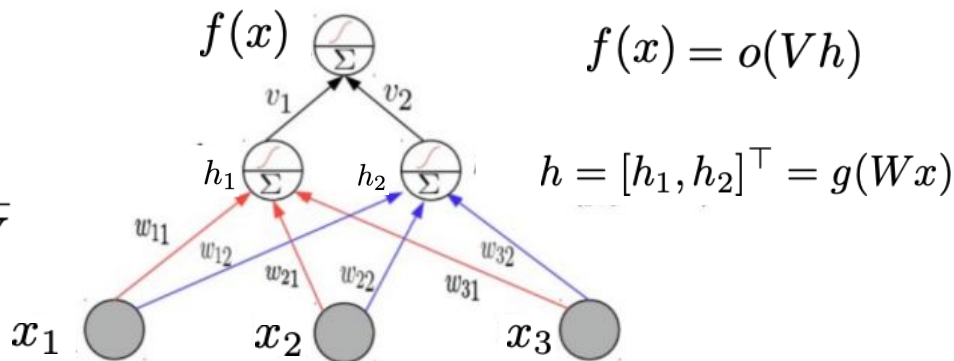
$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial f}$$



Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial V} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial o(Vh)}{\partial V}$$
$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial W} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial o(Vh)}{\partial h} \frac{\partial h}{\partial W}$$

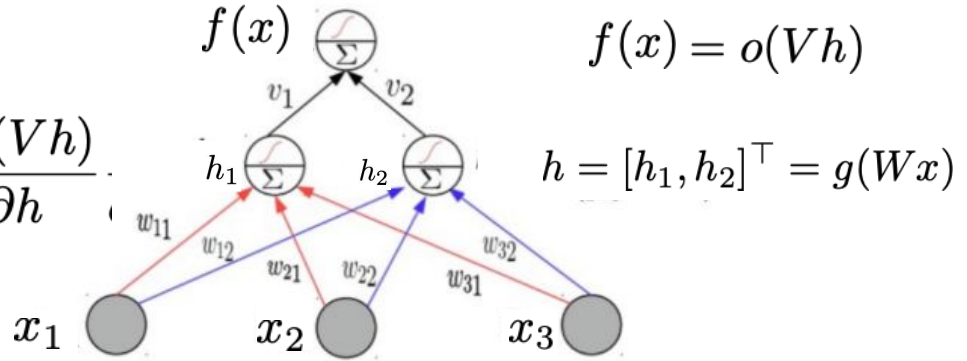


Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial f}$$

$$\frac{\partial o(Vh)}{\partial V} \quad \frac{\partial o(Vh)}{\partial h}$$

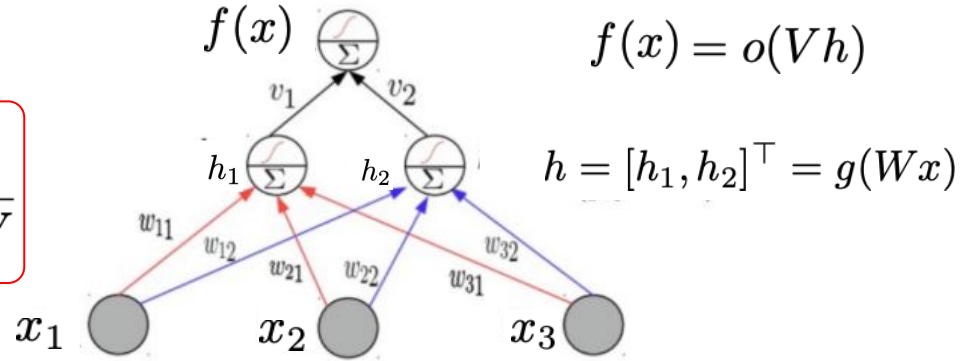


Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial \mathbf{V}} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial o(\mathbf{V}h)}{\partial \mathbf{V}}$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial \mathbf{W}} = \frac{\partial \ell(x^i, y^i, \theta)}{\partial f} \frac{\partial o(\mathbf{V}h)}{\partial h} \frac{\partial h}{\partial \mathbf{W}}$$



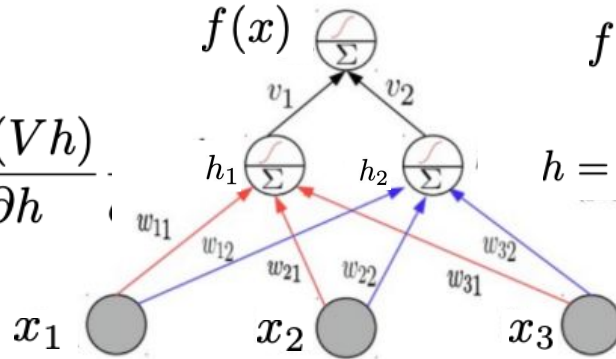
Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial f}$$

$$\frac{\partial o(Vh)}{\partial V} \quad \frac{\partial o(Vh)}{\partial h}$$

$$\frac{\partial h}{\partial W}$$



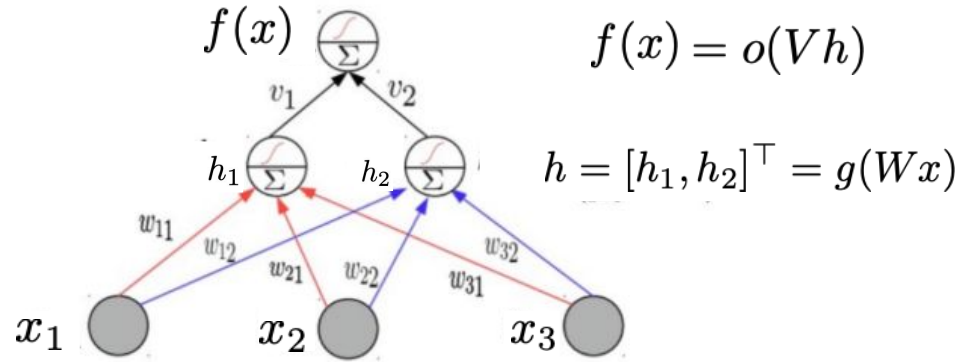
$$f(x) = o(Vh)$$

$$h = [h_1, h_2]^T = g(Wx)$$

Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial f}$$

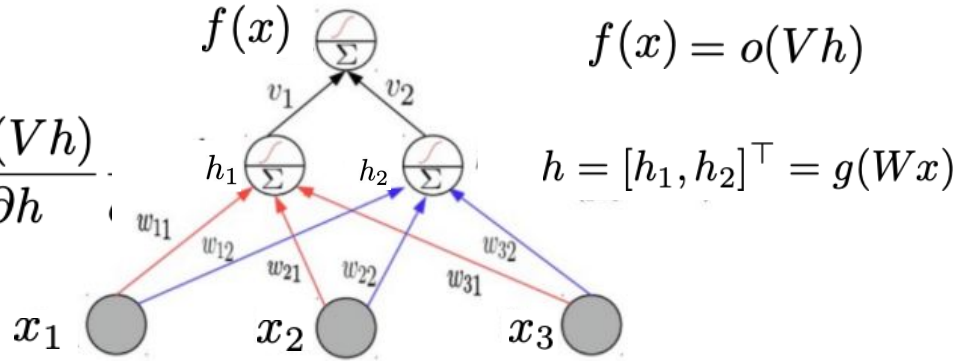


Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial f}$$

$$\frac{\partial o(Vh)}{\partial V} \quad \frac{\partial o(Vh)}{\partial h}$$



Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$

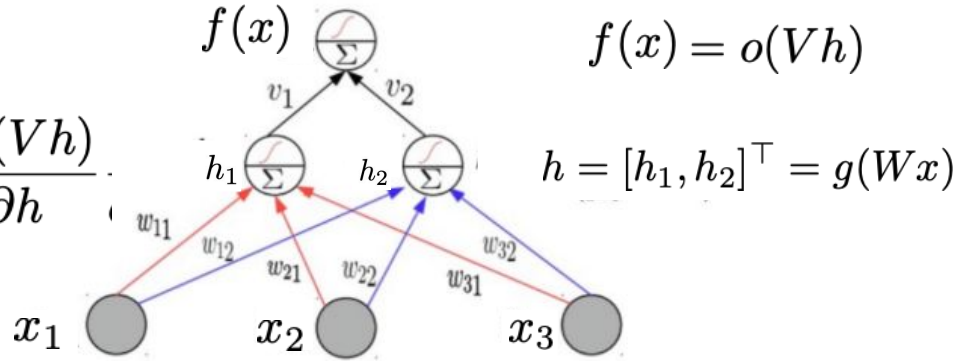
$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial f}$$

$$\frac{\partial f}{\partial o(Vh)}$$

$$\frac{\partial o(Vh)}{\partial V}$$

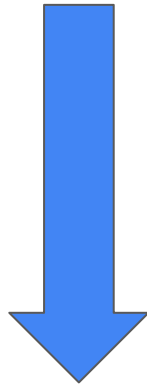
$$\frac{\partial o(Vh)}{\partial h}$$

$$\frac{\partial h}{\partial W}$$



Backpropagation: Chain Rule on Neural Network

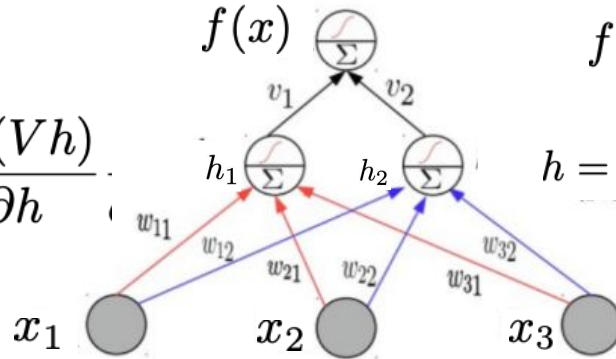
$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$



$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial f}$$

$$\frac{\partial o(Vh)}{\partial V} \frac{\partial o(Vh)}{\partial h}$$

$$\frac{\partial h}{\partial W}$$



$$f(x) = o(Vh)$$

$$h = [h_1, h_2]^T = g(Wx)$$

Backward Pass

(Stochastic) Gradient Descent

- Initialize parameter θ^0
- Sample $\{x^i, y^i\}_{i=1}^B$
- Do
$$\theta^{t+1} \leftarrow \theta^t - \eta \sum_{i=1}^B \nabla_{\theta} \ell(x^i, y^i, \theta^t) - (\lambda \nabla \Omega(\theta^t))$$

Auto-differentiation Packages

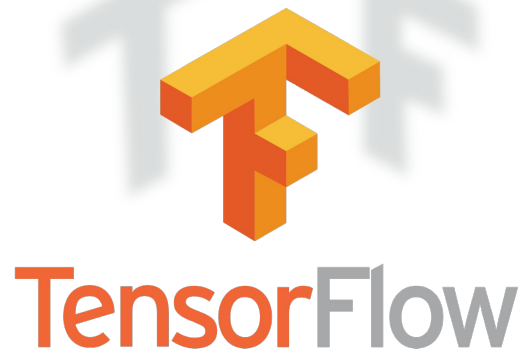
PyTorch



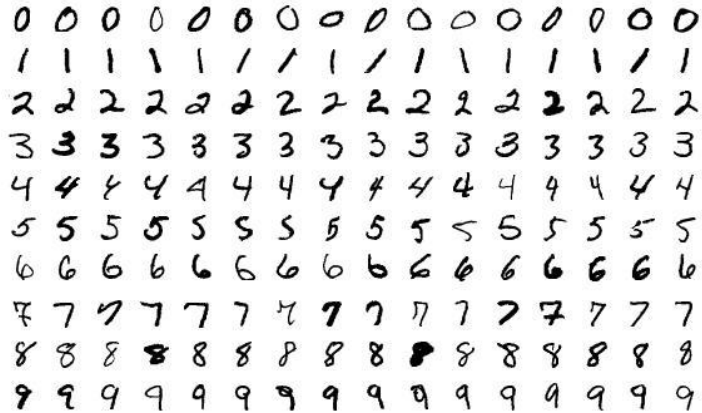
JAX



Tensorflow

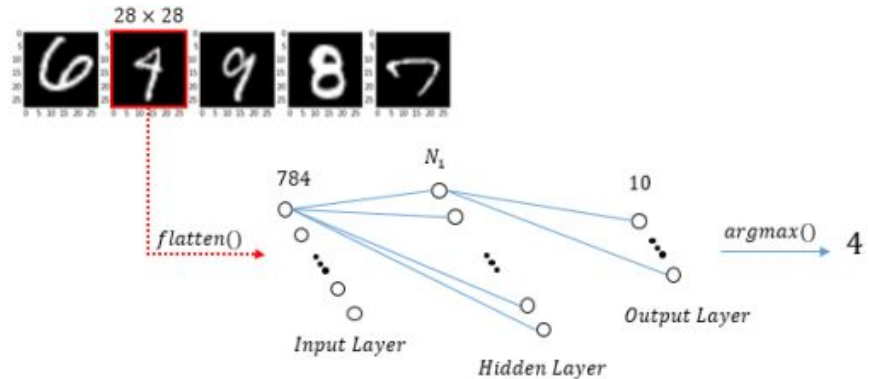


MLP example: MNIST



- 60,000 images
- 28x28 pixels = 784
- Grayscale, from 0 to 255 → Converted to [0,1]

MNIST hand-written character recognition



PyTorch

```
▶ #@title Define model class

class Net(nn.Module):
    def __init__(self, input_size, hidden_size, num_classes):
        super(Net,self).__init__()
        self.fc1 = nn.Linear(input_size, hidden_size)
        self.relu = nn.ReLU()
        self.fc2 = nn.Linear(hidden_size, num_classes)

    def forward(self,x):
        out = self.fc1(x)
        out = self.relu(out)
        out = self.fc2(out)
        return out
```

```
#@title Define loss-function & optimizer
```

```
loss_function = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam( net.parameters(), lr=lr)
```

PyTorch

```
#@title Training the model

for epoch in range(num_epochs):
    for i ,(images,labels) in enumerate(train_gen):
        images = Variable(images.view(-1,28*28)).cuda()
        labels = Variable(labels).cuda()

        optimizer.zero_grad()
        outputs = net(images)
        loss = loss_function(outputs, labels)
        loss.backward()
        optimizer.step()
```

Q&A