

Data Science

Week 9

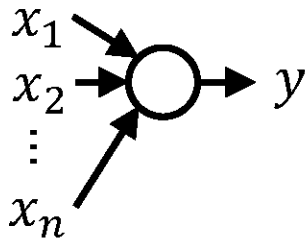
Fitting a Model to Data (2)-

Classification:

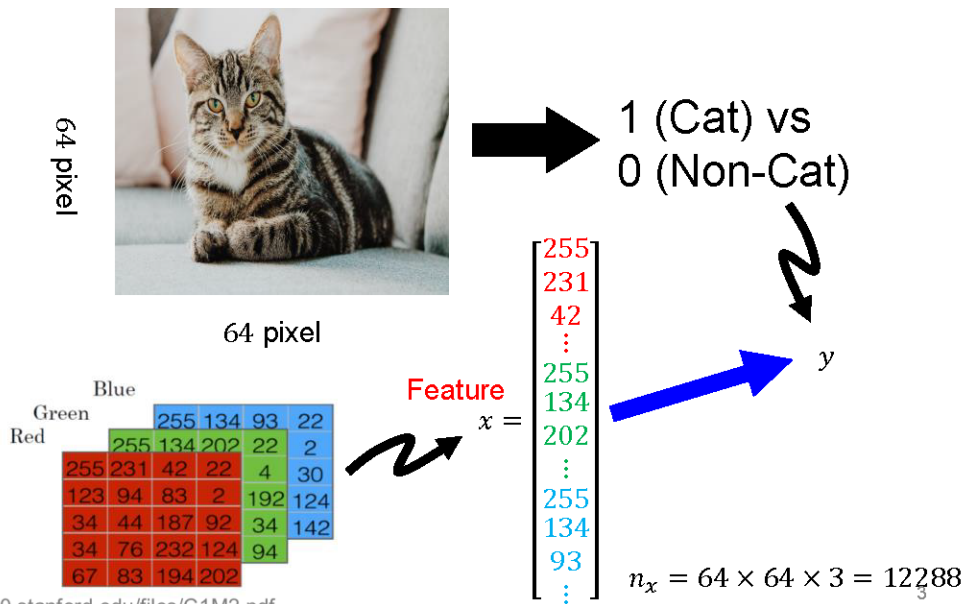
Neural Networks

Basics of Neural Network

Binary Classification



Binary Classification



Basics of Neural Network

Logistic Regression

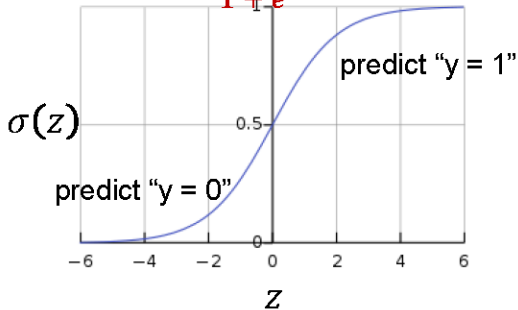
Logistic Regression

Given \mathbf{x}, y , want $\hat{y} = P(y = 1|\mathbf{x})$ infinitely close to y ,

Where, $\mathbf{x} \in \mathbb{R}^{n_x}$, $0 \leq \hat{y} \leq 1$, output predicted value $\hat{y} = \sigma(\mathbf{w}^T \mathbf{x} + b)$

Unknown Parameters: $\mathbf{w} \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

Sigmoid: $\sigma(z) = \frac{1}{1 + e^{-z}}$



$$z \rightarrow +\infty, \sigma(z) \approx \frac{1}{1 + e^{-\infty}} = 1$$

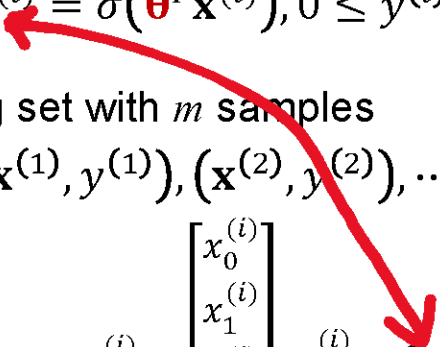
$$z \rightarrow -\infty, \sigma(z) \approx \frac{1}{1 + e^{+\infty}} = 0$$

Basics of Neural Network

Logistic Regression Loss Function

Logistic Regression Loss Function

Logistic Regression $\hat{y}^{(i)} = \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)}), 0 \leq \hat{y}^{(i)} \leq 1, \quad \sigma(z) = \frac{1}{1 + e^{-z}}$



Training set with m samples

$$(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})$$

$$\mathbf{x}^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}, x_0^{(i)} = 1, y^{(i)} \in \{0, 1\}$$

How to choose parameters $\boldsymbol{\theta}$ to make $\hat{y}^{(i)}$ infinitely close to $y^{(i)}$?

Logistic Regression Loss Function

Loss (error) function:

$$L(\hat{y}^{(i)}, y^{(i)}) = -y^{(i)}\ln(\hat{y}^{(i)}) - (1 - y^{(i)})\ln(1 - \hat{y}^{(i)})$$

$L(\hat{y}^{(i)}, y^{(i)}) \geq 0$

if $y^{(i)} = 1, (1 - y^{(i)}) = 0$

when $L(\hat{y}^{(i)}, y^{(i)}) = -\ln(\hat{y}^{(i)}) \rightarrow 0, \hat{y}^{(i)} \rightarrow 1$

if $y^{(i)} = 0, y^{(i)} = 0$

when $L(\hat{y}^{(i)}, y^{(i)}) = -\ln(1 - \hat{y}^{(i)}) \rightarrow 0, \hat{y}^{(i)} \rightarrow 0$

Logistic Regression

$$\begin{aligned} J(\boldsymbol{\theta}) &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \\ &= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \ln(\hat{y}^{(i)}) + (1 - y^{(i)}) \ln(1 - \hat{y}^{(i)}) \right] \end{aligned}$$

Learning: find parameter $\boldsymbol{\theta}$ to $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

Prediction: given new x output $\hat{y} = \sigma(\boldsymbol{\theta}^T \mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$

Basics of Neural Network

Gradient Descent
for Logistic Regression

Gradient Descent

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \ln(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \ln(1 - h_{\theta}(x^{(i)})) \right]$$

Goal: $\min_{\theta} J(\theta)$

Good news: Convex function!

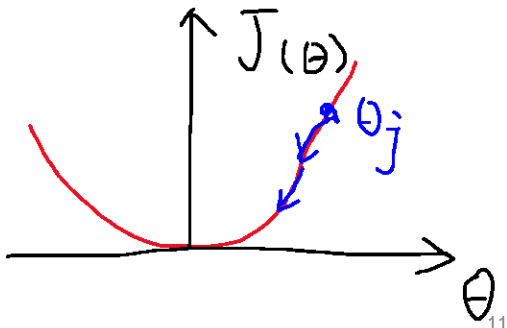
Bad news: No analytical solution

Repeat

{

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

}



Gradient Descent

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \ln(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \ln(1 - h_{\theta}(x^{(i)})) \right]$$

Goal: $\min_{\theta} J(\theta)$

Good news: Convex function!

Bad news: No analytical solution

Repeat

{
 $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$ (Simultaneously update all θ_j)
}

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Basics of Neural Network

Implementation
of Logistic Regression

Implementing Logistic Regression

$w_1 = 0, w_2 = 0, b = 0, \alpha = 4$

for $t = 1$ to k :

$J = 0, dw_1 = 0, dw_2 = 0, db = 0$

for $i = 1$ to m :

$$z^{(i)} = w_1 x_1^{(i)} + w_2 x_2^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \ln a^{(i)} + (1 - y^{(i)}) \ln(1 - a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J = J/m, dw_1 = dw_1/m, dw_2 = dw_2/m, db = db/m$$

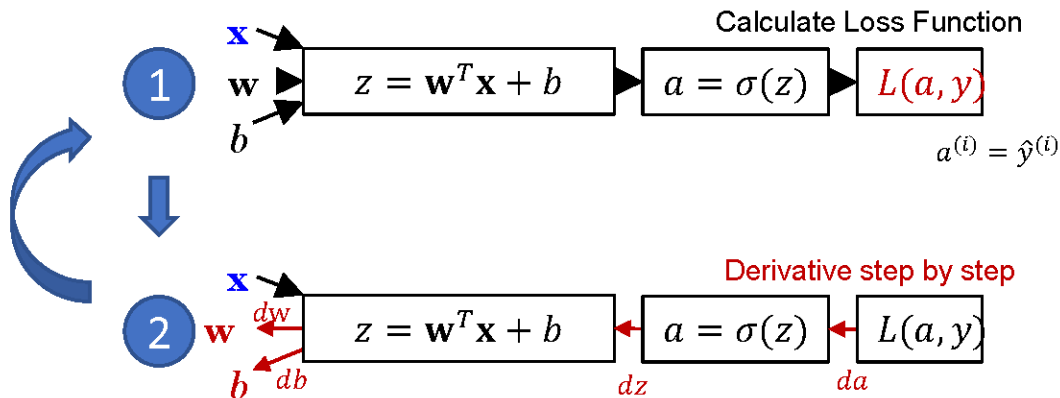
$$w_1 = w_1 - \alpha dw_1, w_2 = w_2 - \alpha dw_2, b = b - \alpha db$$

k is the number of iterations

m is the number of sample data

Training using Logistic Regression

Gradient Descent: process to estimate the parameter (\mathbf{w} , b)



$$d\mathbf{w}_1 = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)}) x_1^{(i)}$$

$$d\mathbf{w}_2 = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)}) x_2^{(i)}$$

...

$$db = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)})$$

Gradient Descent

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

$$-\frac{y^{(i)}}{a^{(i)}} + \frac{1-y^{(i)}}{1-a^{(i)}} \quad a^{(i)}(1-a^{(i)}) \quad x_1^{(i)}$$

$$dw_1 = \sum_{i=1}^m \frac{dL}{da^{(i)}} \times \frac{da^{(i)}}{dz^{(i)}} \times \frac{dz^{(i)}}{dw_1} = \sum_{i=1}^m (a^{(i)} - y^{(i)}) x_1^{(i)}$$

$$dw_2 = \sum_{i=1}^m \frac{dL}{da^{(i)}} \times \frac{da^{(i)}}{dz^{(i)}} \times \frac{dz^{(i)}}{dw_2} = \sum_{i=1}^m (a^{(i)} - y^{(i)}) x_2^{(i)}$$

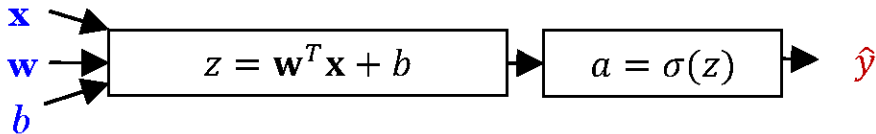
$$\vdots$$

$$db = \sum_{i=1}^m \frac{dL}{da^{(i)}} \times \frac{da^{(i)}}{dz^{(i)}} \times \frac{dz^{(i)}}{db} = \sum_{i=1}^m (a^{(i)} - y^{(i)})$$

$$dw_1 = \frac{1}{m} dw_1, \quad dw_2 = \frac{1}{m} dw_2, \dots \quad db = \frac{1}{m} db$$

Test New Data using Logistic Model

Prediction (Test): given new \mathbf{x} output $\hat{y} = \sigma(\mathbf{w}^T \mathbf{x} + b) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x} + b}}$



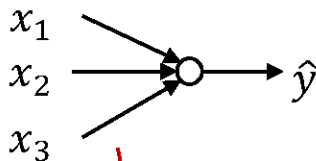
Shallow Neural Networks

Network Representation

“Neuron” Representation of Logistic Regression

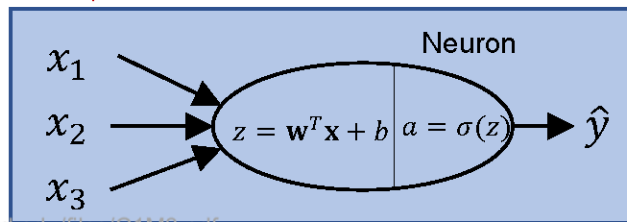
- Logistic Regression

Given \mathbf{x}, y , want parameters \mathbf{w} and b to make \hat{y} infinitely close to y



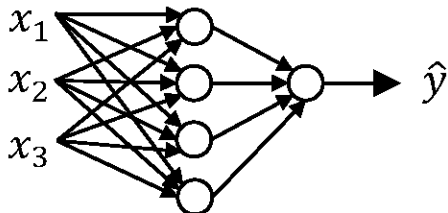
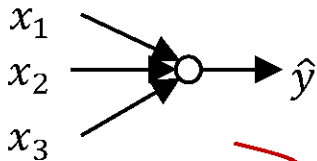
- Where, $\mathbf{x} \in \mathbb{R}^{n_x}$, $0 \leq \hat{y} \leq 1$,
 $\hat{y} = \sigma(\mathbf{w}^T \mathbf{x} + b)$ is the probability of the given \mathbf{x} belongs to the class 1 (binary classification case)

- Parameters: $\mathbf{w} \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$



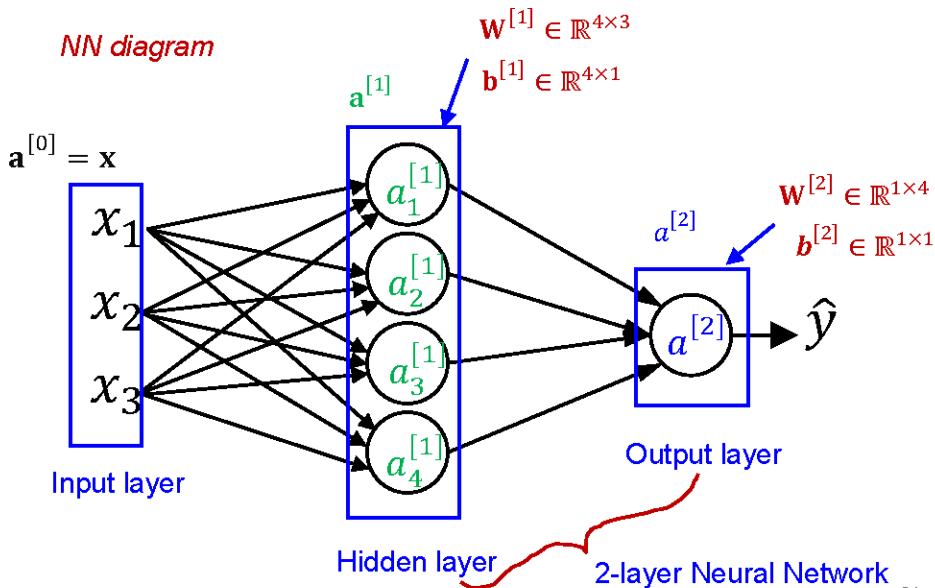
What is a Neural Network?

Logistic Regression



Neural Network

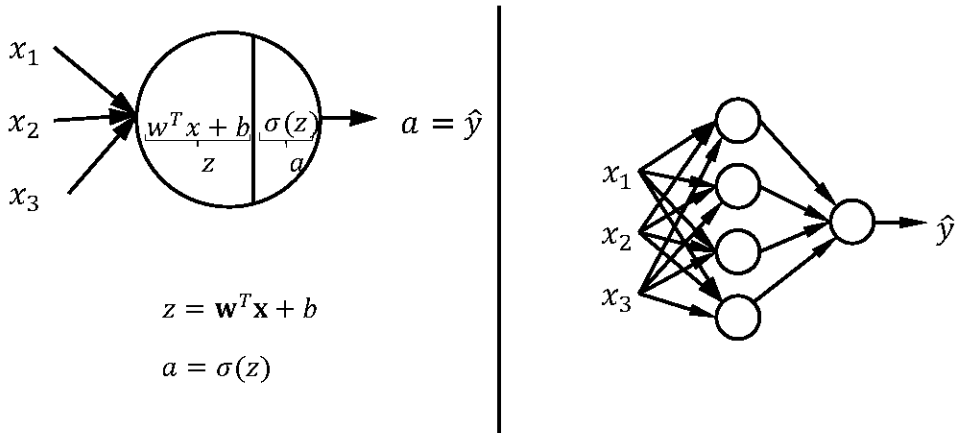
Neural Network (NN) Representation



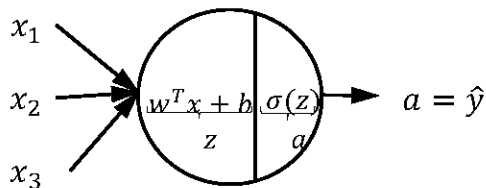
Shallow Neural Networks

Computing Neural Network's Output

Neural Network Representation

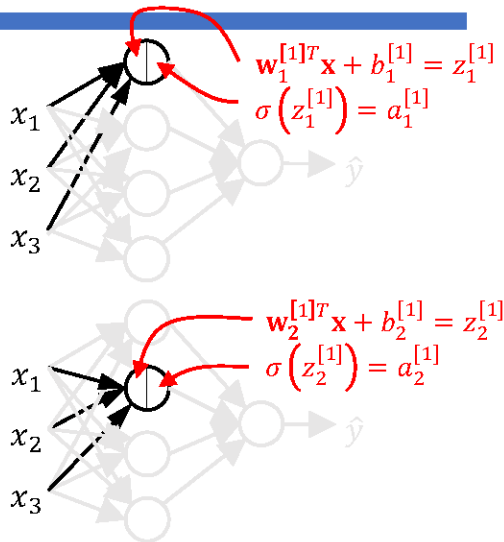


Calculation in Neural Unit

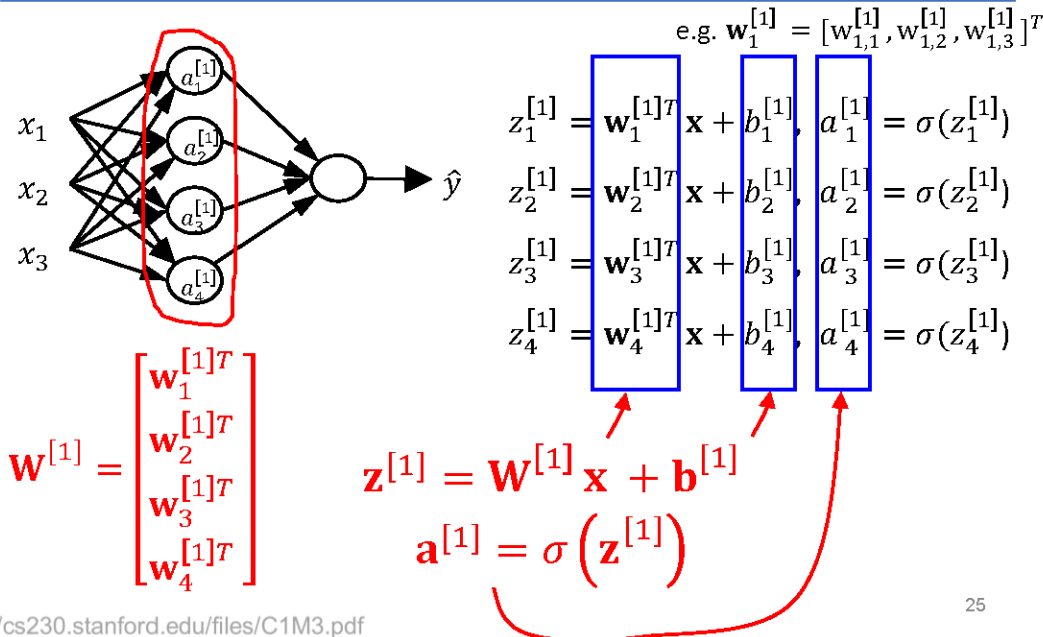


$$z = \mathbf{w}^T \mathbf{x} + b$$

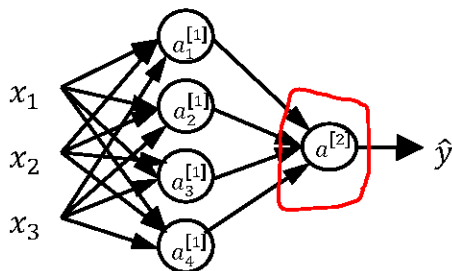
$$a = \sigma(z)$$



Calculation in Neural Network



Calculation in Neural Network



$$z^{[2]} = \mathbf{w}^{[2]} \mathbf{a}^{[1]} + b^{[2]}$$

$$\sigma(\mathbf{w}^{[2]} \mathbf{a}^{[1]} + b^{[2]}) = a^{[2]} = \hat{y}$$

Given input \mathbf{x} :

$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{x} + \mathbf{b}^{[1]}$$

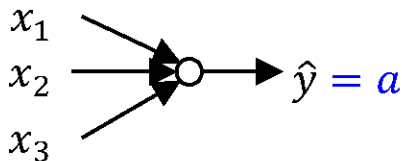
$$\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]})$$

$$z^{[2]} = \mathbf{w}^{[2]} \mathbf{a}^{[1]} + b^{[2]}$$

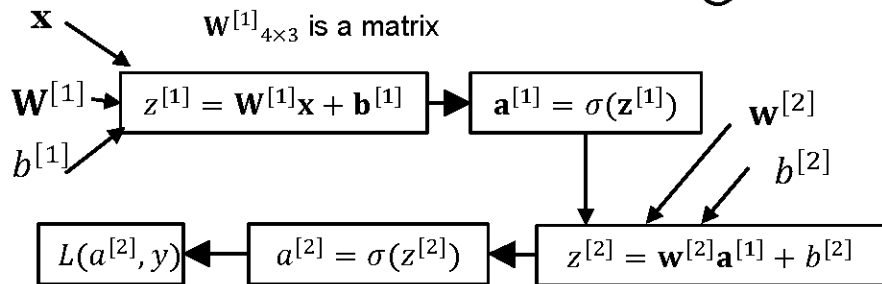
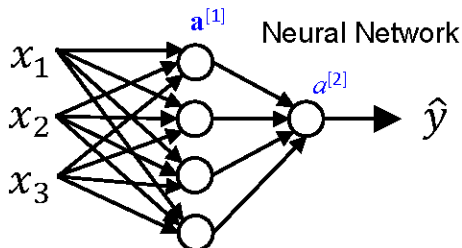
$$a^{[2]} = \sigma(z^{[2]})$$

Neural Network Representation and Output

Logistic Regression



Neural Network



$$a^{[2]} = \hat{y}^{(i)}$$

$$L(\hat{y}^{(i)}, y^{(i)}) = -y^{(i)} \ln(\hat{y}^{(i)}) - (1 - y^{(i)}) \ln(1 - \hat{y}^{(i)})$$

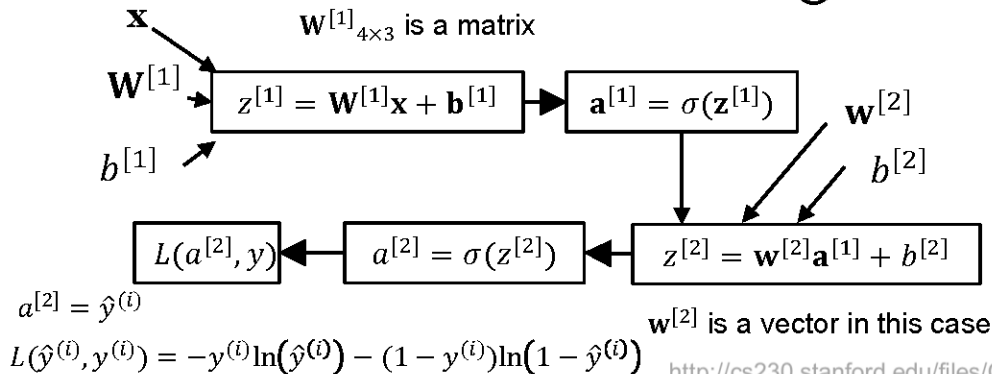
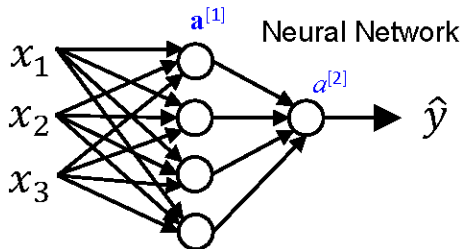
$\mathbf{w}^{[2]}$ is a vector in this case

Shallow Neural Networks

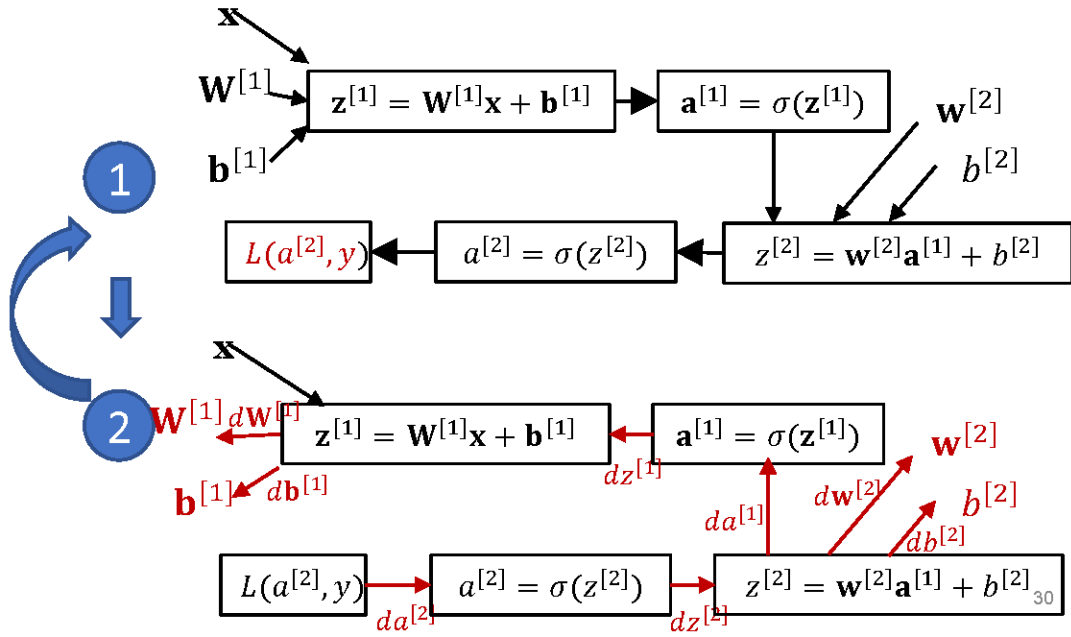
Gradient descent for neural networks

Neural Network

To find parameters \mathbf{W} ($\mathbf{W}^{[1]}$, $\mathbf{w}^{[2]}$...) and \mathbf{b} ($\mathbf{b}^{[1]}$, $b^{[2]}$...) to make $\hat{y}^{(i)}$ infinitely close to $y^{(i)}$.



Training in Neural Network



Gradient Descent for Training Neural Networks

Parameters: $\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \mathbf{w}^{[2]}, b^{[2]}$

Cost function: $J(\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \mathbf{w}^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(a^{[2]}, y)$

Gradient descent:

Repeat

{

Compute predicts ($\hat{y}^{(i)}$)

$$d\mathbf{W}^{[1]} = \frac{\partial J}{\partial \mathbf{W}^{[1]}}, d\mathbf{b}^{[1]} = \frac{\partial J}{\partial \mathbf{b}^{[1]}}$$

$$d\mathbf{W}^{[2]} = \frac{\partial J}{\partial \mathbf{W}^{[2]}}, d\mathbf{b}^{[2]} = \frac{\partial J}{\partial \mathbf{b}^{[2]}}$$

$$\mathbf{W}^{[1]} = \mathbf{W}^{[1]} - \alpha \cdot d\mathbf{W}^{[1]}, \mathbf{b}^{[1]} = \mathbf{b}^{[1]} - \alpha \cdot d\mathbf{b}^{[1]}$$

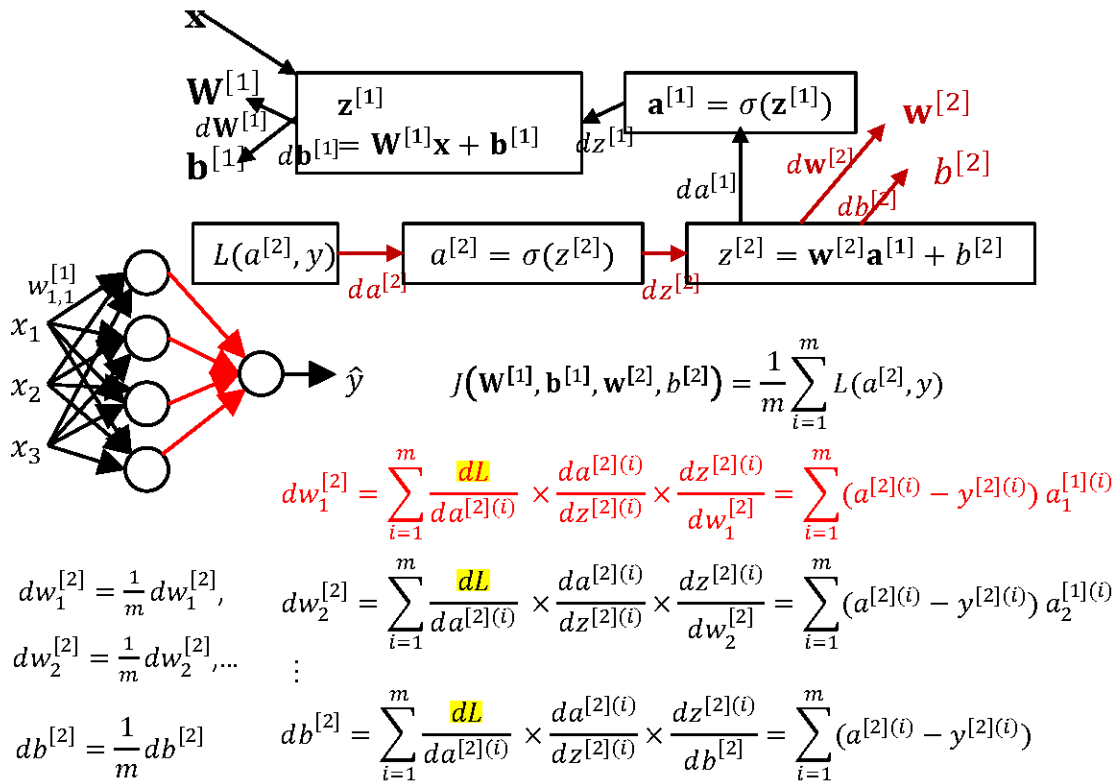
$$\mathbf{W}^{[2]} = \mathbf{W}^{[2]} - \alpha \cdot d\mathbf{W}^{[2]}, \mathbf{b}^{[2]} = \mathbf{b}^{[2]} - \alpha \cdot d\mathbf{b}^{[2]}$$

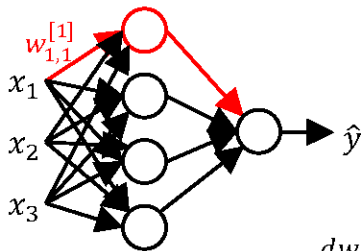
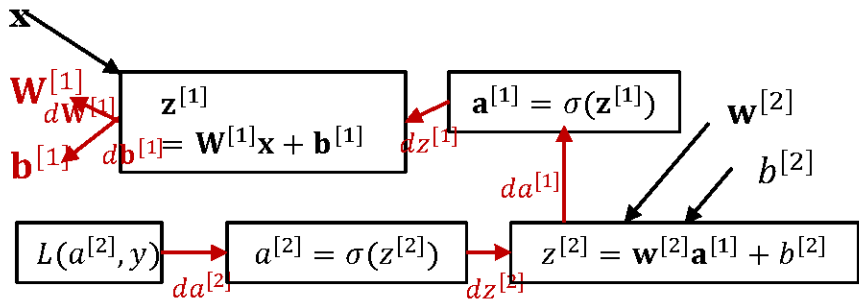
}

Forward Propagation

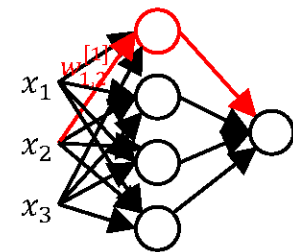
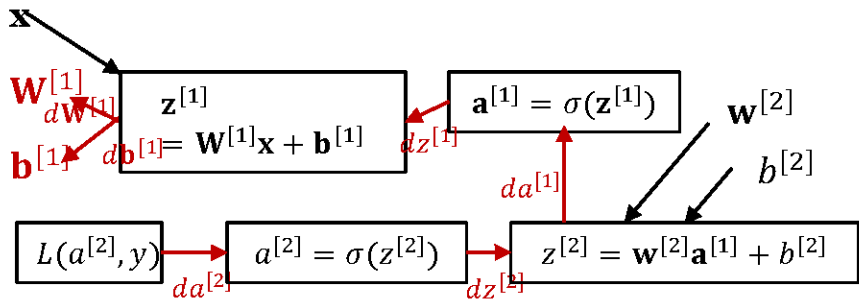
Backpropagation

Updating





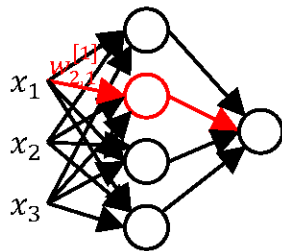
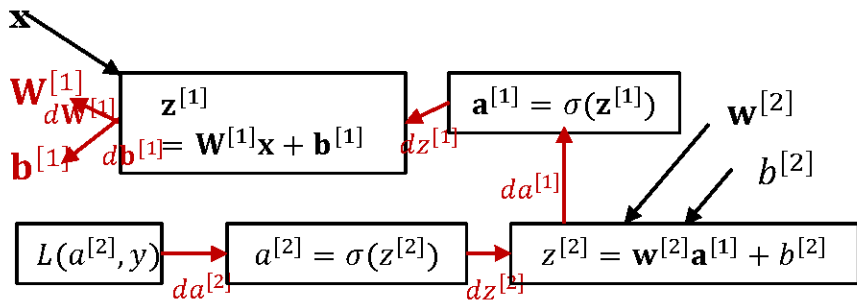
$$\begin{aligned}
 & dw_{1,1}^{[1]} \\
 &= \sum_{i=1}^m \frac{dL}{da^{[2]}(i)} \times \frac{da^{[2]}(i)}{dz^{[2]}(i)} \times \frac{dz^{[2]}(i)}{da_1^{[1]}(i)} \times \frac{da_1^{[1]}(i)}{dz_1^{[1]}(i)} \times \frac{dz_1^{[1]}(i)}{dw_{1,1}^{[1]}} \\
 &= \sum_{i=1}^m (a^{[2]}(i) - y^{[2]}(i)) \times w_1^{[2]}(i) \times a_1^{[1]}(i) (1 - a_1^{[1]}(i)) \times x_1^{(i)} \\
 & dw_{1,1}^{[1]} = \frac{1}{m} dw_{1,1}^{[1]},
 \end{aligned}$$



$$dw_{1,2}^{[1]} = \frac{1}{m} dw_{1,2}^{[1]},$$

$$db_1^{[1]} = \frac{1}{m} db_1^{[1]},$$

$$\begin{aligned} \hat{y} &= \sum_{i=1}^m \frac{dL}{da^{[2](i)}} \times \frac{da^{[2](i)}}{dz^{[2](i)}} \times \frac{dz^{[2](i)}}{da_1^{[1](i)}} \times \frac{da_1^{[1](i)}}{dz_1^{[1](i)}} \times \frac{dz_1^{[1](i)}}{dw_{1,2}^{[1]}} \\ &= \sum_{i=1}^m (a^{[2](i)} - y^{[2](i)}) \times w_1^{[2](i)} \times a_1^{[1](i)} (1 - a_1^{[1](i)}) \times x_2^{(i)} \\ db_1^{[1]} &= \sum_{i=1}^m \frac{dL}{da^{[2](i)}} \times \frac{da^{[2](i)}}{dz^{[2](i)}} \times \frac{dz^{[2](i)}}{da_1^{[1](i)}} \times \frac{da_1^{[1](i)}}{dz_1^{[1](i)}} \times \frac{dz_1^{[1](i)}}{db_1^{[1]}} \\ &= \sum_{i=1}^m (a^{[2](i)} - y^{[2](i)}) \times w_1^{[2](i)} \times a_1^{[1](i)} (1 - a_1^{[1](i)}) \end{aligned}$$



$$\hat{y} = \sum_{i=1}^m \frac{dL}{da^{[2](i)}} \times \frac{da^{[2](i)}}{dz^{[2](i)}} \times \frac{dz^{[2](i)}}{da_2^{[1](i)}} \times \frac{da_2^{[1](i)}}{dz_2^{[1](i)}} \times \frac{dz_2^{[1](i)}}{dw_{2,1}^{[1]}}$$

$$= \sum_{i=1}^m (a^{[2](i)} - y^{[2](i)}) \times w_2^{[2](i)} \times a_2^{[1](i)} (1 - a_2^{[1](i)}) \times x_1^{(i)}$$

$$dw_{2,1}^{[1]} = \frac{1}{m} dw_{2,1}^{[1]},$$

$$db_2^{[1]} = \frac{1}{m} db_2^{[1]},$$

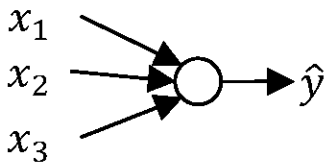
$$db_2^{[1]} = \sum_{i=1}^m \frac{dL}{da^{[2](i)}} \times \frac{da^{[2](i)}}{dz^{[2](i)}} \times \frac{dz^{[2](i)}}{da_2^{[1](i)}} \times \frac{da_2^{[1](i)}}{dz_2^{[1](i)}} \times \frac{dz_2^{[1](i)}}{db_2^{[1]}}$$

$$= \sum_{i=1}^m (a^{[2](i)} - y^{[2](i)}) \times w_2^{[2](i)} \times a_2^{[1](i)} (1 - a_2^{[1](i)})$$

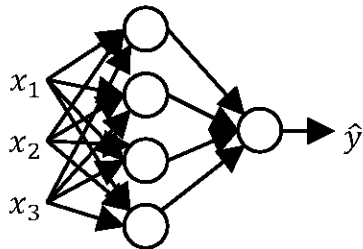
Deep Neural Networks

Deep L-layer Neural network

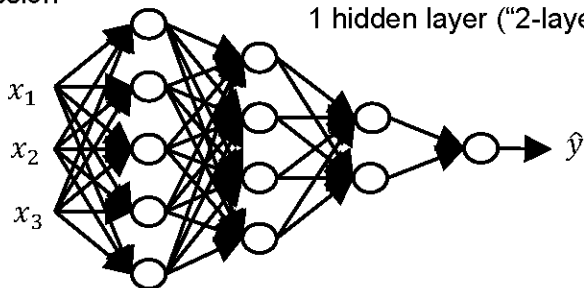
What is a deep neural network?



logistic regression

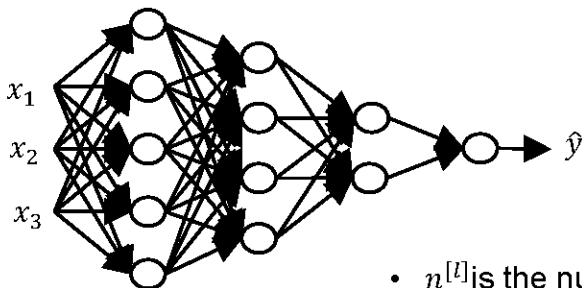


1 hidden layer ("2-layer NN")



3 hidden layers ("4-layer NN")

Deep neural network notation

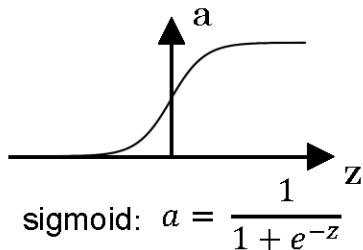
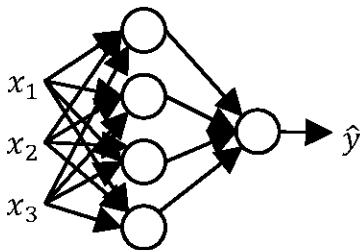


- 3 hidden layers (“4-layer NN”)
- 0-layer: input x_1, x_2, x_3 ($n^{[0]} = 3$)
- 1-layer: 5 neural units ($n^{[1]} = 5$)
- 2-layer: 4 neural units ($n^{[2]} = 4$)
- 3-layer: 2 neural units ($n^{[3]} = 2$)
- 4-layer (output layer): 1 neural units ($n^{[4]} = 1$)
- $n^{[l]}$ is the number of the neural unit in layer l
- $\mathbf{a}^{[l]}$ = activations in layer l
- $\mathbf{a}^{[l]} = g(\mathbf{z}^{[l]})$
- $\mathbf{W}^{[l]}, \mathbf{b}^{[l]}$ are weights for $\mathbf{z}^{[l]}$

Deep Neural Networks

Activation functions

Activation Functions



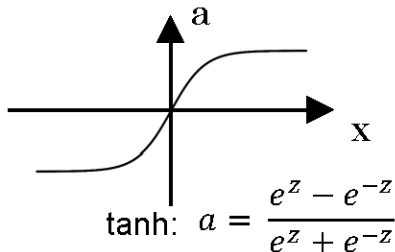
Given input \mathbf{x} :

$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]}\mathbf{x} + \mathbf{b}^{[1]}$$





$$\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]}) \quad g(\mathbf{z}^{[1]})$$

$$\mathbf{z}^{[2]} = \mathbf{w}^{[2]}\mathbf{a}^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(\mathbf{z}^{[2]}) \quad g(\mathbf{z}^{[2]})$$

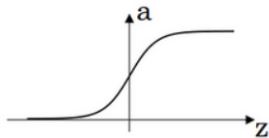


Other activation functions

Name	Plot	Equation	Derivative
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
Rectified Linear Unit (ReLU)		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Parametric Rectified Linear Unit (PReLU) [2]		$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$

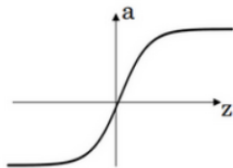
Other activation functions

1) sigmoid activation function



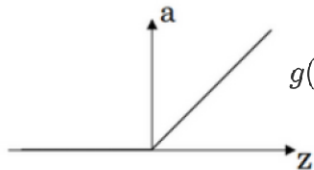
$$g(z) = \frac{1}{1 + e^{-z}}$$

2) Tanh activation function



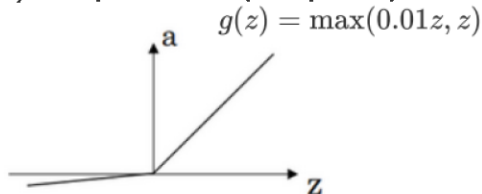
$$g(z) = \tanh(z)$$

3) Rectified Linear Unit (ReLU)



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

4) Leaky linear unit (Leaky ReLU)

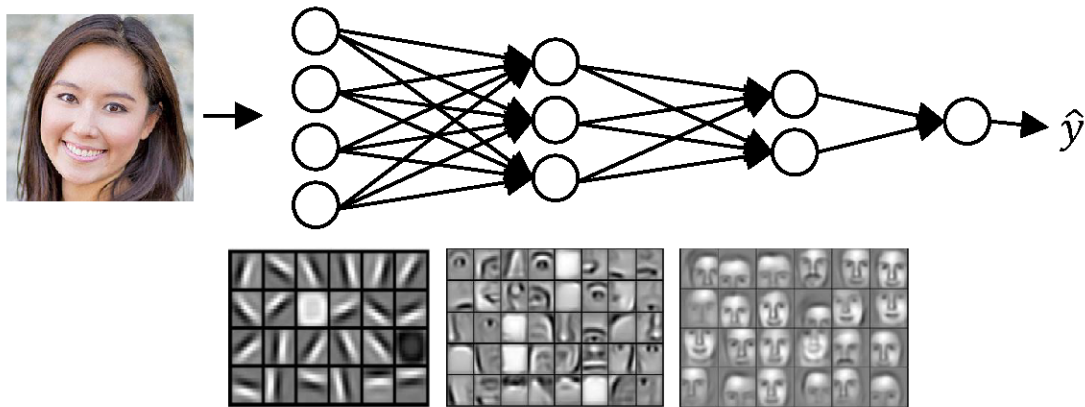


$$g(z) = \max(0.01z, z)$$

Deep Neural Networks

Why deep representations?

Intuition about deep representation



What are hyperparameters?

Parameters: $\mathbf{W}^{[1]}$, $\mathbf{b}^{[1]}$, $\mathbf{W}^{[2]}$, $\mathbf{b}^{[2]}$, $\mathbf{W}^{[3]}$, $\mathbf{b}^{[3]}$...

- HyperParameters:

- Learning rate

- The number of iterations

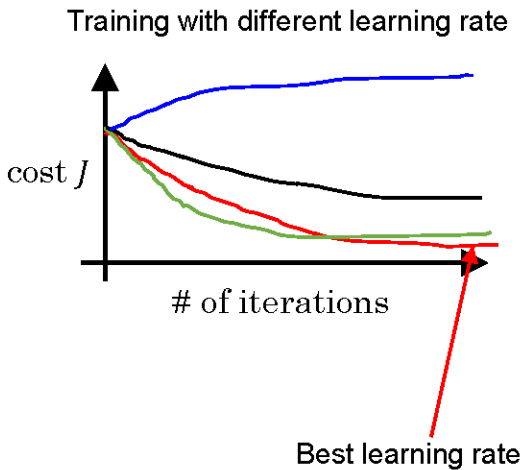
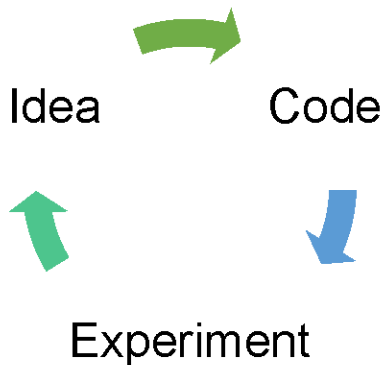
- The number of hidden layers

- The number of neural units in each layer

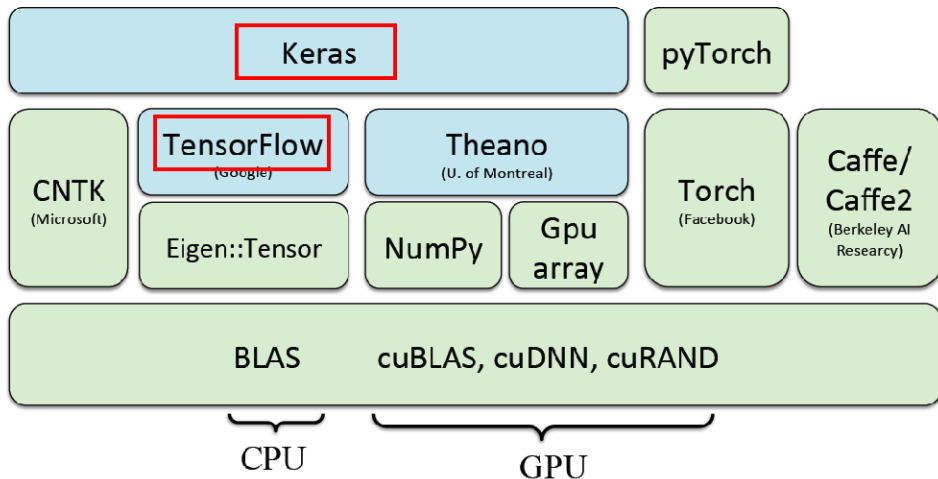
- Choice of activation function

- mini batch size...

Empirical process in training



Deep Learning Software



Data Science

Week 10

Fitting a Model to Data (3)-
Regression