

Deep Learning Specialization

Neural Networks and Deep Learning

Learning Notes

Du Ang
du2ang233@gmail.com

September 29, 2017

Contents

1	Welcome	3
2	Introduction to Deep Learning	3
2.1	What is a neural network?	3
2.1.1	Single neural network	3
2.1.2	Multiple neural network	4
2.2	Supervised Learning	4
2.2.1	Structured vs unstructured data	4
2.3	Why is Deep Learning taking off?	5
3	Neural Network Basics	6
3.1	Logistic Regression as a Neural Network	6
3.1.1	Binary Classification	6
3.1.2	Notation	6
3.1.3	Logistic Regression	7
3.2	Python and Vectorization	10
3.2.1	Vectorization	10
3.2.2	More Vectorization Examples	10
3.2.3	Vectorizing Logistic Regression	11
3.2.4	Broadcasting in Python	13
3.2.5	A Note on Python/NumPy Vectors	13
3.2.6	Explanation of Logistic Regression Cost Function (Optional)	14
4	Shallow Neural Network	14
4.1	Neural Network Overview	14
4.2	Neural Network Representation	15
4.3	Computing a Neural Network's Output	15
4.4	Vectorizing across Multiple Examples	16
4.5	Explanation for Vectorized Implementation	17
4.6	Activation Functions	18
4.6.1	Pros and cons of activation functions	18
4.7	Why do you need non-linear activation functions?	19
4.8	Derivatives of activation functions	19
4.8.1	Sigmoid activation function	19
4.8.2	Tanh activation function	19
4.9	ReLU	19
4.10	Leaky ReLU	19
4.11	Gradient descent for Neural Networks	20

4.12	Backpropagation intuition (Optional)	20
4.12.1	Summary of gradient descent	21
4.13	Random Initialization	21
5	Deep Neural Networks	22
5.1	Deep L -layer neural network	22
5.2	Forward Propagation in a Deep Network	22
5.3	Getting your matrix dimensions right	24
5.4	Why deep representations?	25
5.4.1	Intuition about deep representation	25
5.4.2	Circuit theory and deep learning	26
5.5	Building blocks of deep neural networks	26
5.6	Parameters vs Hyperparameters	26
5.6.1	What are hyperparameters?	26
5.6.2	Applied deep learning is a very empirical process	27

1 Welcome

AI is the new Electricity. — Andrew Ng

Electricity had once transformed countless industries: transportation, manufacturing, healthcare, communications, and more.

AI will now bring about an equally big transformation.

Courses in this sequence (Specialization):

1. Neural Networks and Deep Learning
2. Improving Deep Neural Networks: Hyperparameter tuning, Regularization and Optimization
3. Structuring Machine Learning Projects
4. Convolutional Neural Networks
5. Sequence Models

2 Introduction to Deep Learning

2.1 What is a neural network?

It is powerful learning algorithm inspired by how the brain works.

2.1.1 Single neural network

Given data about the size of houses on the real estate market and you want to fit a function that will predict their price. It is a linear regression problem because the price as a function of size is a continuous output.

We know the prices can never be negative so we are creating a function called Rectified Linear Unit (ReLU) which starts at zero.

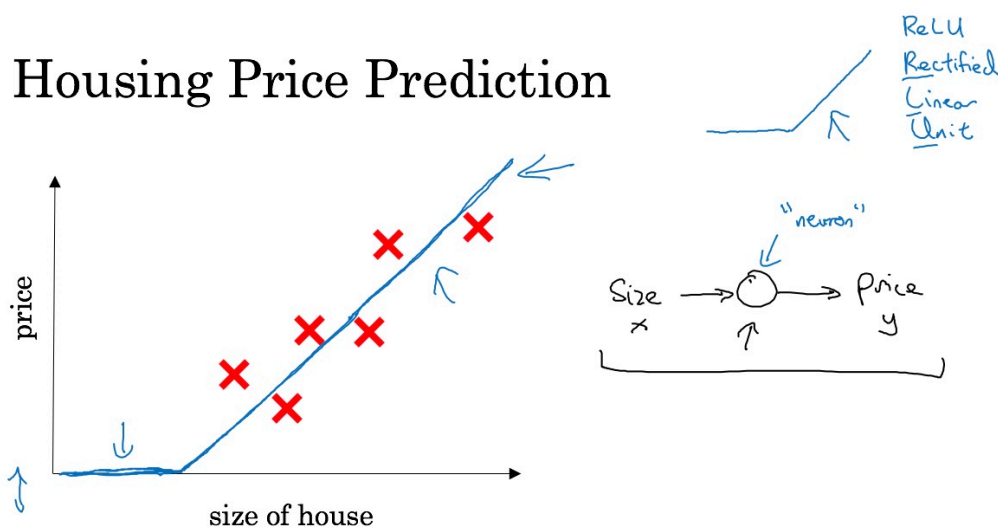


Figure 1: The “Housing Price Prediction” problem. The input is the size of the house (x); The output is the price (y); The “neuron” implements the function ReLU.

2.1.2 Multiple neural network

The price of a house can be affected by other features such as size, number of bedrooms, zip code and wealth. The role of the neural network is to predicted the price and it will automatically generate the hidden units. We only need to give the inputs x and the output y .

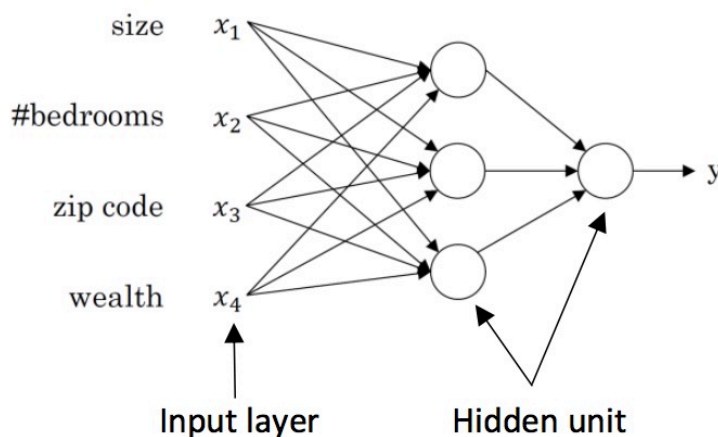


Figure 2: Multiple neural network

2.2 Supervised Learning

In supervised learning, we are given a data set and already know about what our correct output should look like, having the idea that there is a relationship between the input and the output.

Supervised learning problems are categorized into “regression” and “classification” problems. In a regression problem, we are trying to predict results with a continuous output, meaning that we are trying to map input variables to some continuous function. In a classification problem, we are instead trying to predict results in a discrete output. In other words, we are trying to map input variables into discrete categories.

Here are some example of supervised learning:

Input(x)	Output(y)	Application	NN Types
Home features	Price	Real Estate	Standard NN
Ad, user info	Click on ad? (0/1)	Online Advertising	Standard NN
Image	Object (1, ..., 1000)	Photo tagging	CNN
Audio	Text transcript	Speech Recognition	RNN
English	Chinese	Machine translation	RNN
Image, Radar info	Position of other cars	Autonomous driving	Custom/Hybrid

There are different types of neural network, for example, Convolutional Neural Network (CNN) used often for image application and Recurrent Neural Network (RNN) used for one-dimensional sequence data such as translating English to Chinese or a temporal component such as text transcript. As for the autonomous driving, it is a hybrid neural network architecture.

2.2.1 Structured vs unstructured data

Structured data refers to things that has a defined meaning such as price, age whereas unstructured data refers to thing like pixel, raw audio, text.

2.3 Why is Deep Learning taking off?

Deep learning is taking off due to a large amount of data available through the digitization of the society, faster computation and innovation in the development of neural network algorithm.

Scale drives deep learning progress

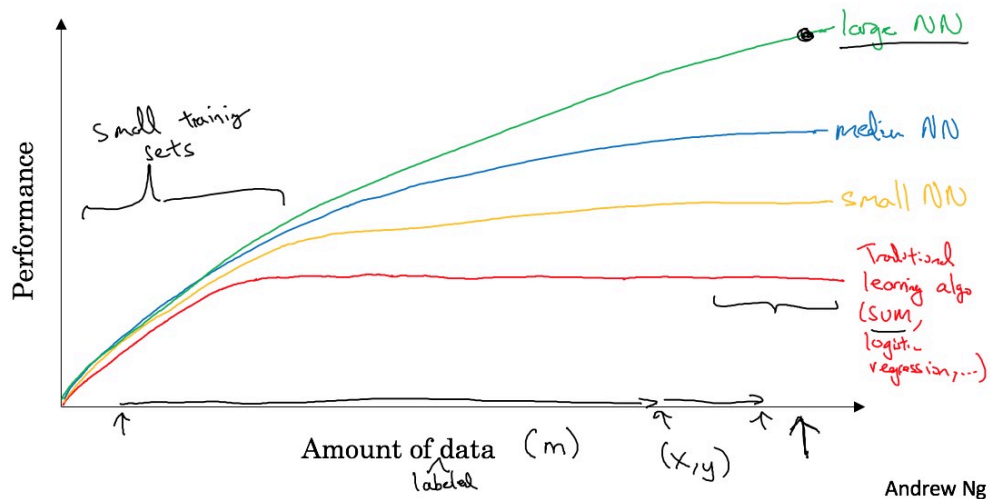
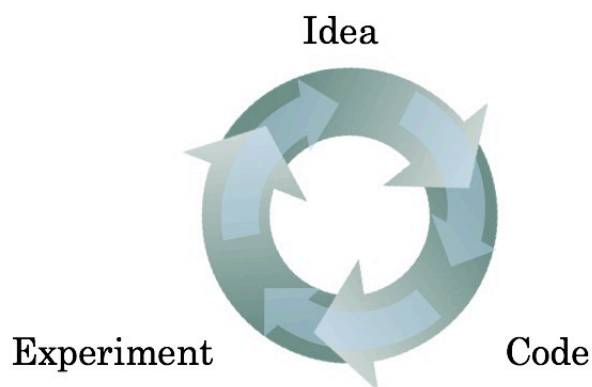


Figure 3: Scale drives deep learning progress

Two things have to be considered to get to the high level of performance:

1. Being able to train a big enough neural network
2. Huge amount of labelled data

The process of training a neural network is iterative:



It could take a good amount of time to train a neural network, which affects your productivity. Faster computation helps to iterate and improve new algorithm.

3 Neural Network Basics

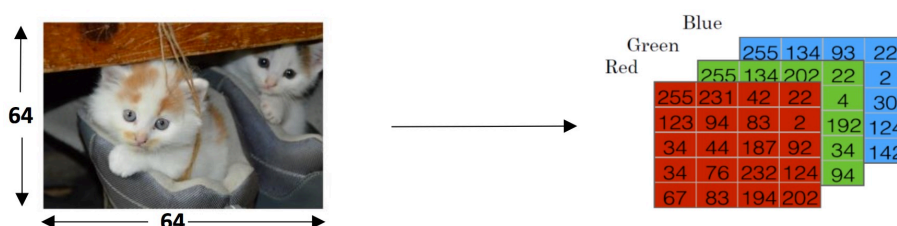
3.1 Logistic Regression as a Neural Network

3.1.1 Binary Classification

In a binary classification problem, the result is a discrete value output. For example:

- account hacked (1) or compromised (0)
- a tumor malign (1) or benign (0)

Example: Cat vs Non-Cat The goal is to train a classifier that the input is an image represented by a feature vector, x , and predicts whether the corresponding label y is 1 or 0. In this case, whether this is a cat image (1) or a non-cat image (0).



An image is stored in the computer in three separate matrices corresponding to the Red, Green, and Blue color channels of the image. The three matrices have the same size as the image, for example, the resolution of the cat image is 64 pixels \times 64 pixels, the three matrices (RGB) are 64 \times 64 each.

The value in a cell represents the pixel intensity which will be used to create a feature vector of n -dimension. In pattern recognition and machine learning, a feature vector represents an object, in this case, a cat or no cat.

To create a feature vector, x , the pixel intensity values will be “unroll” or “reshape” for each color. The dimension of the input feature vector x , is $n_x = 64 \times 64 \times 3 = 12288$.

$$x = \begin{bmatrix} 255 \\ 231 \\ 42 \\ \vdots \\ 255 \\ 134 \\ 202 \\ \vdots \\ 255 \\ 134 \\ 93 \\ \vdots \end{bmatrix} \begin{matrix} \text{red} \\ \text{green} \\ \text{blue} \end{matrix}$$

3.1.2 Notation

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

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

$m = m_{train}$ m_{test} = the number of test examples

$$\mathbf{X} = \begin{bmatrix} | & | & & | \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \\ | & | & & | \end{bmatrix} \quad \mathbf{X} \in \mathbb{R}^{n_x \times m}$$

$$\mathbf{Y} = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix} \quad \mathbf{Y} \in \mathbb{R}^{1 \times m}$$

In Python/NumPy, $\mathbf{X}.\text{shape} = (\mathbf{n_x}, \mathbf{m})$, $\mathbf{Y}.\text{shape} = (1, \mathbf{m})$.

3.1.3 Logistic Regression

Logistic regression is a learning algorithm used in a supervised learning problem when the output y are all either zero or one. The goal of logistic regression is to minimize the error between its predictions and training data.

Example: Cat vs Non-Cat Given an image represented by a feature vector \mathbf{x} , the algorithm will evaluate the probability of a cat being in that image.

Given \mathbf{x} , want $\hat{y} = P(y = 1|\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^{n_x}$, $0 \leq \hat{y} \leq 1$

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

Output: $\hat{y} = \mathbf{w}^T \mathbf{x} + b$ $\hat{y} = \sigma(\mathbf{w}^T \mathbf{x} + b)$, where $\sigma(z) = \frac{1}{1+e^{-z}}$.

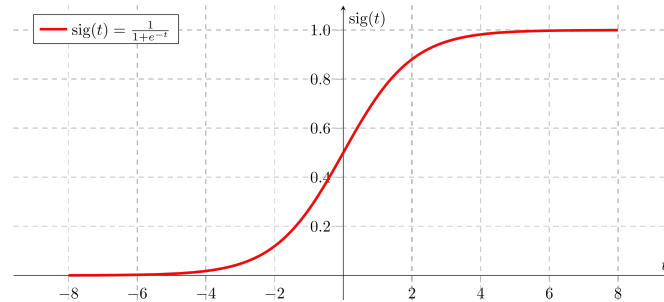


Figure 4: Sigmoid

$(\mathbf{w}^T \mathbf{x} + b)$ is a linear function ($a\mathbf{x} + b$), but since we are looking for a probability constraint between $[0, 1]$, the sigmoid function is used. The function is bounded between $[0, 1]$ as shown in the graph 4.

Some observations from the sigmoid function graph:

- If z is a large positive number, then $\sigma(z) = 1$
- If z is small or large negative number, then $\sigma(z) = 0$
- If $z = 0$, then $\sigma(z) = 0.5$

Logistic Regression cost function $\hat{y}^{(i)} = \sigma(\mathbf{w}^T \mathbf{x}^{(i)} + b)$, where $\sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}}$

Given $\{(\mathbf{x}^{(1)}, y^{(2)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$.

Loss (error) function: (for single example)

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

If $y = 1$: $\mathcal{L}(\hat{y}, y) = -\log \hat{y} \leftarrow$ want $\log \hat{y}$ large, want \hat{y} large.

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

The loss function measures the discrepancy between the prediction ($\hat{y}^{(i)}$) and the desired output ($y^{(i)}$). In other words, the loss function computes the error for a single training example.

Note: Why can't use $\mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$? Because $\mathcal{L}(\hat{y}, y) = -(y \log \hat{y} + (1 - y) \log(1 - \hat{y}))$ is convex, however, $\mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$ is non-convex.

Cost function: (for entire training set)

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

The cost function is the average of the loss function of the entire training set. We are going to find the parameters \mathbf{w} and b that minimize the overall cost function.

Gradient Descent

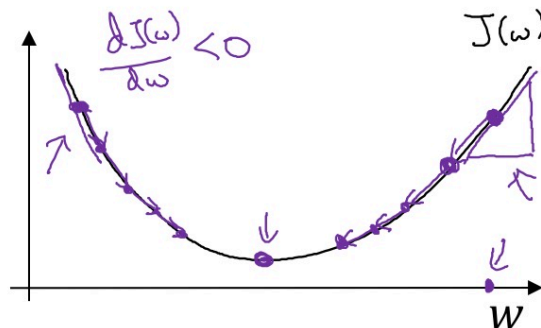


Figure 5: Gradient descent. (Ignore b here)

```
repeat
     $\mathbf{w} := \mathbf{w} - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial \mathbf{w}}$ 
     $b := b - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial b}$ 
until iteration = MaxIteration;
```

Gradient Descent Note: We often use “ $d\mathbf{w}$ ” to denote $\frac{\partial J(\mathbf{w}, b)}{\partial \mathbf{w}}$, use “ db ” to denote $\frac{\partial J(\mathbf{w}, b)}{\partial b}$. In other words, we use “ $dvar$ ” to denote $\frac{dFinalOutputVar}{dvar}$ or $\frac{\partial FinalOutputVar}{\partial var}$.

Gradient on single example:

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

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

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

According to the computation graph 6, go backwards to compute the derivatives:

$$da = \frac{d\mathcal{L}(a, y)}{da} = -\frac{y}{a} + \frac{1 - y}{1 - a}$$

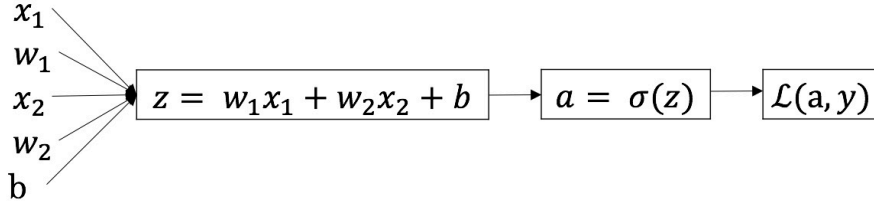


Figure 6: The computation graph of logistic regression gradient descent

$$dz = \frac{d\mathcal{L}(a, y)}{dz} = \frac{d\mathcal{L}}{da} \cdot \frac{da}{dz} = \left(-\frac{y}{a} + \frac{1-y}{1-a}\right) \cdot a(1-a) = a - y$$

$$dw_1 = \frac{\partial \mathcal{L}}{\partial w_1} = x_1 dz = x_1(a - y)$$

$$dw_2 = \frac{\partial \mathcal{L}}{\partial w_2} = x_2 dz = x_2(a - y)$$

$$db = \frac{\partial \mathcal{L}}{\partial b} = dz = a - y$$

Gradient on m examples:

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

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

$$\frac{\partial J(\mathbf{w}, b)}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m \frac{\partial \mathcal{L}(a^{(i)}, y^{(i)})}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m dw_1^{(i)}$$

The Algorithm 1 is the algorithm of logistic regression gradient descent on m examples. It uses two for-loops, so it's less efficient than use vectorization.

```

// Initialization
J = 0;
for j = 1 to n do
  | dw_j = 0
end
db = 0;
// Compute cost and derivatives
for i = 0 to m do
  | z^(i) = w^T x^(i) + b;
  | a^(i) = σ(z^(i));
  | J = J - [y^(i) log a^(i) + (1 - y^(i)) log(1 - a^(i))];
  | dz^(i) = a^(i) - y^(i);
  | for j = 1 to n do
  | | dw_j = dw_j + x_j^(i) dz^(i);
  | end
  | db = db + dz^(i)
end
// Get the average
J = J/m;
for j = 1 to n do
  | dw_j = dw_j/m
end
db = db/m;
// Gradient descent
for j = 1 to n do
  | dw_j = dw_j - α dw_j
end
b = b - α db

```

Algorithm 1: Logistic regression gradient descent on m examples

3.2 Python and Vectorization

3.2.1 Vectorization

What is vectorization?

$$z = \mathbf{w}^T \mathbf{x} + b \quad \mathbf{w} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}, \quad \mathbf{w} \in \mathbb{R}^{n_x}, \quad \mathbf{x} \in \mathbb{R}^{n_x}$$

Non vectorized:

```

z = 0
for i in range(n_x):
    z += w[i] * x[i]
z += b

```

Vectorized:

```

z = np.dot(w, x) + b

```

$$\begin{cases} \text{CPU} \\ \text{GPU} \end{cases} \quad \text{SIMD - Single Instruction Multiple Data}$$

3.2.2 More Vectorization Examples

Neural network programming guideline Whenever possible, avoid explicit for-loops.

For example, compute $\mathbf{u} = \mathbf{A}\mathbf{v}$. $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{v} \in \mathbb{R}^n$.

Non-vectorized: $u_i = \sum_j A_{ij} v_j$

```
u = np.zeros((n, 1))
for i in range(A.shape[0]):
    for j in range(A.shape[1]):
        u[i] += A[i][j] * v[j]
```

Vectorized:

```
u = np.dot(A, v)
```

Vectors and matrix valued functions Say you need to apply the exponential operation on every element of a matrix/vector.

$$\mathbf{v} = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \longrightarrow \mathbf{u} = \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

Non-vectorized:

```
u = np.zeros((n, 1))
for i in range(n):
    u[i] = math.exp(v[i])
```

Vectorized:

```
u = np.exp(v)
```

J = 0

// $dw_1 = 0, dw_2 = 0, \dots, dw_{n_x}$

$d\mathbf{w} = \text{np.zeros}((n_x, 1))$

db = 0

for i = 1 **to** m **do**

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

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

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

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

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

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

 // ...

 // $dw_{n_x} = dw_{n_x} + x_{n_x}^{(i)} dz^{(i)}$

$d\mathbf{w} = d\mathbf{w} + \mathbf{x}^{(i)} dz^{(i)}$

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

end

J = J/m

// $dw_1 = dw_1/m$

// $dw_2 = dw_2/m$

// ...

// $dw_{n_x} = dw_{n_x}/m$

$d\mathbf{w} = d\mathbf{w}/m$

db = db/m

Algorithm 2: Replace the inner for-loop by vectorization in logistic regression

Logistic regression derivatives

3.2.3 Vectorizing Logistic Regression

Inference process:

$$z^{(1)} = \mathbf{w}^T \mathbf{x}^{(1)} + b \quad z^{(2)} = \mathbf{w}^T \mathbf{x}^{(2)} + b \quad \dots \quad z^{(m)} = \mathbf{w}^T \mathbf{x}^{(m)} + b$$

$$a^{(1)} = \sigma(z^{(1)}) \quad a^{(2)} = \sigma(z^{(2)}) \quad \dots \quad a^{(m)} = \sigma(z^{(m)})$$

$$\mathbf{X} = \left[\begin{array}{c|c|c|c} & & & \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \\ & & & \end{array} \right] \quad \mathbf{X} \in \mathbb{R}^{n_x \times m}$$

$$\mathbf{Z} = \begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = \mathbf{w}^T \mathbf{X} + \underbrace{\begin{bmatrix} b & b & \dots & b \end{bmatrix}}_{1 \times m}$$

$$\mathbf{A} = \begin{bmatrix} a^{(1)} & a^{(2)} & \dots & a^{(m)} \end{bmatrix} = \sigma(\mathbf{Z})$$

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

$$d\mathbf{Z} = \begin{bmatrix} dz^{(1)} & dz^{(2)} & \dots & dz^{(m)} \end{bmatrix}$$

$$\mathbf{Y} = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}$$

$$d\mathbf{Z} = \mathbf{A} - \mathbf{Y} = \begin{bmatrix} a^{(1)} - y^{(1)} & a^{(2)} - y^{(2)} & \dots & a^{(m)} - y^{(m)} \end{bmatrix}$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)} \longrightarrow db = \frac{1}{m} \text{np.sum}(d\mathbf{Z})$$

$$d\mathbf{w} = \frac{1}{m} \left[\begin{array}{c|c|c|c} & & & \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \\ & & & \end{array} \right] \begin{bmatrix} dz^{(1)} \\ dz^{(2)} \\ \vdots \\ dz^{(m)} \end{bmatrix} = \frac{1}{m} \underbrace{\begin{bmatrix} \mathbf{x}^{(1)} dz^{(1)} + \mathbf{x}^{(2)} dz^{(2)} + \dots + \mathbf{x}^{(m)} dz^{(m)} \end{bmatrix}}_{n \times 1} \longrightarrow d\mathbf{w} = \frac{1}{m} \mathbf{X} d\mathbf{Z}^T$$

Vectorized (in a single iteration):

$$\mathbf{Z} = \mathbf{w}^T \mathbf{X} + b \longrightarrow \text{np.dot}(\mathbf{w.T}, \mathbf{X}) + b$$

$$\mathbf{A} = \sigma(\mathbf{Z})$$

$$d\mathbf{Z} = \mathbf{A} - \mathbf{Y}$$

$$d\mathbf{w} = \frac{1}{m} \mathbf{X} (d\mathbf{Z}^T)$$

$$db = \frac{1}{m} \text{np.sum}(d\mathbf{Z})$$

Then update \mathbf{w} and b for each iteration:

$$\mathbf{w} = \mathbf{w} - \alpha d\mathbf{w}$$

$$b = b - \alpha db$$

	Apples	Beef	Eggs	Potatoes
Carb	56.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
Fat	1.8	135.0	99.0	0.9

Figure 7: Calories from Carbs, Proteins, Fats in 100g of different foods

3.2.4 Broadcasting in Python

According to Figure 7, denote the matrix in the figure as \mathbf{A} , calculate the percentages of calories from Carbs, Proteins, Fats for each of four foods. Try to do that without explicit for-loop.

```
cal = A.sum(axis=0) # sum vertically, pass axis=1 to sum horizontally
# broadcasting, A.shape = (3, 4), cal.shape = (1, 4)
percentage = 100 * A / (cal.reshape(1, 4))
```

Broadcasting examples

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix}$$

(Note: In MATLAB/Octave, the `bsxfun` seems to have the same feature, it can apply element-wise operation to two arrays with implicit expansion enabled.)

General principle of broadcasting

$$\underbrace{\text{Matrix}}_{(m,n)} \begin{matrix} + \\ - \\ * \\ / \end{matrix} \underbrace{\text{Vector}}_{(1,n) \text{ or } (m,1)} \longrightarrow \underbrace{\text{Matrix}}_{(m,n)}$$

3.2.5 A Note on Python/NumPy Vectors

For broadcasting:

$\left\{ \begin{array}{ll} \text{Strength:} & \text{expressivity, flexibility} \\ \text{Weakness:} & \text{may introduce subtle bugs} \end{array} \right.$

```
>>> import numpy as np
>>> a = np.random.randn(5)
>>> a # rank 1 array, do not use in this course
array([-0.40700705, -1.27728182,  0.10255894,  0.06871343, -1.01614418])
>>> a.shape
(5,)
>>> a = a.reshape((5, 1)) # column vector
```

```

>>> a
array([[ -0.40700705],
       [ -1.27728182],
       [  0.10255894],
       [  0.06871343],
       [ -1.01614418]])
>>> a = np.random.randn(5, 1)  # column vector
>>> a
array([[ 1.32238445],
       [-2.09932992],
       [ 0.29040232],
       [-1.31295795],
       [-1.05893313]])
>>> a.shape
(5, 1)
>>> a = np.random.randn(1, 5)  # row vector
>>> a
array([[ -0.7020727 , -0.54867968,  1.38088584,  0.69785511, -0.2510519 ]])
>>> a.shape
(1, 5)
>>> assert(a.shape == (1, 5))
>>> assert(a.shape == (5, 1))
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
AssertionError

```

3.2.6 Explanation of Logistic Regression Cost Function (Optional)

Logistic regression cost function:

$$\hat{y} = \sigma(\mathbf{w}^T \mathbf{x} + b), \quad \text{where } \sigma(z) = \frac{1}{1 + e^{-z}}.$$

$$\left. \begin{array}{ll} \text{If } y = 1: & P(y|x) = \hat{y} \\ \text{If } y = 0: & P(y|x) = 1 - \hat{y} \end{array} \right\} P(y|x) = \hat{y}^y (1 - \hat{y})^{(1-y)}$$

$$\log P(y|x) = y \log \hat{y} + (1 - y) \log(1 - \hat{y}) = -\mathcal{L}(\hat{y}, y)$$

Cost on m examples: (Maximum Likelihood Estimation)

$$\log P(\text{labels in training set}) = \log \prod_{i=1}^m P(y^{(i)}|x^{(i)}) = \underbrace{\sum_{i=1}^m \log P(y^{(i)}|x^{(i)})}_{-\mathcal{L}(\hat{y}^{(i)}, y^{(i)})} = - \sum_{(i=1)}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

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

minimize

4 Shallow Neural Network

4.1 Neural Network Overview

Figure 8 shows an example of a shallow neural network.

Notation: Use superscript $[i]$ to denote the i -th layer; use the superscript (i) to denote the i -th example in dataset.

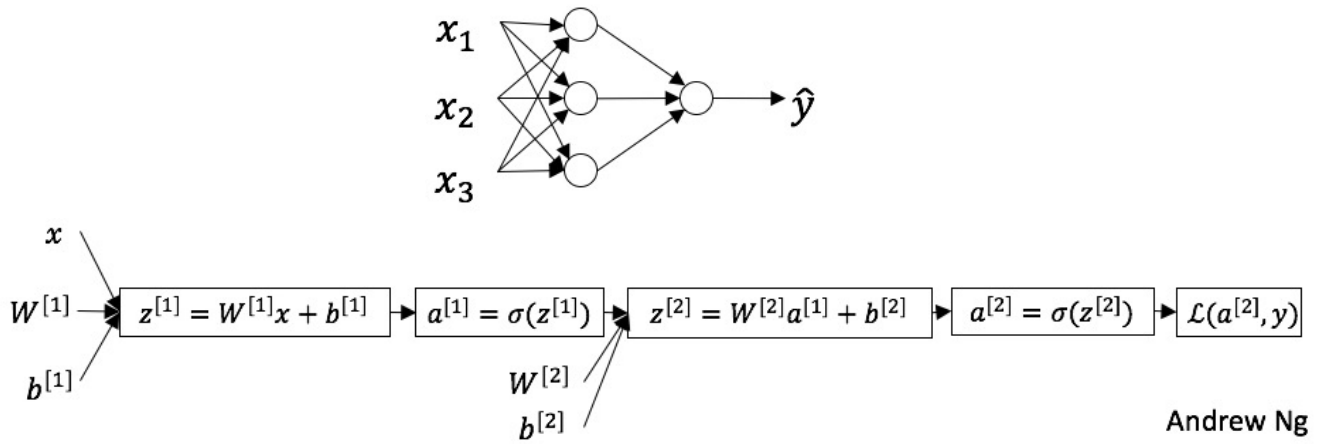


Figure 8: The example of a shallow neural network

4.2 Neural Network Representation

As is shown in Figure 9, $a^{[i]}$ denotes the activation of the i -th layer.

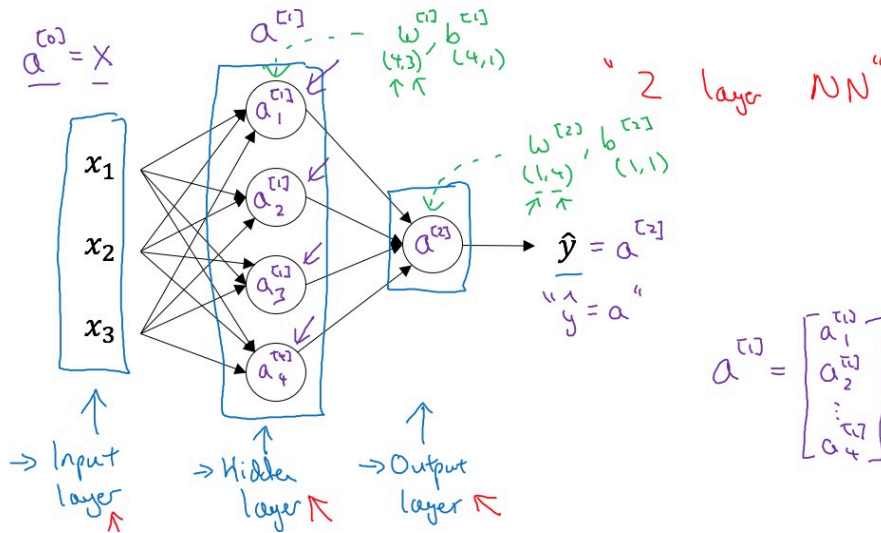


Figure 9: The representation of a neural network

4.3 Computing a Neural Network's Output

As the Figure 10 shows, the computation of a single neuron includes two parts: the first part, $z = w^T x + b$; the second part, $a = \sigma(z)$.

The Figure 11 shows the computation of the neural network's first layer. Cause it uses a for-loop, it is less efficient. According to it, we can generalize a vectorized version:

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix} \quad W^{[1]} = \begin{bmatrix} \text{---} w_1^{[1]T} \text{---} \\ \text{---} w_2^{[1]T} \text{---} \\ \text{---} w_3^{[1]T} \text{---} \\ \text{---} w_4^{[1]T} \text{---} \end{bmatrix} \quad x = a^{[0]} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad b = \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}$$

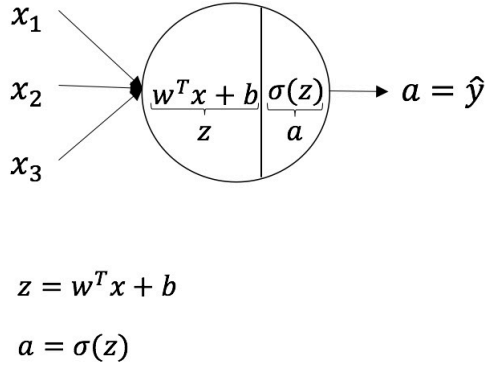


Figure 10: The computation of a single neuron

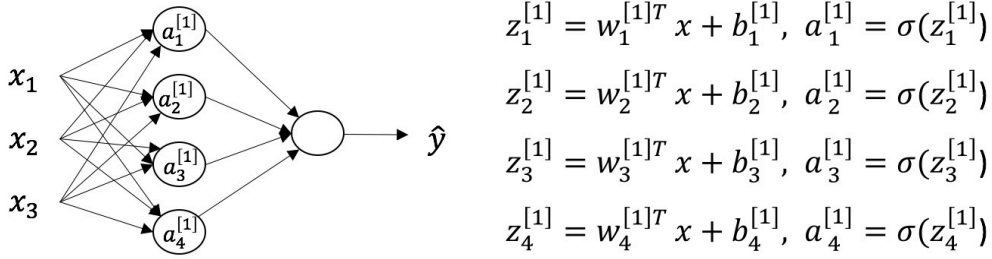


Figure 11: The computation of the first layer of a neural network

Given input $\mathbf{x}(\mathbf{a}^{[0]})$:

$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]} + \mathbf{b}^{[1]} \quad \mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]}) \quad \mathbf{z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]} + \mathbf{b}^{[2]} \quad \mathbf{a}^{[2]} = \sigma(\mathbf{z}^{[2]})$$

4.4 Vectorizing across Multiple Examples

For single example:

$$\mathbf{x} \longrightarrow \mathbf{a}^{[2]} = \hat{\mathbf{y}}$$

For m examples:

$$\mathbf{x}^{(1)} \longrightarrow \mathbf{a}^{[2](1)} = \hat{\mathbf{y}}^{(1)} \quad \mathbf{x}^{(2)} \longrightarrow \mathbf{a}^{2} = \hat{\mathbf{y}}^{(2)} \quad \dots \quad \mathbf{x}^{(m)} \longrightarrow \mathbf{a}^{[2](m)} = \hat{\mathbf{y}}^{(m)}$$

Non-vectorized:

```

for  $i = 1$  to  $m$  do
     $\mathbf{z}^{[1](i)} = \mathbf{W}^{[1]} \mathbf{x}^{(i)} + \mathbf{b}^{[1]}$ 
     $\mathbf{a}^{[1](i)} = \sigma(\mathbf{z}^{[1](i)})$ 
     $\mathbf{z}^{[2](i)} = \mathbf{W}^{[2]} \mathbf{a}^{[1](i)} + \mathbf{b}^{[2]}$ 
     $\mathbf{a}^{[2](i)} = \sigma(\mathbf{z}^{[2](i)})$ 
end

```


$$\text{Let } \mathbf{X} = \underbrace{\begin{bmatrix} | & | & | & | \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \\ | & | & | & | \end{bmatrix}}_{(n_x, m)}$$

$$\text{then } \mathbf{Z}^{[1]} = \underbrace{\begin{bmatrix} | & | & | & | \\ z^{1} & z^{[1](2)} & \dots & z^{[1](m)} \\ | & | & | & | \end{bmatrix}}_{(\text{hidden units, training examples})} \quad \mathbf{A}^{[1]} = \underbrace{\begin{bmatrix} | & | & | & | \\ \mathbf{a}^{1} & \mathbf{a}^{[1](2)} & \dots & \mathbf{a}^{[1](m)} \\ | & | & | & | \end{bmatrix}}_{(\text{hidden units, training examples})}$$

Vectorized:

$$\mathbf{Z}^{[1]} = \mathbf{W}^{[1]} \mathbf{A}^{[0]} + \mathbf{b}^{[1]} = \mathbf{W}^{[1]} \mathbf{X} + \mathbf{b}^{[1]}$$

$$\mathbf{A}^{[1]} = \sigma(\mathbf{Z}^{[1]})$$

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

$$\mathbf{A}^{[2]} = \sigma(\mathbf{Z}^{[2]})$$

4.5 Explanation for Vectorized Implementation

$$z^{1} = \mathbf{W}^{[1]} \mathbf{x}^{(1)} + \mathbf{b}^{[1]}, \quad z^{[1](2)} = \mathbf{W}^{[1]} \mathbf{x}^{(2)} + \mathbf{b}^{[1]}, \quad \dots, \quad z^{[1](m)} = \mathbf{W}^{[1]} \mathbf{x}^{(m)} + \mathbf{b}^{[1]}$$

$$\mathbf{W}^{[1]} = \begin{bmatrix} \text{---} w_1^{[1]T} \text{---} \\ \text{---} w_2^{[1]T} \text{---} \\ \text{---} w_3^{[1]T} \text{---} \\ \text{---} w_4^{[1]T} \text{---} \end{bmatrix}$$

$$\mathbf{W}^{[1]} \mathbf{x}^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad \mathbf{W}^{[1]} \mathbf{x}^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad \dots \quad \mathbf{W}^{[1]} \mathbf{x}^{(m)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$\mathbf{W}^{[1]} \begin{bmatrix} | & | & | & | \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \\ | & | & | & | \end{bmatrix} + \mathbf{b}^{[1]} = \begin{bmatrix} \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \end{bmatrix} + \mathbf{b}^{[1]} = \begin{bmatrix} | & | & | & | \\ \mathbf{z}^{1} & \mathbf{z}^{[1](2)} & \dots & \mathbf{z}^{[1](m)} \\ | & | & | & | \end{bmatrix} = \mathbf{Z}^{[1]}$$

So, $\mathbf{W}^{[1]} \mathbf{x}^{(i)} + \mathbf{b}^{[1]} = z^{[1](i)}$, $\mathbf{Z}^{[1]} = \mathbf{W}^{[1]} \mathbf{X} + \mathbf{b}^{[1]}$.

4.6 Activation Functions

Given \mathbf{x} :

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

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

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

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

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

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

We use $g(z)$ to denote activation function. Sometimes we can use other activation function like $\tanh(z)$, ReLU ($\max(0, z)$), etc., instead of $\sigma(z)$ function. Different layers of the neural network may have different activation function, so we can use $g^{[i]}(z)$ to denote the i -th layer's activation function of a neural network.

4.6.1 Pros and cons of activation functions

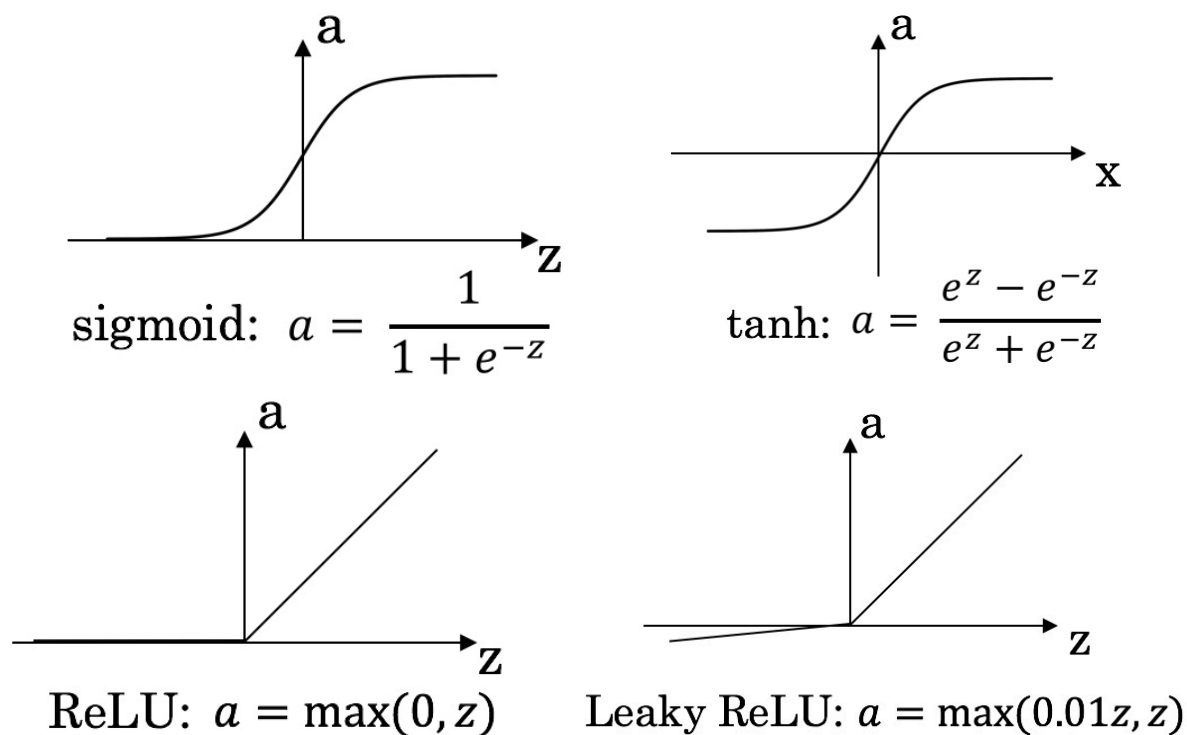


Figure 12: Activation functions

- One of downside of both the sigmoid and tanh activation function is that if z is either very large or very small, then the gradient (derivative, slope) of the function becomes very small, this slows down the gradient descent.
- Never use sigmoid activation function, except for the output layer if you are doing binary classification. Tanh is pretty superior than sigmoid.
- If you don't have an idea to choose which activation function, choose ReLU as the default activation function, or you can try Leaky ReLU.

4.7 Why do you need non-linear activation functions?

If $g(z) = z$, sometimes we will call it “linear activation function” or “identity activation function”.
Given \mathbf{x} :

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

$$\mathbf{a}^{[1]} = g^{[1]}(\mathbf{z}^{[1]}) = \mathbf{z}^{[1]} = \mathbf{W}^{[1]}\mathbf{x} + \mathbf{b}^{[1]}$$

$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]}\mathbf{a}^{[1]} + \mathbf{b}^{[2]} = \mathbf{W}^{[2]}(\mathbf{W}^{[1]}\mathbf{x} + \mathbf{b}^{[1]}) + \mathbf{b}^{[2]} = \underbrace{(\mathbf{W}^{[2]}\mathbf{W}^{[1]})}_{\mathbf{W}'}\mathbf{x} + \underbrace{(\mathbf{W}^{[2]}\mathbf{b}^{[1]} + \mathbf{b}^{[2]})}_{\mathbf{b}'} = \mathbf{W}'\mathbf{x} + \mathbf{b}'$$

$$\mathbf{a}^{[2]} = g^{[2]}(\mathbf{z}^{[2]}) = \mathbf{z}^{[2]} = \mathbf{W}'\mathbf{x} + \mathbf{b}'$$

If you use linear activation function as the activation function of your neural network, then the neural network is just outputting a linear function of the input. It turns out that if you use a linear function or alternatively if you don't have an activation function, then no matter how many layers your neural network has always doing is just computing a linear activation function, so you might as well not have any hidden layers.

4.8 Derivatives of activation functions

4.8.1 Sigmoid activation function

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

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

4.8.2 Tanh activation function

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

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

4.9 ReLU

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

$$g'(z) = \begin{cases} 0, & \text{if } z < 0; \\ 1, & \text{if } z > 0; \\ \text{undefined}, & \text{if } z = 0. \end{cases} \xrightarrow{\text{implemented in software}} g'(z) = \begin{cases} 0, & \text{if } z < 0; \\ 1, & \text{if } z \geq 0. \end{cases}$$

4.10 Leaky ReLU

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

$$g'(z) = \begin{cases} 0.01, & \text{if } z < 0; \\ 1, & \text{if } z > 0; \\ \text{undefined}, & \text{if } z = 0. \end{cases} \xrightarrow{\text{implemented in software}} g'(z) = \begin{cases} 0.01, & \text{if } z < 0; \\ 1, & \text{if } z \geq 0. \end{cases}$$

4.11 Gradient descent for Neural Networks

Parameters: $n_x = n^{[0]}$, $n^{[1]}$, $n^{[2]} = 1$, $\underbrace{\mathbf{W}^{[1]}}_{(n^{[1]}, n^{[0]})}$, $\underbrace{\mathbf{b}^{[1]}}_{(n^{[2]}, n^{[1]})}$, $\underbrace{\mathbf{W}^{[2]}}_{(n^{[2]}, n^{[1]})}$, $\underbrace{\mathbf{b}^{[2]}}_{(n^{[2]}, 1)}$

$$\text{Data: } \mathbf{X} = \underbrace{\begin{bmatrix} | & | & & | \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \\ | & | & & | \end{bmatrix}}_{(n^{[0]}, m)}, \mathbf{Y} = \underbrace{\begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}}_{(1, m)}$$

Cost Function:

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

repeat

```
// Forward propagation
 $\mathbf{Z}^{[1]} = \mathbf{W}^{[1]} \mathbf{X} + \mathbf{b}^{[1]}$ 
 $\mathbf{A}^{[1]} = g^{[1]}(\mathbf{Z}^{[1]})$ 
 $\mathbf{Z}^{[2]} = \mathbf{W}^{[2]} \mathbf{A}^{[1]} + \mathbf{b}^{[2]}$ 
 $\mathbf{A}^{[2]} = g^{[2]}(\mathbf{Z}^{[2]}) = \sigma(\mathbf{Z}^{[2]})$ 
// Backpropagation
 $d\mathbf{Z}^{[2]} = \mathbf{A}^{[2]} - \mathbf{Y}$ 
 $d\mathbf{W}^{[2]} = \frac{1}{m} d\mathbf{Z}^{[2]} \mathbf{A}^{[1]T}$ 
// keepdims=True ensures that the shape of result is  $(n^{[2]}, 1)$ 
 $d\mathbf{b}^{[2]} = \frac{1}{m} \text{np.sum}(d\mathbf{Z}^{[2]}, \text{axis}=1, \text{keepdims}=\text{True})$ 
// here * means element-wise multiplication
 $d\mathbf{Z}^{[1]} = \mathbf{W}^{[2]T} d\mathbf{Z}^{[2]} * g^{[1]'}(\mathbf{Z}^{[1]})$ 
 $d\mathbf{W}^{[1]} = \frac{1}{m} d\mathbf{Z}^{[1]} \mathbf{X}^T$ 
 $d\mathbf{b}^{[1]} = \frac{1}{m} \text{np.sum}(d\mathbf{Z}^{[1]}, \text{axis}=1, \text{keepdims}=\text{True})$ 
```

until iteration = MaxIteration;

Algorithm 3: The gradient descent algorithm for neural networks

4.12 Backpropagation intuition (Optional)

According to Figure 6, we have got how to compute the gradient of logistic regression. By the chain rule, we can get:

$$\text{If } a = g(z), \text{ then } dz = da \cdot g'(z)$$

Logistic regression can be seen as a neural network with only one layers (exclude input layer), while the neural network in Figure 8 and Figure 9 contains two layers. To compute the derivatives for each single example, we can also use backpropagation and the chain rule.

At first, because $\mathcal{L}(\hat{y}, y) = -(y \log \hat{y} + (1 - y) \log(1 - \hat{y}))$, we can compute $da^{[2]}$:

$$da^{[2]} = \frac{d\mathcal{L}(a^{[2]}, y)}{da^{[2]}} = -\frac{y}{a^{[2]}} + \frac{1 - y}{1 - a^{[2]}}$$

Then, according to the chain rule, we can compute $dz^{[2]}$:

$$dz^{[2]} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} = \frac{d\mathcal{L}}{da^{[2]}} g^{[2]'}(z^{[2]}) = \left(-\frac{y}{a^{[2]}} + \frac{1 - y}{1 - a^{[2]}}\right) [z^{[2]}(1 - z^{[2]})] = a^{[2]} - y$$

Because $z^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]} + b^{[2]}$, once we get $dz^{[2]}$, we can continue to compute $d\mathbf{W}^{[2]}$ and $db^{[2]}$:

$$d\mathbf{W}^{[2]} = dz^{[2]} \mathbf{a}^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

To continue it, we can get $d\mathbf{a}^{[1]} = \mathbf{W}^{[2]T} dz^{[2]}$. In practice, the computation of $d\mathbf{a}^{[1]}$ and $dz^{[1]}$ often collapses into one step. And we can check the dimensions, $\mathbf{W}^{[2]} \in \mathbb{R}^{(n^{[2]}, n^{[1]})}$, $z^{[2]}, dz^{[2]} \in \mathbb{R}^{(n^{[2]}, 1)}$, $\mathbf{z}^{[1]}, d\mathbf{z}^{[1]} \in \mathbb{R}^{(n^{[1]}, 1)}$ (the $*$ here means element-wise multiplication):

$$d\mathbf{z}^{[1]} = d\mathbf{a}^{[1]} * g^{[1]'}(\mathbf{z}^{[1]}) = \underbrace{\underbrace{\mathbf{W}^{[2]T}}_{(n^{[1]}, n^{[2]})} \underbrace{dz^{[2]}}_{(n^{[2]}, 1)}}_{(n^{[1]}, 1)} * \underbrace{g^{[1]'}(\mathbf{z}^{[1]})}_{(n^{[1]}, 1)}$$

At last, we can get $d\mathbf{W}^{[1]}$ and $d\mathbf{b}^{[1]}$:

$$d\mathbf{W}^{[1]} = dz^{[1]} \underbrace{\mathbf{x}^T}_{\mathbf{a}^{[0]T}}$$

$$d\mathbf{b}^{[1]} = dz^{[1]}$$

4.12.1 Summary of gradient descent

For single example:

$$dz^{[2]} = \mathbf{a}^{[2]} - \mathbf{y}$$

$$d\mathbf{W}^{[2]} = dz^{[2]} \mathbf{a}^{[1]T}$$

$$d\mathbf{b}^{[2]} = dz^{[2]}$$

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

$$d\mathbf{W}^{[1]} = dz^{[1]} \mathbf{x}^T$$

$$d\mathbf{b}^{[1]} = dz^{[1]}$$

Vectorized implementation for m examples:

$$d\mathbf{Z}^{[2]} = \mathbf{A}^{[2]} - \mathbf{Y}$$

$$d\mathbf{W}^{[2]} = \frac{1}{m} d\mathbf{Z}^{[2]} \mathbf{A}^{[1]T}$$

$$d\mathbf{b}^{[2]} = \frac{1}{m} \text{np.sum}(d\mathbf{Z}^{[2]}, \text{axis}=1, \text{keepdims=True})$$

$$d\mathbf{Z}^{[1]} = \mathbf{W}^{[2]T} d\mathbf{Z}^{[2]} * g^{[1]'}(\mathbf{Z}^{[1]})$$

$$d\mathbf{W}^{[1]} = \frac{1}{m} d\mathbf{Z}^{[1]} \mathbf{X}^T$$

$$d\mathbf{b}^{[1]} = \frac{1}{m} \text{np.sum}(\mathbf{Z}^{[1]}, \text{axis}=1, \text{keepdims=True})$$

4.13 Random Initialization

When you are training a neural network, it's important to initialize your weights randomly. For logistic regression, it was okay to initialize the weights to 0, but for a neural network of initialized weights and parameters to all 0 and then apply gradient descent, it won't work.

As the Figure 13 shows, there is a shallow neural network with two input features, two hidden units and one output unit. Then we initialized all the weights and parameters to zero:

$$\mathbf{W}^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad \mathbf{b}^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \mathbf{W}^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix} \quad \mathbf{b}^{[2]} = [0]$$

It turns out that initializing the bias terms \mathbf{b} to all zero is exactly okay, but initializing \mathbf{W} to all zero causes a problem. So the problem with this formalization is that for any example you give it, you'll have that $a_1^{[1]} = a_2^{[2]}$, and then when you compute backpropagation, it turns out that $dz_1^{[1]} = dz_2^{[2]}$.

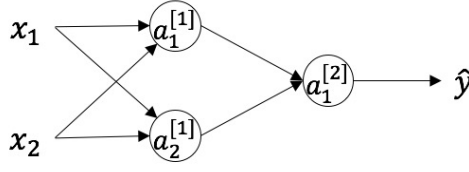


Figure 13: A shallow neural network with all weights and parameters initialized to 0

If we initialize the neural network this way, the hidden unit $a_1^{[1]}$ and $a_2^{[1]}$ will be completely identical (symmetric). It turns out that after every single iteration of training, the two hidden units are still computing exactly the same function, both units have the same influence on the output unit. The $d\mathbf{W}^{[1]}$ and $\mathbf{W}^{[1]}$ will look like as the following:

$$d\mathbf{W} = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \quad \mathbf{W}^{[1]} = \mathbf{W}^{[1]} - \alpha d\mathbf{W} \quad \mathbf{W}^{[1]} = \begin{bmatrix} r & s \\ r & s \end{bmatrix}$$

So in this case there is really no point to have more than one hidden unit, because all units are computing the same thing. However, that's not helpful, because you want different units compute different functions. The solution is to initialize the weights randomly:

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

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

Why we multiply a factor 0.01? That's because when we use the sigmoid or tanh activation function, if the weight is large, then the value input to the activation function will be large, so in that case the slope of the gradient is small and the gradient descent will be very slow.

5 Deep Neural Networks

5.1 Deep L -layer neural network

What is a deep neural network? As Figure 14 shows, logistic regression can be seen as one layer neural network. When we count the layers of the neural network, we don't count the input layer. Technically logistic regression is a quite shallow neural network, but over the last several years, the AI and machine learning community has realized that there are functions that very deep neural networks can learn while shallower are unable to.

Deep neural network notation Figure 15 shows a 4-layer neural network, with the input layer as "layer 0", three hidden layers "layer 1", "layer 2" and "layer 3", the output layer as "layer 4". Some notations should be remembered as the following:

L : # layers. In Figure 15, $L = 4$.

$n^{[l]}$: # units in layer l . In this case, $n^{[1]} = 5, n^{[2]} = 5, n^{[3]} = 3, n^{[4]} = n^{[L]} = 1, n^{[0]} = n_x = 3$.

$\mathbf{a}^{[l]}$: activations in layer l . In propagation, we will implement $\mathbf{a}^{[l]} = g^{[1]}(\mathbf{z}^{[l]})$. For single example, $\mathbf{a}^{[0]} = \mathbf{x}$, $\mathbf{a}^{[L]} = \hat{\mathbf{y}}$

$\mathbf{W}^{[l]}$: weights for $\mathbf{z}^{[l]}$.

5.2 Forward Propagation in a Deep Network

We still look at the Figure 15, see the forward propagation in this neural network.

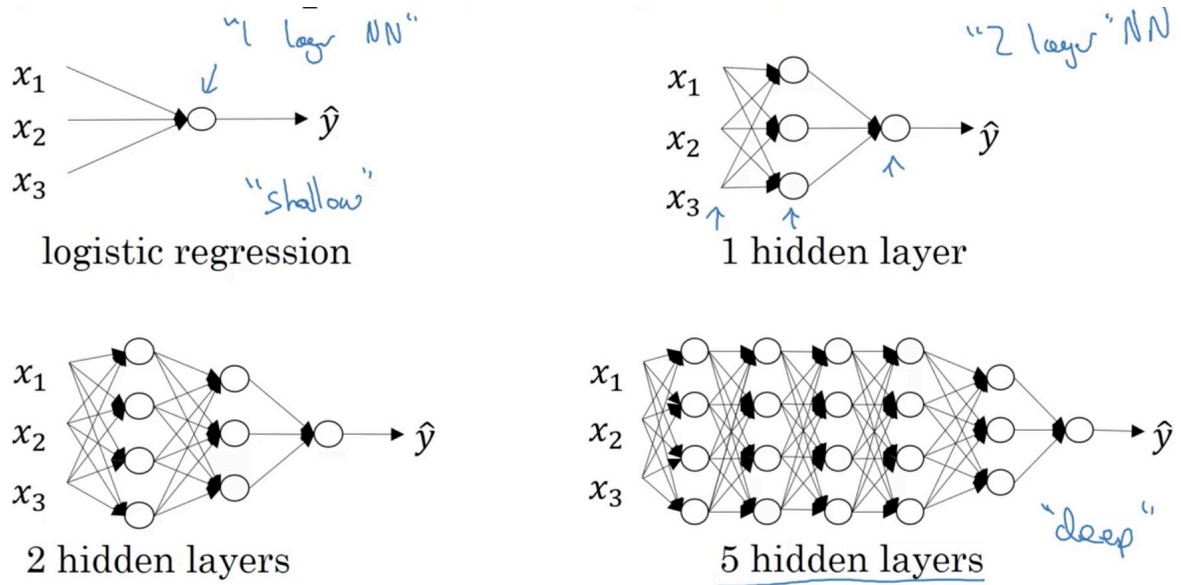


Figure 14: Examples of neural networks

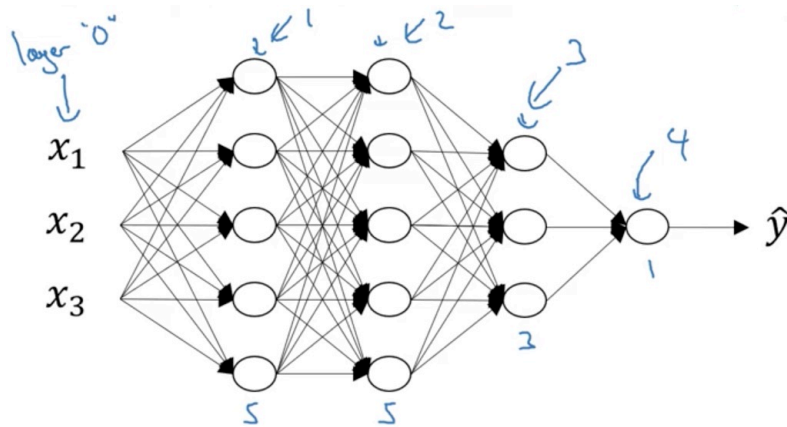


Figure 15: A 4-layer neural network

For every single example,

$$z^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]} + \mathbf{b}^{[1]} = \mathbf{W}^{[1]} \mathbf{x} + \mathbf{b}^{[1]}$$

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

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

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

$$z^{[3]} = \mathbf{W}^{[3]} \mathbf{a}^{[2]} + \mathbf{b}^{[3]}$$

$$\mathbf{a}^{[3]} = g^{[3]}(z^{[3]})$$

$$z^{[4]} = \mathbf{W}^{[4]} \mathbf{a}^{[3]} + \mathbf{b}^{[4]}$$

$$\mathbf{a}^{[4]} = g^{[4]}(\mathbf{z}^{[4]}) = \hat{\mathbf{y}}$$

Then we can get the general rule for forward propagation:

$$\mathbf{z}^{[l]} = \mathbf{W}^{[l]} \mathbf{a}^{[l-1]} + \mathbf{b}^{[l]}$$

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

For the whole training set, in a vectorized way:

$$\mathbf{Z}^{[1]} = \mathbf{W}^{[1]} \mathbf{A}^{[0]} + \mathbf{b}^{[1]} = \mathbf{W}^{[1]} \mathbf{X} + \mathbf{b}^{[1]}$$

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

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

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

$$\mathbf{Z}^{[3]} = \mathbf{W}^{[3]} \mathbf{A}^{[2]} + \mathbf{b}^{[3]}$$

$$\mathbf{A}^{[3]} = g^{[3]}(\mathbf{Z}^{[3]})$$

$$\mathbf{Z}^{[4]} = \mathbf{W}^{[4]} \mathbf{A}^{[3]} + \mathbf{b}^{[4]}$$

$$\mathbf{A}^{[4]} = g^{[4]}(\mathbf{Z}^{[4]}) = \hat{\mathbf{Y}}$$

We can put the left into a explicit for-loop:

for $l = 1$ **to** 4 **do**

$$\mathbf{Z}^{[l]} = \mathbf{W}^{[l]} \mathbf{A}^{[l-1]} + \mathbf{b}^{[l]}$$

$$\mathbf{A}^{[l]} = g^{[l]}(\mathbf{Z}^{[l]})$$

end

5.3 Getting your matrix dimensions right

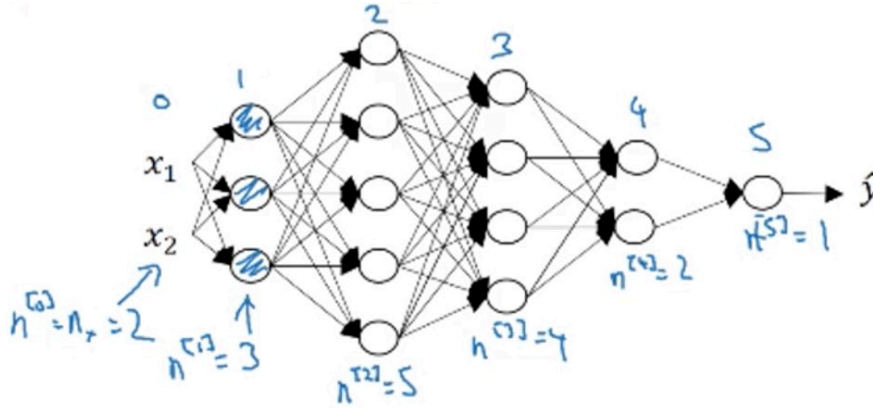


Figure 16: A 5-layer neural network

In Figure 16, there is a five-layer neural network, with $n^{[0]} = n_x = 2$, $n^{[1]} = 3$, $n^{[2]} = 5$, $n^{[3]} = 4$, $n^{[4]} = 2$ and $n^{[5]} = 1$. The checking of the matrix dimensions can be seen as following.

For every single example

$$\underbrace{\mathbf{z}^{[1]}}_{(n^{[1]},1)=(3,1)} = \underbrace{\mathbf{W}^{[1]}}_{(n^{[1]},n^{[0]})=(3,2)} \underbrace{\mathbf{x}}_{(n^{[0]},1)=(2,1)} + \underbrace{\mathbf{b}^{[1]}}_{(n^{[1]},1)=(3,1)}$$

(3,1)

$$\underbrace{\mathbf{z}^{[2]}}_{(n^{[2]},1)=(5,1)} = \underbrace{\mathbf{W}^{[2]}}_{(n^{[2]},n^{[1]})=(5,3)} \underbrace{\mathbf{a}^{[1]}}_{(3,1)} + \underbrace{\mathbf{b}^{[2]}}_{(n^{[2]},1)=(5,1)}$$

.....

$$\mathbf{W}^{[3]} : (n^{[3]}, n^{[2]}) = (4, 5) \quad \mathbf{W}^{[4]} : (n^{[4]}, n^{[3]}) = (2, 4) \quad \mathbf{W}^{[5]} : (n^{[5]}, n^{[4]}) = (1, 2)$$

$$\mathbf{z}^{[3]}, \mathbf{a}^{[3]}, \mathbf{b}^{[3]} : (n^{[3]}, 1) = (4, 1) \quad \mathbf{z}^{[4]}, \mathbf{a}^{[4]}, \mathbf{b}^{[4]} : (n^{[4]}, 1) = (3, 1) \quad \mathbf{z}^{[5]}, \mathbf{a}^{[5]}, \mathbf{b}^{[5]} : (n^{[5]}, 1) = (1, 1)$$

The general formula to check is that when you're implementing the matrix \mathbf{W} for layer l , the dimension of that matrix will be $(n^{[l]}, n^{[l-1]})$. And if you're implementing the back-propagation, the matrix $d\mathbf{W}$ will have the same dimension with \mathbf{W} , the vector $d\mathbf{b}$ will have the same dimension with \mathbf{b} . So for each single example, the formula is:

$$\mathbf{W}^{[l]} : (n^{[l]}, n^{[l-1]}) \quad d\mathbf{W}^{[l]} : (n^{[l]}, n^{[l-1]})$$

$$\mathbf{z}^{[l]}, \mathbf{a}^{[l]}, \mathbf{b}^{[l]} : (n^{[l]}, 1) \quad d\mathbf{z}^{[l]}, d\mathbf{a}^{[l]}, d\mathbf{b}^{[l]} : (n^{[l]}, 1)$$

Vectorized implementation

$$\underbrace{\mathbf{Z}^{[1]}}_{(n^{[1]},m)} = \underbrace{\mathbf{W}^{[1]}}_{(n^{[1]},n^{[0]})} \underbrace{\mathbf{X}}_{(n^{[0]},m)} + \underbrace{\mathbf{b}^{[1]}}_{(n^{[1]},1) \rightarrow (n^{[1]},m)}$$

.....

Stack all m training examples together, then we can get the vectorized version of the formula:

$$\mathbf{W}^{[l]} : (n^{[l]}, n^{[l-1]}) \quad d\mathbf{W}^{[l]} : (n^{[l]}, n^{[l-1]})$$

$$\mathbf{b}^{[l]} : (n^{[l]}, 1) \quad d\mathbf{b}^{[l]} : (n^{[l]}, 1)$$

$$\mathbf{Z}^{[l]}, \mathbf{A}^{[l]} : (n^{[l]}, m) \quad d\mathbf{Z}^{[l]}, d\mathbf{A}^{[l]} : (n^{[l]}, m)$$

5.4 Why deep representations?

5.4.1 Intuition about deep representation

Like Figure 17 shows, if you are building a system for face recognition or face detection, here's what a deep neural network could be doing. Perhaps you input a picture of a face then the first layer of the neural network you can think of as maybe being a feature detector or an edge detector.

In this example, there are about 20 units in the hidden layer. So the 20 hidden units are visualized by these little square boxes. So for example, this little visualization in the first hidden layer represents the hidden unit is trying to figure out where the edges of that orientation in the image.

Now, let's think about where the edges in this picture by grouping together pixels to form edges. It can then detect the edges and group edges together to form parts of the faces. So for example, you might have a neuron trying to see if it's finding an eye, or a different neuron trying to find that part of the nose. So by putting together lots of edges, it can start to detect different parts of faces.

Finally, by putting together different parts of faces, like an eye or a nose or an ear or a chin, it can try to recognize or detect different types of faces.

So intuitively, you can think of the earlier layers of the neural network as detecting simple functions, like edges. And then composing them together in the later layers of a neural network so that it can learn more and more complex functions. These visualizations will make more sense when we talk about convolutional nets.

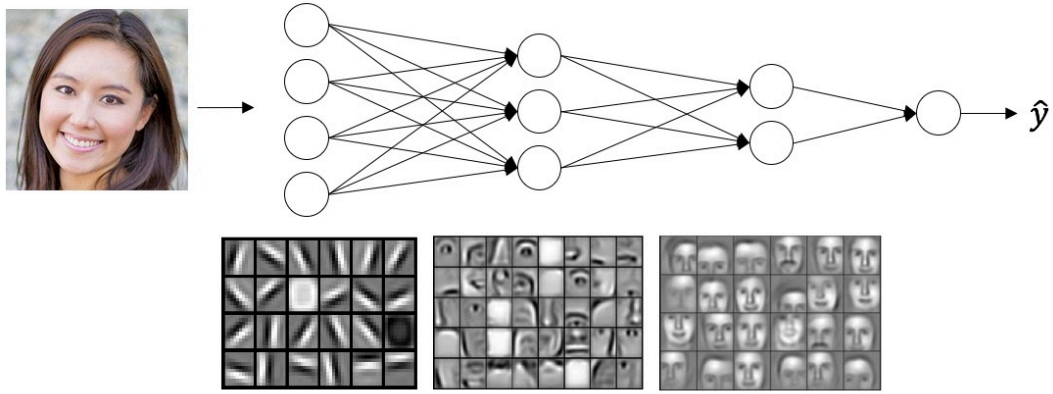


Figure 17: Intuition about deep representation

5.4.2 Circuit theory and deep learning

Informally: There are functions you can compute with a “small” (hidden units) L -layer deep neural network that shallower networks require exponentially more hidden units to compute.

5.5 Building blocks of deep neural networks

For the layer l in the deep neural network:

Parameters: $\mathbf{W}^{[l]}$, $\mathbf{b}^{[l]}$.

Forward: input: $\mathbf{a}^{[l-1]}$; output: $\mathbf{a}^{[l]}$; cache $(\mathbf{z}^{[l-1]})$.

$$\mathbf{z}^{[l]} = \mathbf{W}^{[l]} \mathbf{a}^{[l-1]} + \mathbf{b}^{[l]}$$

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

Vectorized version:

$$\mathbf{Z}^{[l]} = \mathbf{W}^{[l]} \mathbf{A}^{[l-1]} + \mathbf{b}^{[l]}$$

$$\mathbf{A}^{[l]} = g^{[l]}(\mathbf{Z}^{[l]})$$

Backward: input: $d\mathbf{a}^{[l]}$ (get from cache $(\mathbf{z}^{[l-1]})$); output: $d\mathbf{a}^{[l-1]}$, $d\mathbf{W}^{[l]}$, $d\mathbf{b}^{[l]}$;

$$d\mathbf{z}^{[l]} = d\mathbf{a}^{[l]} * g^{[l]'}(\mathbf{z}^{[l]})$$

$$d\mathbf{W}^{[l]} = d\mathbf{z}^{[l]} \mathbf{a}^{[l-1]}$$

$$d\mathbf{b}^{[l]} = d\mathbf{z}^{[l]}$$

$$d\mathbf{a}^{[l-1]} = \mathbf{W}^{[l]T} d\mathbf{z}^{[l]}$$

Vectorized version:

$$d\mathbf{Z}^{[l]} = d\mathbf{A}^{[l]} * g^{[l]'}(\mathbf{Z}^{[l]})$$

$$d\mathbf{W}^{[l]} = \frac{1}{m} d\mathbf{Z}^{[l]} \mathbf{A}^{[l-1]T}$$

$$d\mathbf{b}^{[l]} = \frac{1}{m} \text{np.sum}(d\mathbf{Z}^{[l]}, \text{axis}=1, \text{keepdims}=\text{True})$$

$$d\mathbf{A}^{[l-1]} = \mathbf{W}^{[l]T} d\mathbf{Z}^{[l]}$$

The process above can be concluded as a block in Figure 19.

5.6 Parameters vs Hyperparameters

5.6.1 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 α ; # iterations; # hidden layer L ; # hidden units $n^{[1]}$, $n^{[2]}$, ...; choice of activation function.

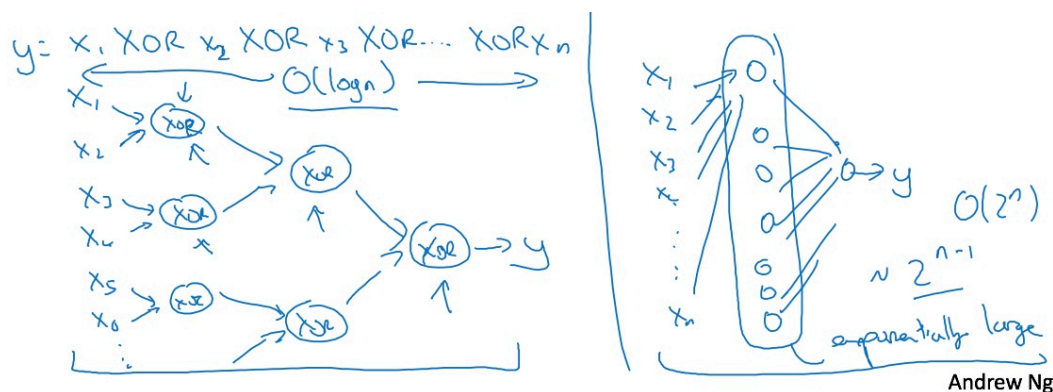


Figure 18: Use deep neural networks and shallower neural networks to implement the XOR of n input. For deep neural network (the left one), to compute the XOR, the depth of the network will be on the order of $O(\log n)$, so the number of nodes, or the number of circuit components or the number of gates in this neural network is not that large. But if you are forced to compute the function with just one hidden layer, then in order to get the parity of XOR to compute this XOR function, the size of this hidden layer will need to be exponentially large (2^{n-1}), because essentially you need to exhaustively enumerate all 2 to the n possible configurations.

Hyperparameters are parameters that control parameters \mathbf{W} and \mathbf{b} . In later course, we will learn other hyperparameters such as momentum, mini-batch size and regularization parameters.

5.6.2 Applied deep learning is a very empirical process

Applied deep learning today is a very empirical process. For example, in Figure 20 you have an idea that set $\alpha = 0.01$, then you implement it, try out and see how that works, and then based on the outcome you might say you know that I have changed my mind, then you do the process again to try $\alpha = 0.05$. You might plot the cost J of the number of iterations to see which is the best learning rate.

The 'empirical process' in the title means that you have to try a lot of things and see what works. When starting on a new problem, just try out a range of values and see what works. In the next course, we'll see some systematic ways for trying out a range of values of hyperparameters.

Secondly, if you're working on one application for a long time, maybe you're working on online advertising as you make progress on the problem is quite possible there the best value for the learning rate, a number of hidden units and so on might change, even if you tune your system to the best value of hyperparameters today, it's possible that you find that the best value might change year from now.

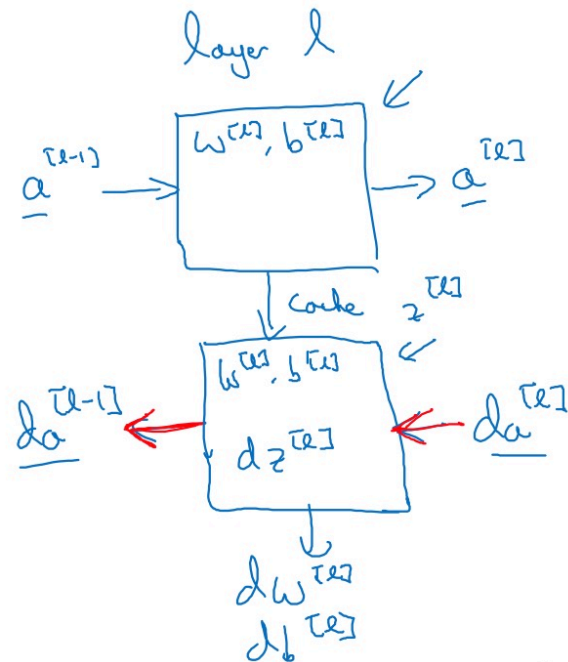


Figure 19: A block of forward and backward process in a deep Neural network

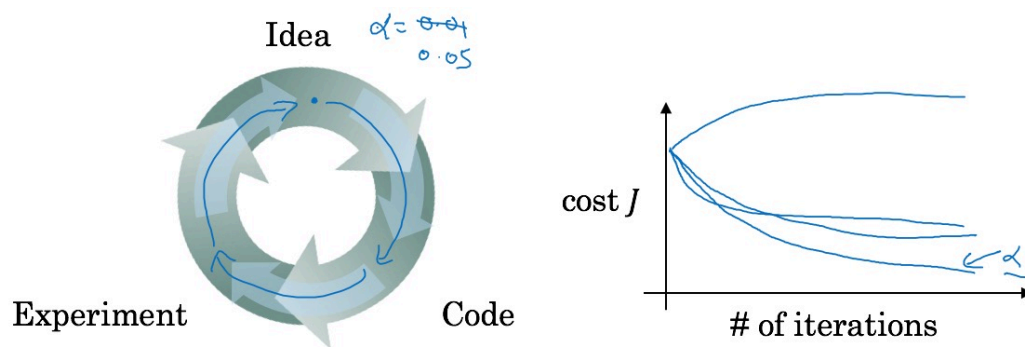


Figure 20: The process of try different hyperparamanters in a deep Neural network