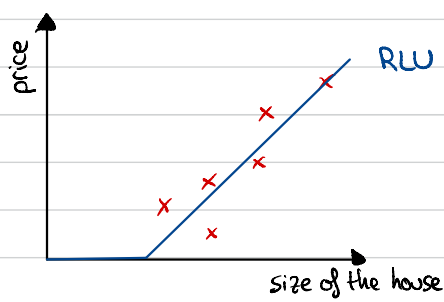


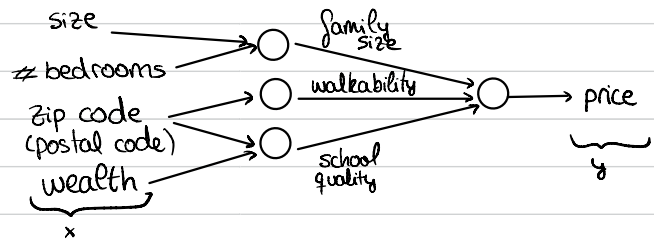
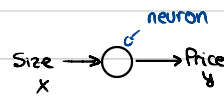
INTRODUCTION TO DEEP LEARNING

What is a Neural Network

Deep learning \rightarrow Training Neural Networks

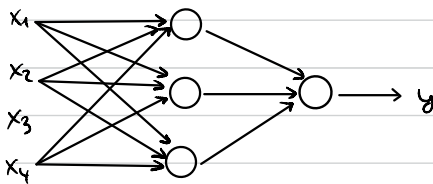


RLU \sim Rectified Linear Unit

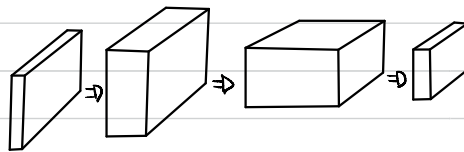


Supervised Learning with Neural Networks

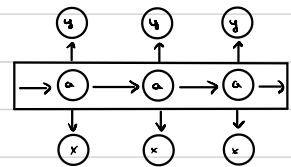
Input (x)	Output (y)	Application	Method
Home features	Price	Real state	standard NN
Ad, user info	Click on ad? (0/1)	Online advertising	CNN \sim Convolution NN
Image	Object (1, ..., 1000)	Photo tagging	RNN \sim Recurrent NN
Audio	Text transcript	Speech recognition	Custom/hybrid NN
English	Chinese	Machine translation	
Image, Radar Info	Position	Autonomous driving	



Standard NN

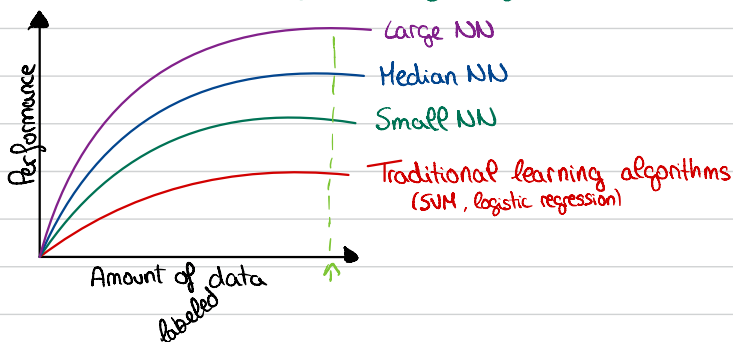


Convolutional NN



Recurrent NN

Why is Deep Learning taking off?



LOGISTIC REGRESSION AS A NEURAL NETWORK

Binary Classification

Output label (y) is 1 or 0

Notation

$$(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$$

m training examples: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$m = m_{\text{train}}$

$m_{\text{test}} = \#$ test examples

$$X = \begin{bmatrix} 1 & & & \\ x^{(1)} & & & \\ & x^{(2)} & & \\ & & \dots & \\ & & & x^{(m)} \\ & & & & 1 \end{bmatrix} \begin{matrix} \uparrow \\ n_x \\ \downarrow \\ m \end{matrix}$$

$$X \in \mathbb{R}^{n_x \times m}$$

\hookrightarrow Python \rightarrow `X.shape(n_x, m)`

$$y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

$$y \in \mathbb{R}^{1 \times m}$$

\hookrightarrow Python \rightarrow `y.shape(1, m)`

Logistic Regression

Given x , want $\hat{y} = P(y=1|x) \rightarrow 0 \leq \hat{y} \leq 1$

$x \in \mathbb{R}^{n_x} \rightarrow$ Parameters: $w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$

Output: $\hat{y} = \sigma(\underbrace{w^T x + b}_z)$
 \uparrow sigmoid function

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

If z large $\sigma(z) \approx \frac{1}{1+0} = 1$

If z large negative $\sigma(z) \approx \frac{1}{1+\infty} = 0$

Logistic Regression cost function

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

If $y=1: L(\hat{y}, y) = -\log \hat{y} \rightarrow$ want \hat{y} large
If $y=0: L(\hat{y}, y) = -\log(1-\hat{y}) \rightarrow$ want \hat{y} small

Cost function

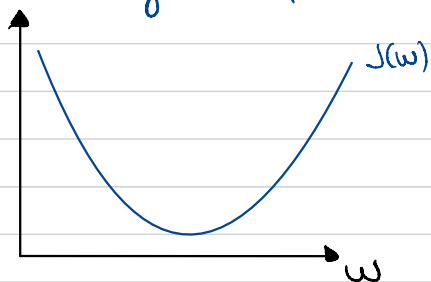
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

Loss function \rightarrow Error for a single training example

Cost function \rightarrow Average of the loss functions of the entire training set

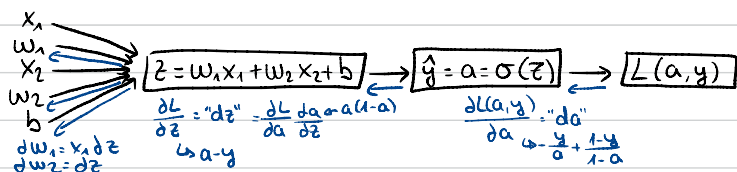
Gradient Descent

Want to find w, b that minimize $J(w, b)$



Repeat $\{$ $w := w - \alpha \frac{\partial J(w, b)}{\partial w}$ $\}$ ^{learning rate} $dw \sim \text{slope}$
 $b := b - \alpha \frac{\partial J(w, b)}{\partial b} db$

Logistic Regression Gradient Descent



Logistic regression \rightarrow modify parameters (w and b) in order to minimize the loss

Gradient Descent on m Examples

Initialize $\rightarrow J=0, dw_1=0, dw_2=0, db=0$

for $i=1$ to m $\{$

$$z^{(i)} = w^T x^{(i)} + b$$

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

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

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

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

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

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

$\}$

$J/=m$

$dw_1/=m; dw_2/=m; db/=m$

Update:

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

PYTHON AND VECTORIZATION

Vectorization

$$z = w^T X + b \rightarrow \text{Python} \rightarrow z = \text{np.dot}(w, x)$$

Vectorizing Logistic Regression

$$z^{(1)} = w^T X^{(1)} + b$$

$$z^{(2)} = w^T X^{(2)} + b$$

$$z^{(3)} = w^T X^{(3)} + b$$

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

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

$$a^{(3)} = \sigma(z^{(3)})$$

$$X = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & \dots & 1 \end{bmatrix}_{n \times m}$$

$$Z = [z^{(1)}, z^{(2)}, \dots, z^{(m)}] = w^T X + [b, b, \dots, b] = [w^T x^{(1)} + b, w^T x^{(2)} + b, \dots, w^T x^{(m)} + b]$$

$$\rightarrow \text{python} \rightarrow Z = \text{np.dot}(w.T, X) + b \quad \text{"broadcasting"}$$

$$A = [a^{(1)}, a^{(2)}, \dots, a^{(m)}] = \sigma(Z)$$

Vectorizing Logistic Regression's Gradient Output

$$dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)}$$

$$dZ = [dz^{(1)}, dz^{(2)}, \dots, dz^{(m)}]_{1, m}$$

$$dZ = A - Y$$

$$A = [a^{(1)}, a^{(2)}, \dots, a^{(m)}] \quad Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)} \rightarrow \text{Python} \rightarrow \frac{1}{m} \text{np.sum}(dZ)$$

$$dw = \frac{1}{m} X dZ^T = \frac{1}{m} \begin{bmatrix} x^{(1)} & \dots & x^{(m)} \\ 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix} = \frac{1}{m} [x^{(1)} dz^{(1)} \dots x^{(m)} dz^{(m)}]$$

$$w := w - \alpha dw$$

$$b := b - \alpha db$$

Broadcasting in Python

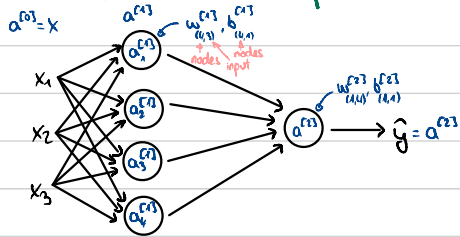
$$ab + ac - b - c$$

$$\begin{array}{l} (m, n) + (1, n) \sim (m, n) \\ (m, n) * (m, 1) \sim (m, n) \\ (m, n) / (1, n) \sim (m, n) \\ (m, n) // (1, n) \sim (m, n) \end{array}$$

for iter in range(1000):

SHALLOW NEURAL NETWORKS

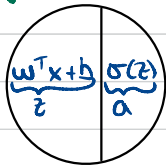
Neural Networks Representation



2 layer NN
 Layer 0 - input layer

input layer hidden layer output layer

Computing a Neural Network's Output



z ← layer
 a_i ← node in layer

Vectorizing Across Multiple Examples

$$Z^{(1)} = W^{(1)} X + b^{(1)}, A^{(1)} = \sigma(Z^{(1)})$$

$$Z^{(2)} = W^{(2)} X + b^{(2)}, A^{(2)} = \sigma(Z^{(2)})$$

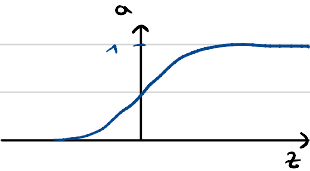
training examples

$$X = \begin{bmatrix} | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix} \quad Z = \begin{bmatrix} | & | & \dots & | \\ z^{(1)(1)} & z^{(1)(2)} & \dots & z^{(1)(m)} \\ | & | & \dots & | \end{bmatrix} \quad A^{(1)} = \begin{bmatrix} | & | & \dots & | \\ a^{(1)(1)} & a^{(1)(2)} & \dots & a^{(1)(m)} \\ | & | & \dots & | \end{bmatrix} \quad \text{hidden units}$$

Activation functions

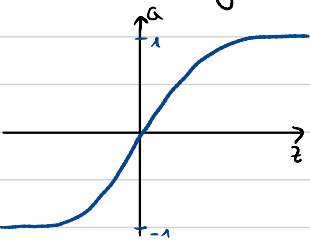
Can be different for different layers

Sigmoid



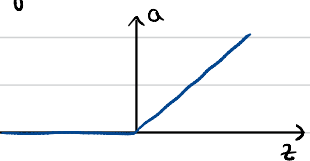
$$a = \frac{1}{1 + e^{-z}}$$

Hyperbolic tangent



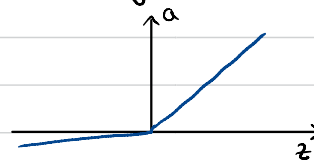
$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Rectified Linear Unit (ReLU)



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

Leaky ReLU



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

Derivatives of Activation Functions

Sigmoid

$$g(z) = \frac{1}{1+e^{-z}} \rightarrow g'(z) = \frac{d}{dz} g(z) = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}}\right) = g(z)(1-g(z)) = a(1-a)$$

Hyperbolic tangent

$$g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \rightarrow g'(z) = \frac{d}{dz} g(z) = 1 - (\tanh(z))^2 = 1 - g(z)^2 = 1 - a^2$$

ReLU

$$g(z) = \max(0, z) \rightarrow g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Leaky ReLU

$$g(z) = \max(0.01z, z) \rightarrow g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Gradient Descent for Neural Networks (2 layers)

Parameters:

$$W^{[1]} \quad b^{[1]} \quad W^{[2]} \quad b^{[2]}$$

$(n^{[1]}, n^{[2]}) \quad (n^{[1]}, 1) \quad (n^{[2]}, n^{[2]}) \quad (n^{[2]}, 1)$

$$n_x = n^{[0]}, \quad n^{[1]}, \quad n^{[2]} = 1$$

Cost function (binary classification)

$$J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}_i, y_i) \alpha^{[2]}$$

Gradient descent

Repeat ξ

Compute predicts $(\hat{y}^{(i)}, i=1, \dots, m)$

$$dW^{[1]} = \frac{\partial J}{\partial W^{[1]}}, \quad db^{[1]} = \frac{\partial J}{\partial b^{[1]}} \dots$$

$$W^{[1]} = W^{[1]} - \alpha dW^{[1]}$$

$$b^{[1]} = b^{[1]} - \alpha db^{[1]}$$

$$W^{[2]} = W^{[2]} - \alpha dW^{[2]}$$

$$b^{[2]} = b^{[2]} - \alpha db^{[2]}$$

}

Forward propagation

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

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

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

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

Backward propagation

$$dz^{[2]} = A^{[2]} - y \quad \alpha \text{ me fonction cou } y - A^{[2]}$$

$$dW^{[2]} = 1/m dz^{[2]} A^{[1]T}$$

$$db^{[2]} = 1/m \text{ np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims}=\text{True})$$

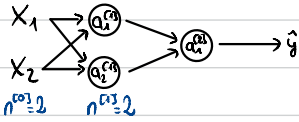
$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

\uparrow element-wise product

$$dW^{[1]} = 1/m dz^{[1]} X^T$$

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

Random Initialization



Zero initialization \rightarrow Symmetry \rightarrow After every iteration, the two hidden units are still computing exactly the same function.

$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$\hookrightarrow a_1 = a_2$ $\hookrightarrow dz_1 = dz_2$

Random initialization \rightarrow Avoids the symmetry

$$W^{[1]} = \text{np.random.randn}((2,2)) * 0.01$$

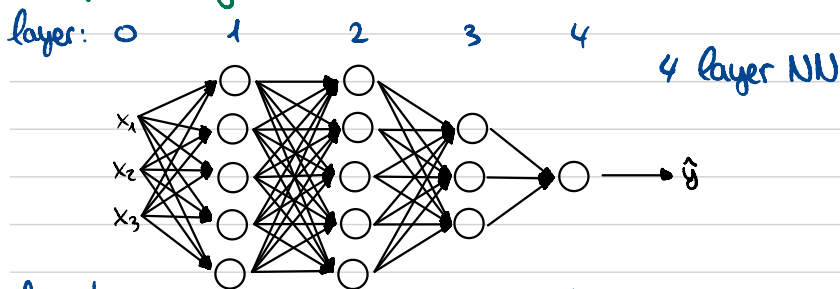
$$b^{[1]} = \text{np.zeros}((2,1))$$

$$W^{[2]} = \text{np.random.randn}((1,2)) * 0.01$$

$$b^{[2]} = 0$$

DEEP NEURAL NETWORK

Deep L-layer Neural Network



elements: 5 5 3 1

$L=4 \sim \# \text{ layers}$

$n^{[l]} = \# \text{ nodes/units in layer } l$

$a^{[l]} = \text{activations in layer } l$

$a^{[l]} = g^{[l]}(z^{[l]})$

$W^{[l]} = \text{weights for } z^{[l]}$

$b^{[l]}$

$x = a^{[0]}$

$y = a^{[L]}$

$$n^{[1]} = 5 \quad n^{[3]} = 3$$

$$n^{[2]} = 5 \quad n^{[4]} = n^{[L]} = 1$$

$$n^{[0]} = n_x = 3$$

shallow NN
deep NN

Forward propagation in a Deep Network

$$z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]} \quad \text{for } l=1:L$$

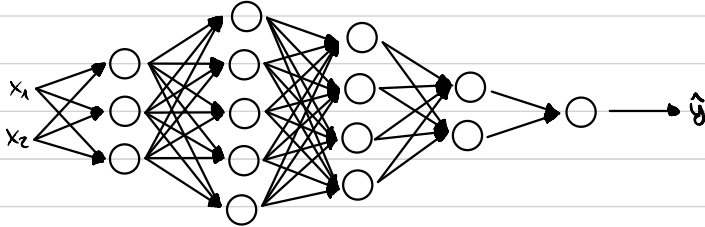
$$A^{[l]} = g^{[l]}(z^{[l]})$$

All are column vectors stacked horizontally

Getting your Matrix Dimensions Right

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

layer: 0 1 2 3 4 5



elements: 3 5 4 2 1

$$z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

$(n^{[l]}, m)$ $(n^{[l]}, n^{[l-1]})$ $(n^{[l-1]}, m)$ $(n^{[l]}, 1)$ ← broadcasting

$$dW^{[l]} \quad (n^{[l]}, n^{[l-1]})$$

$$db^{[l]} \quad (n^{[l]}, 1)$$

Why Deep representations

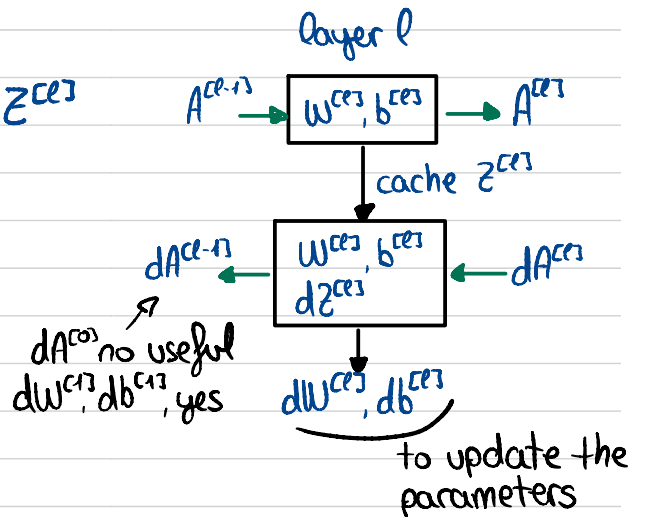
Circuit theory and deep learning (informally) → There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

Building blocks of Deep Neural networks

Layer l : $W^{[l]}, b^{[l]}$

Forward: Input $A^{[l-1]}$, Output $A^{[l]}$
 $z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$ → cache the value of $z^{[l]}$
 $A^{[l]} = g^{[l]}(z^{[l]})$

Backward: Input $dA^{[l]}$, Output $dA^{[l-1]}$
 (also) $dW^{[l]}, db^{[l]}$



Parameters vs Hyperparameters

Parameters: w^{ces} , b^{ces}

Hyperparameters: α , # iterations, # hidden layers (L), # hidden units (n^{ces}),
choice of activation function (sigmoid, ReLU, tanh)

To choose the correct α :

