

Deep Learning

Lecture 3

Recap

- Classification
 - Logistic Regression
 - Sigmoid function
 - Cross entropy loss function
- Neural Networks
 - Perceptron
 - Layers
 - Non-Linear activation function

Agenda

- Neural Network Learning
 - Feed forward
 - Back propagation
- Practical Machine Learning
 - Data distribution
 - Overfitting
 - Regularization
 - Data Normalization
 - Hyper parameter Tuning

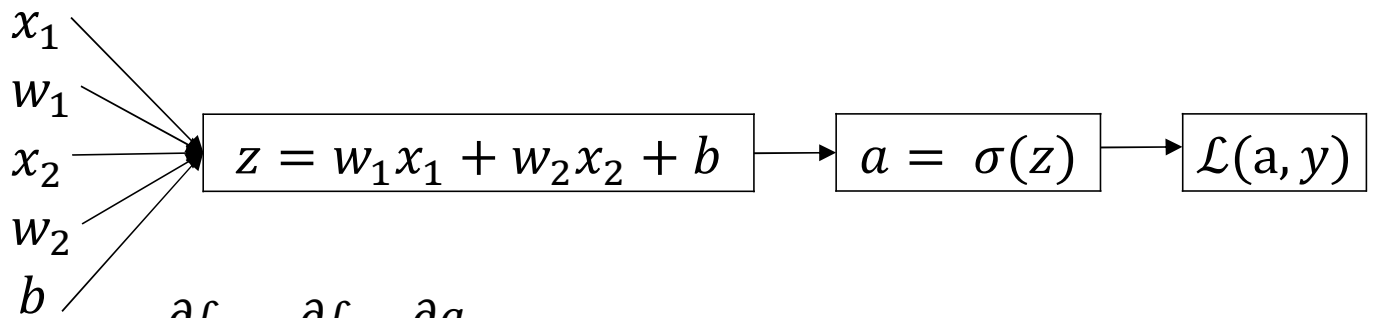
Computation Graph: Logistic Regression and Gradient

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$$\mathcal{L}(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$$

Logistic regression derivatives



$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial a} * \frac{\partial a}{\partial z} = a - y$$

$$\frac{\partial a}{\partial z} = a(1 - a) \qquad \frac{\partial \mathcal{L}}{\partial a} = -\frac{y}{a} + \frac{1 - y}{1 - a}$$

Logistic regression on m examples

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^i, y^i)$$

$$a^i = \hat{y} = \sigma(z^i) = \sigma(w^T x^i + b)$$

$$\frac{\partial}{\partial w_i} J(\theta, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_i} \mathcal{L}(a^i, y^i)$$

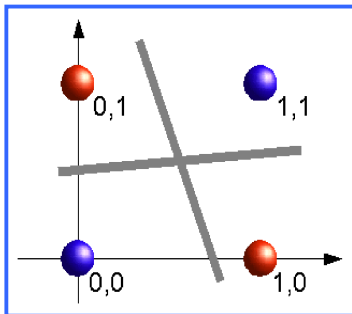
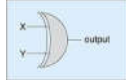
Artificial Neural Networks

Learning highly non-linear functions

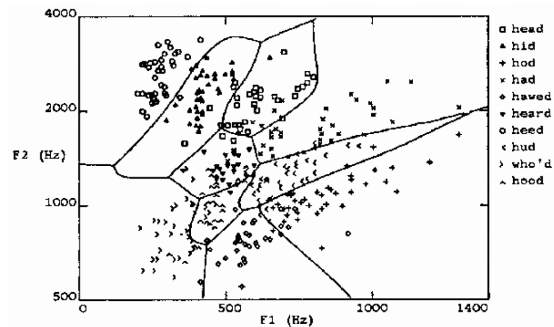
$f: X \rightarrow Y$

- f might be non-linear function
- X (vector of) continuous and/or discrete vars
- Y (vector of) continuous and/or discrete vars

The XOR gate



Speech recognition



Background

A Recipe for Machine Learning

1. Given training data 1:m

$$\{\mathbf{x}_i, \mathbf{y}_i\}$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

3. Define goal:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1} \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

4. Train with SGD:

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

A Recipe for Machine Learning

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

3. Define goal:

2. Choose each of the

– Decision function

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

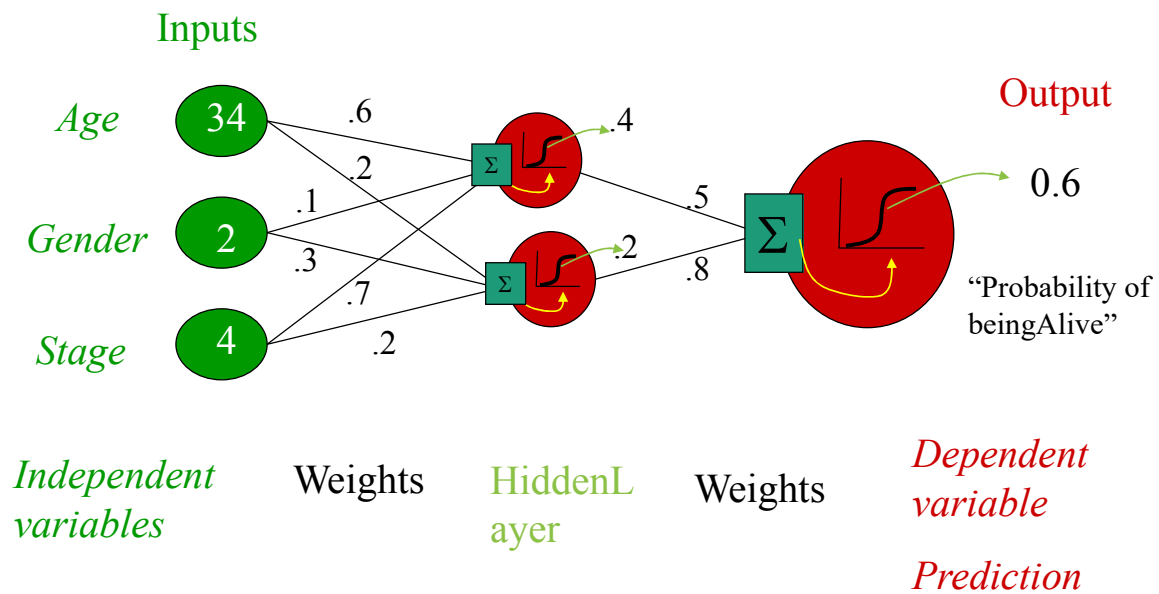
Gradients

Backpropagation can compute this gradient!

the gradient)

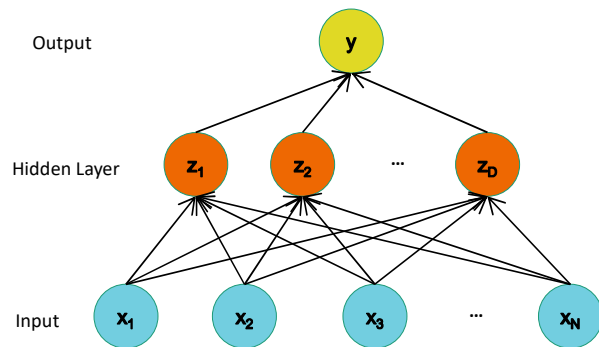
$$\boldsymbol{\theta}^{(t)} \rightarrow \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

Neural Network Model



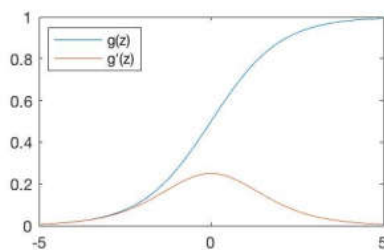
Activation Functions

Neural Network with arbitrary
nonlinear activation functions



Common Activation Functions

Sigmoid Function



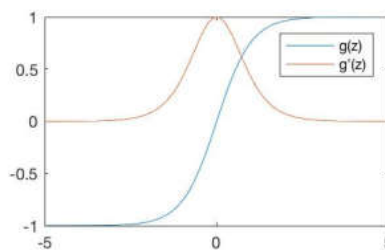
$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$



tf.nn.sigmoid(z)

Hyperbolic Tangent



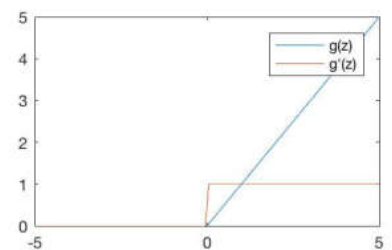
$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$



tf.nn.tanh(z)

Rectified Linear Unit (ReLU)



$$g(z) = \max(0, z)$$

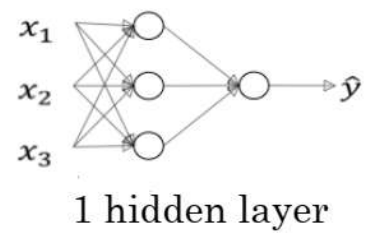
$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

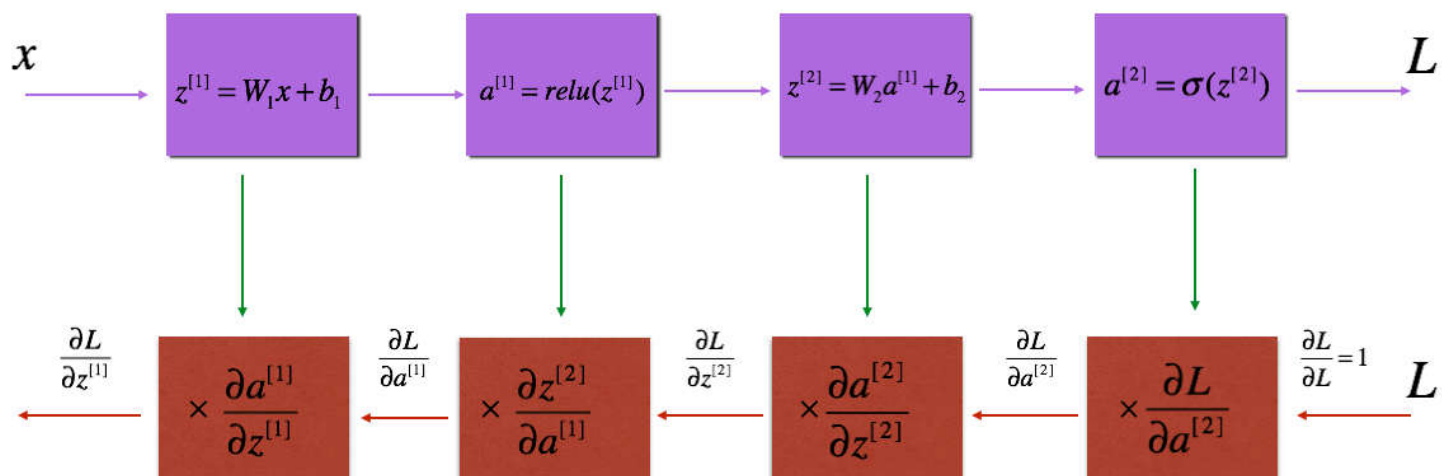


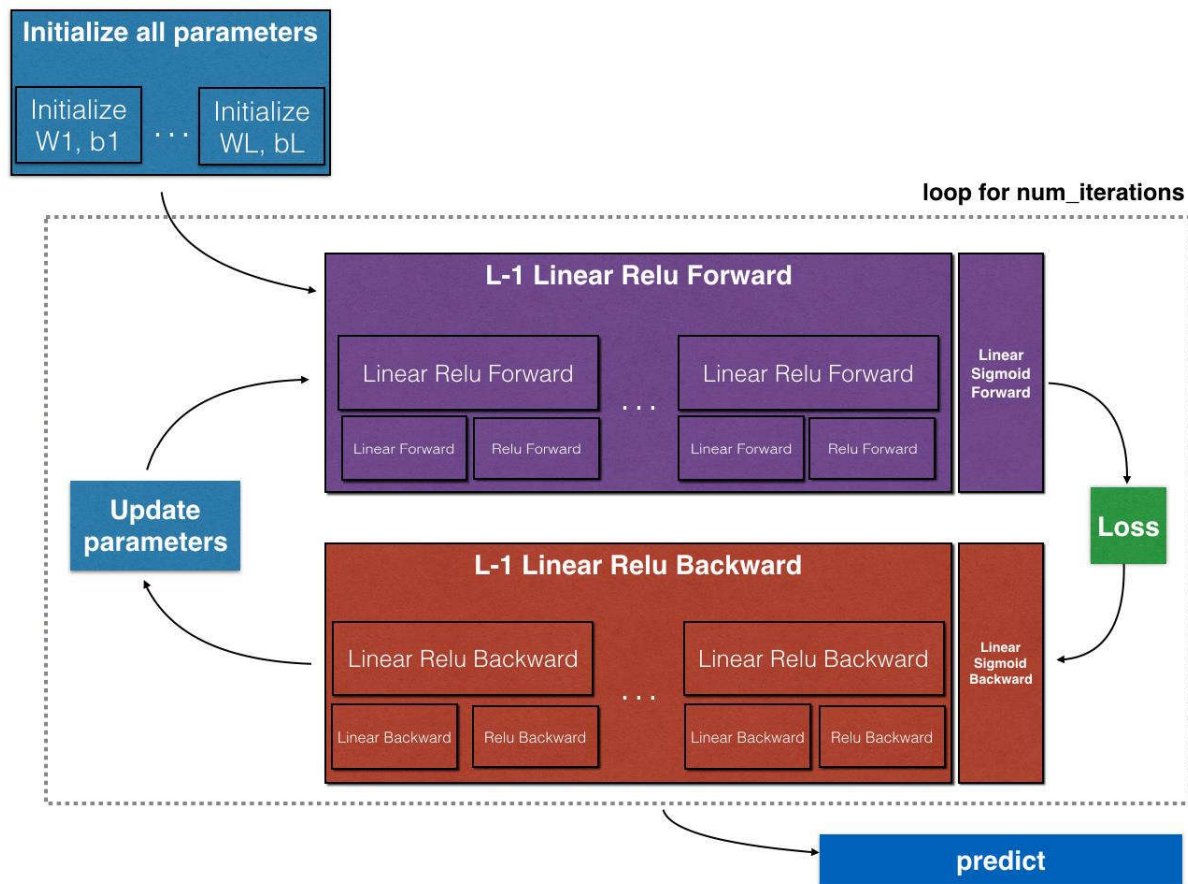
tf.nn.relu(z)

Parameters $W^{[l]}$ and $b^{[l]}$

- $z^{[1]} = W^{[1]}X + b^{[1]}$
- $a^{[1]} = g^{[1]}(z^{[1]})$
- $z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$
- $\hat{y} = a^{[2]} = g^{[2]}(z^{[2]})$







Forward and backpropagation

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(Z^{[2]})$$

$$\vdots$$

$$A^{[L]} = g^{[L]}(Z^{[L]}) = \hat{Y}$$

$$dZ^{[L]} = A^{[L]} - Y$$

$$dW^{[L]} = \frac{1}{m} dZ^{[L]} A^{[L]T}$$

$$db^{[L]} = \frac{1}{m} np.sum(dZ^{[L]}, axis = 1, keepdims = True)$$

$$dZ^{[L-1]} = dW^{[L]T} dZ^{[L]} g'^{[L]}(Z^{[L-1]})$$

$$\vdots$$
$$dZ^{[1]} = dW^{[L]T} dZ^{[2]} g'^{[1]}(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} A^{[1]T}$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True)$$

Summary of Neural Network

1. Neural Networks...

- provide a way of learning features
- are highly nonlinear prediction functions
- (can be) a highly parallel network of logistic regression classifiers
- discover useful hidden representations of the input

2. Backpropagation...

- provides an efficient way to compute gradients
- is a special case of reverse-mode automatic differentiation