

Deep Learning

Labix

September 24, 2024

Abstract

Contents

1	Foundations of Deep Learning	2
1.1	Neural Networks	2
1.2	Binary Classification	4
1.3	Random Initialization	5
1.4	Forward Propagation	6
1.5	Lost and Cost Function	7
1.6	Backwards Propagation	8
1.7	Putting Everything Together	13
2	Improving the Adaptability of the Prediction	14
2.1	Hypervariables	14
2.2	Data Sets	14
2.3	Bias and Variance	14
2.4	Regularization	15
3	Improving the Speed of the Algorithm	18
3.1	Normalization	18
3.2	Weight Initialization	18
3.3	Minibatch	19
3.4	Optimization of the Gradient Descent	19
3.5	Learning Rates	19
4	Convolutional Neural Networks	20
4.1	Matrix Convolutions	20
4.2	Convolutions in Neural Networks	20
4.3	Types of Convolutions	21
5	Deep Neural Networks	22
6	Recurrent Neural Networks	23
7	Large Language Models	24
7.1	Data Preprocessing	24
7.2	Attention	25
7.3	Transformer Blocks	25

1 Foundations of Deep Learning

Deep learning involves mostly neural networks and using given inputs and outputs to find out the intermediate process that determines the relation. In other words, we are given a set of inputs and a set of outputs, and we want to find out the function relating these sets.

Deep learning is used in a lot of different places such as photo recognition, speech recognition and real state price estimations. There are different neural network models that allows a more accurate prediction for different data structures and maybe specific to the questions itself.

Inputs can be roughly separated into two types: Structured data and unstructured data. Structured data usually can be aligned into a table while unstructured data include audio, images and texts.

Deep learning is now taking off mainly because of the size of data that we possess as compared to a few years ago.

1.1 Neural Networks

We begin with the notion of a node.

Definition 1.1.1: Nodes

A node takes a vector x of n inputs and outputs a single value a using the follow equation

$$a = \varphi(wx + b)$$

where w, b, φ are called parameters of the node. w is a $1 \times n$ matrix and b is a real number. $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is called an activation function. The vector x is usually the output of the previous layer.

We will explain what the activation functions mean. For now, please ignore it and simply consider it to be part of the necessary data of one node.

Definition 1.1.2: Layers

A layer is a collection of ordered nodes of the same depth. The depth is denoted l , which means that it takes l steps to transform all the initial inputs to a node in layer l .

Layer 0, consisting of all the initial inputs are called the input layer. The final layer L in which there are no further outputs for all nodes are called the output layer is called the output layer. The rest are called the hidden layer.

A neural network consists of inputs, intermediate nodes and outputs. It is also further separated into layers. The input nodes are grouped to called the input layer, similarly for the output layer. The rest of the nodes are part of the hidden layer. Nodes connected immediately from the input nodes are the first layer, nodes connected immediately from the first layer forms the second layer and so on. They are usually called hidden layer 1, hidden layer 2 and so on.

Definition 1.1.3: Neural Networks

A neural network consists of the following data.

- The total number of layers including the input and output layers
- The total number of nodes $n^{[l]}$ for each layer l
- The connections between the nodes of a previous layer and the current layer

Suppose that on layer l , there are $n^{[l]}$ nodes. Each node has its own w and b labeled $w_1^{[l]}, \dots, w_{n^{[l]}}^{[l]}$ and

similarly for b , which allows us to form the matrices

$$W^{[l]} = \begin{pmatrix} | & \cdots & | \\ w_1^{[l]} & \cdots & w_{n^{[l]}}^{[l]} \\ | & \cdots & | \end{pmatrix}^T \quad \text{and} \quad b^{[l]} = \begin{pmatrix} b_1^{[l]} & \cdots & b_1^{[l]} \\ \vdots & \ddots & \vdots \\ b_{n^{[l]}}^{[l]} & \cdots & b_{n^{[l]}}^{[l]} \end{pmatrix}$$

respectively. (We will see the dimensions of the matrices below)

We denote the training examples as

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

where m is the total number of training data. Each $x^{(i)}$ is the input and $y^{(i)}$ is the expected output for $1 \leq i \leq m$. If the neural network has $n^{[0]}$ inputs, then each $x^{(i)}$ is a $n^{[0]} \times 1$ vector. Similarly, if the neural network has $n^{[L]}$ outputs then each $y^{(i)}$ is a $n^{[L]} \times 1$ vector.

We can concatenate these $n^{[0]} \times 1$ vectors into the training data matrix:

$$A^{[0]} = \begin{pmatrix} | & \cdots & | \\ x^{(1)} & \cdots & x^{(m)} \\ | & \cdots & | \end{pmatrix}$$

Through this, we can easily batch process an entire training data set using the matrix multiplication

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

and further applying the activation function to get

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

$A^{[1]}$ in this case is of dimension $n^{[1]} \times 1$. If we inductively apply the process, we can write the following:

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

and

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

The dimensions of the matrices are as follows:

Theorem 1.1.4

Let l be a layer of a given neural network. Suppose that there are $n^{[l]}$ nodes in layer l . Let m be the size of the training data. Then the following are true.

- $W^{[l]}$ has dimension $n^{[l]} \times n^{[l-1]}$
- $b^{[l]}$ has dimension $n^{[l]} \times m$
- $Z^{[l]}$ and $A^{[l]}$ both has dimension $n^{[l]} \times m$

Deep learning involves providing a specified neural network (giving it L layers and $n^{[l]}$ nodes for each layer l) some training data and outputs the linear regression matrix $W^{[l]}$ and bias $b^{[l]}$ on each layer l .

In other words, the learning process is as follows. We provide with the machine a set of training data. These training data come in a pair (x_i, y_i) . x_i refers to the input of the machine and y_i refers to the expected output of the machine. The “learning” process involves adjusting the w and b values of each node, so that inputting x_i will give a result very close to y_i (In reality, it will be very hard for the machine to achieve the exact answer y_i). To do this, the “learning” process will be repeated many times and so the values of w and b will be constantly changing.

But how should we measure whether our machine outputs accurate predictions? We need a ruler to for us to evaluate its effectiveness. This is called the loss function.

Definition 1.1.5: Loss Function

The loss function of a neural network calculates the difference between the predicted value of the network and the actual value of in the training data. It is a function which depends on $W^{[l]}$ and $b^{[l]}$ for $1 \leq l \leq L$. It is denoted as

$$\mathcal{L}(W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}, Y)$$

Below is a series of code used to implement the calculation of our predicted value. The following code implements the linear part of the of each node.

Algorithm 1.1.6: Linear Forward

Implements the linear part of a layer's forward propagation.

Arguments:

A – activations from previous layer (or input data): (size of previous layer, number of examples)

W – weights matrix: numpy array of shape (size of current layer, size of previous layer)

b – bias vector, numpy array of shape (size of the current layer, 1)

Returns:

Z – the input of the activation function, also called pre-activation parameter

cache – a python tuple containing "A", "W" and "b" ; stored for computing the backward pass efficiently

```
def linear_forward(A, W, b):
    Z = np.dot(W, A) + b
    cache = (A, W, b)
    return Z, cache
```

The goal is simple. The loss function is a function that depends on w and b value of each node, measuring how different the predicted and actual value are. Therefore, for every iteration of the learning process, we want to use the results of the loss function, to adjust our values of w and b of each node, so that eventually the collection of w and b values for each node will be able to give a prediction value close to the training set for every input data in the training set.

1.2 Binary Classification

Binary classification refers to the fact that the outputs are binary. One example is to use deep learning to identify cats. In this case, the output consists of either 1, which means that the picture is a cat, or 0 if not.

Linear regression is the standard way of finding the values of w and b in statistics. Recall that in linear regression, we calculate our estimation of the output \hat{y} using the following equation:

$$\hat{y} = w \cdot x + b$$

where $w \in M_{1 \times n_{[0]}}(\mathbb{R})$ is some vector that we are trying to improvise, and b is the bias which we are also trying to improvise. Assuming that our outputs are binary, \hat{y} should be a probability

$$\hat{y} = \mathbb{P}(y = 1|x)$$

so that it generates an estimation between 0 and 1.

Unfortunately this does not work using the regular equation for linear regression. Instead, we input the result $z = w \cdot x + b$ further into an activation function so that

1. the output is between 0 and 1 (It becomes a probability)
2. the neural network does not reduce to a simple case of that of 1 node.

To elaborate, notice that if we take a linear function for the activation function, say the identity, by composing $A^{[l]}$, we get

$$A^{[L]} = W^{[L]}W^{[L-1]} \dots W^{[1]}A^{[0]} + B$$

This is clearly just the same expression as $z = w \cdot x + b$ which means that no matter how many layers we define the neural network, the network reduces to a one layer network. Therefore, we can no longer use the standard techniques of linear regression in statistics.

Definition 1.2.1: Activation Function

An activation function of a node is a function that calculates the output of the node using the input data from the previous layer, a choice of weight and a bias.

Common choices of activation functions include the following

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

$$\text{ReLU}(z) = \max\{0, z\}$$

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

We will mostly use ReLU because its derivative is 0 for negative values (as we will see, lost functions are optimized by using gradient descent, which depends on the derivative). However, in the case of binary classification, we may want to use the sigmoid function so that its outputs are between 0 and 1. \tanh are sometimes useful because the function lies between -1 and 1 so that the mean of the data can be closer to 0.

The remainder of the codes of the section will be dedicated to how to program an L layer with binary classification at the last layer. The regression model for a binary classification model is called a logistic regression.

1.3 Random Initialization

W and b have to be initialized in order for the first iteration to be successful. However, if the weights are all set to 0, the layer in which the nodes are in will reduce to just being 1 node due to the fact that they have the same weight, and so what each node does is identical.

The code for random initialization is given below:

Algorithm 1.3.1: Initialize Parameters Deep

Program to initialize all weights and bias.

Arguments:

layer_dims – python array (list) containing the dimensions of each layer in our network

Returns:

parameters – python dictionary containing your parameters "W1", "b1", ..., "WL", "bL":

Wl – weight matrix of shape (layer_dims[l], layer_dims[l-1])

bl – bias vector of shape (layer_dims[l], 1)

```
def initialize_parameters_deep(layer_dims):

    np.random.seed(3)
    parameters = {}
    L = len(layer_dims) # number of layers in the network

    for l in range(1, L):
        parameters["W" + str(l)] = np.random.randn(layer_dims[l], layer_dims[l-1]) * 0.01
        parameters["b" + str(l)] = np.zeros((layer_dims[l], 1))

        assert(parameters["W" + str(l)].shape == (layer_dims[l], layer_dims[l - 1]))
        assert(parameters["b" + str(l)].shape == (layer_dims[l], 1))

    return parameters
```

The factor of 0.01 is there to make sure that the weights are close to 0. As seen in the graph of sigmoid, for large weights, z becomes large and thus the gradient becomes close to 0, making gradient descent (the optimization technique) slow).

Moreover, bias take no effect in random initialization nor are they penalized in the algorithm.

1.4 Forward Propagation

The following codes the function of a single node, using the linear forward module.

Algorithm 1.4.1: Linear Activation Forward

Implements the forward propagation for the LINEAR to ACTIVATION layer.

Arguments:

A_prev – activations from previous layer (or input data): (size of previous layer, number of examples)

W – weights matrix: numpy array of shape (size of current layer, size of previous layer)

b – bias vector, numpy array of shape (size of the current layer, 1)

activation – the activation to be used in this layer, stored as a text string: "sigmoid" or "relu"

Returns:

A – the output of the activation function, also called the post-activation value

cache – a python tuple containing "linear_cache" and "activation_cache"; stored for computing the backward pass efficiently

```
def linear_activation_forward(A_prev, W, b, activation):

    if activation == "sigmoid":
        Z, linear_cache = linear_forward(A_prev, W, b)
        A, activation_cache = sigmoid(Z)

    elif activation == "relu":
        Z, linear_cache = linear_forward(A_prev, W, b)
        A, activation_cache = relu(Z)

    cache = (linear_cache, activation_cache)

    return A, cache
```

Forward propagation refers action of inputting training data and outputting the calculated values using the neural network. The goal is then to compare the calculated values with the expected values

and to try and make sense of the difference.

We then use linear activation forward to calculate the final prediction result using the following module:

Algorithm 1.4.2: L Model Forward

Implements forward propagation for the $[\text{LINEAR} \rightarrow \text{RELU}]^*(L-1) \rightarrow \text{LINEAR} \rightarrow \text{SIGMOID}$ computation.

Arguments:

X – data, numpy array of shape (input size, number of examples)

parameters – output of initialize_parameters_deep()

Returns:

AL – activation value from the output (last) layer

caches – list of caches containing: every cache of linear_activation_forward() (there are L of them, indexed from 0 to $L-1$)

```
def L_model_forward(X, parameters):

    caches = []
    A = X
    L = len(parameters) // 2          # number of layers in the neural network

    # The first n-1 layers using ReLU
    for l in range(1, L):
        A_prev = A
        A, cache = linear_activation_forward(A_prev, parameters["W" + str(l)], parameters["b" + str(l)], "relu")
        caches.append(cache)

    # The output layer using sigmoid
    AL, cache = linear_activation_forward(A, parameters["W" + str(L)], parameters["b" + str(L)], "sigmoid")
    caches.append(cache)

    return AL, caches
```

1.5 Lost and Cost Function

Loss functions are calculated independently from each training set. The square error

$$\mathcal{L} = \frac{1}{2}(\hat{Y} - Y)^2$$

is often not useful in linear regression due to having no vertex. It is common to use the loss function

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

for logistic regression.

It is easy to see that if a training data $y = 1$, then $\mathcal{L} = -\log(\hat{y})$ which means that our optimization should be based on \hat{y} being as large as possible. Since $0 \leq \hat{y} \leq 1$ we must have $\hat{y} = 1$ for best optimization. Similarly, if $y = 0$, $\mathcal{L} = -\log(1 - \hat{y})$ will result in optimization towards $\hat{y} = 0$.

Since every training set contains multiple pairs of input and output, we cannot only consider the loss function of one particular choice of input and output pair. Instead, we need to consider all the loss of EACH input output pair all at once. This is what the cost function does.

Definition 1.5.1: Cost Function

The cost function

$$J(W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}, Y)$$

of a neural network combines the loss function output for all input output pairs so that we only have to minimize one function (the cost function), instead of needing to minimize the loss function for each individual input output pair.

One common choice of cost function in statistics is the average square error:

$$J = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{Y}_i, Y_i)$$

We will use this cost function from now on, unless otherwise specified. The cost function for a single iteration is calculated as below, given the prediction:

Algorithm 1.5.2: Compute Cost

Implements the cost function defined as above.

Arguments:

AL – probability vector corresponding to your label predictions, shape (1, number of examples)

Y – true "label" vector (for example: containing 0 if non-cat, 1 if cat), shape (1, number of examples)

Returns:

cost – cross-entropy cost

```
def compute_cost(AL, Y):
    m = Y.shape[1]
    cost = (-1/m) * np.sum(np.multiply(Y, np.log(AL)) + np.multiply(1-Y, np.log(1-AL)))
    cost = np.squeeze(cost)      # To make sure your cost's shape is what we expect
    return cost
```

1.6 Backwards Propagation

Optimization of the loss function involves finding the minimum of the function. Since J is a convex function, we can use the method of gradient descent.

Definition 1.6.1: Gradient Descent

For a function J with partial variable w , we find the local optima with respect to w using the formula

$$w_{n+1} = w_n - \alpha \frac{\partial J(w)}{\partial w} \Big|_{w=w_n}$$

where $0 \leq \alpha \leq 1$ is called the learning rate. In matrix form this is written as

$$W_{n+1}^{[l]} = W_n^{[l]} - \alpha \frac{\partial J(W^{[l]})}{\partial W^{[l]}} \Big|_{W^{[l]}=W_n^{[l]}}$$

n here indicates the n th iteration.

We will repeatedly obtain new values of W and b for every node and run more iterations to fine tune their values, hence explaining why iteration appears in the definition.

The process of finding $A^{[l]}$ for each layer is called forward propagation for its iterative use of $A^{[l-1]}$, while finding $\frac{\partial J(w)}{\partial w}|_{w=w_n}$ is called backwards propagation for its iterative use of $\frac{\partial J(w)}{\partial w}|_{w=w_{n-1}}$ AFTER computing $A^{[L]}$.

We establish the following notation.

Definition 1.6.2

Suppose there is a neural network with L layers. Suppose there are $n^{[l]}$ nodes on the l th layer. Denote

$$\frac{\partial \mathcal{L}}{\partial W^{[l]}} = \begin{pmatrix} \frac{\partial \mathcal{L}}{\partial w_1^{[l]}} & \cdots & \frac{\partial \mathcal{L}}{\partial w_{n^{[l]}}^{[l]}} \end{pmatrix}^T$$

for the gradient of \mathcal{L} with respect to $w^{[l]}$, where

$$\frac{\partial \mathcal{L}}{\partial w_k^{[l]}} = \begin{pmatrix} \frac{\partial \mathcal{L}}{\partial (w_k^{[l]})_1} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial (w_k^{[l]})_{n^{[l-1]}}} \end{pmatrix} = \nabla_{w_k^{[l]}} \mathcal{L}$$

In particular, the dimension of $\frac{\partial \mathcal{L}}{\partial W^{[l]}}$ is $(n^{[l]}, n^{[l-1]})$. The notation is similar for the gradient of J .

In one layer neural networks, we have the following gradient:

Proposition 1.6.3

For the loss function \mathcal{L} in a neural network with one layer with m training data, the gradient is given by

$$\frac{\partial J}{\partial W} = \frac{1}{m} \frac{\partial \mathcal{L}}{\partial Z} (A^{[0]})^T$$

The gradient for the bias vector is given by

$$\frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^m \left(\frac{\partial \mathcal{L}}{\partial Z} \right)_i$$

where the sum runs through the elements of the $1 \times m$ vector. (Sanity check: $\frac{\partial \mathcal{L}}{\partial Z}$ is a $1 \times m$ matrix, and $\frac{\partial J}{\partial W}$ has the same dimension as W).

In multiple layers, the gradient is given partially by the following equations:

Proposition 1.6.4

For the cost function \mathcal{L} in a neural network with L layers and m training data, the gradients on the l th layer with activation function φ are given by

$$\frac{\partial J}{\partial W^{[l]}} = \frac{1}{m} \frac{\partial \mathcal{L}}{\partial Z^{[l]}} (A^{[l-1]})^T$$

$$\frac{\partial J}{\partial b^{[l]}} = \frac{1}{m} \sum_{i=1}^m \left(\frac{\partial \mathcal{L}}{\partial Z^{[l]}} \right)_i$$

where the sum are taken over the rows so that $\frac{\partial J}{\partial b^{[l]}}$ has dimension $n^{[l]} \times 1$ and

$$\frac{\partial \mathcal{L}}{\partial Z^{[l]}} = \frac{\partial J}{\partial A^{[l]}} * \varphi'(Z^{[l]})$$

with $*$ being elementwise multiplication and

$$\frac{\partial J}{\partial A^{[l-1]}} = (W^{[l]})^T \frac{\partial \mathcal{L}}{\partial Z^{[l]}}$$

(Sanity check: Recall that both $Z^{[l]}$ and $A^{[l]}$ are both $n^{[l]} \times m$ matrices. $\frac{\partial J}{\partial A^{[l-1]}}$ would have dimension $n^{[l]} \times m$ as well)

The following code calculates part of the gradient (without $\frac{\partial J}{\partial z}$) at each iteration. It must be done after forward propagation since it makes use of its intermediate values.

Algorithm 1.6.5

Implements the linear portion of backward propagation for a single layer (layer l).

Arguments:

dZ – Gradient of the cost with respect to the linear output (of current layer l)

cache – tuple of values (A_prev, W, b) coming from the forward propagation in the current layer

Returns:

dA_prev – Gradient of the cost with respect to the activation (of the previous layer $l - 1$), same shape as A_prev

dW – Gradient of the cost with respect to W (current layer l), same shape as W

db – Gradient of the cost with respect to b (current layer l), same shape as b

```
def linear_backward(dZ, cache):
    A_prev, W, b = cache
    m = A_prev.shape[1]

    dW = (1/m)*(np.dot(dZ, A_prev.T))
    db = (1/m)*(np.sum(dZ, axis=1, keepdims=True))
    dA_prev = np.dot(W.T, dZ)

    return dA_prev, dW, db
```

Helper functions:

Algorithm 1.6.6

We can now combine the backward propagation step of each activation function into the linear activation backward module, which computes our required gradient:

Algorithm 1.6.7

Implements the backward propagation for the LINEAR→ACTIVATION layer.

Arguments:

dA – post-activation gradient for current layer l

cache – tuple of values (linear_cache, activation_cache) we store for computing backward

propagation efficiently

activation – the activation to be used in this layer, stored as a text string: "sigmoid" or "relu"

Returns:

dA_prev – Gradient of the cost with respect to the activation (of the previous layer $l - 1$), same shape as A_prev

dW – Gradient of the cost with respect to W (current layer l), same shape as W

db – Gradient of the cost with respect to b (current layer l), same shape as b

```
def linear_activation_backward(dA, cache, activation):

    linear_cache, activation_cache = cache

    if activation == "relu":
        dZ = relu_backward(dA, activation_cache)
        dA_prev, dW, db = linear_backward(dZ, linear_cache)

    elif activation == "sigmoid":
        dZ = sigmoid_backward(dA, activation_cache)
        dA_prev, dW, db = linear_backward(dZ, linear_cache)

    return dA_prev, dW, db
```

Recalling that the last layer uses the sigmoid activation function and the rest using ReLU, we implement the backwards propagation using the above module:

Algorithm 1.6.8: L Model Backwards

Implements the backward propagation for the [LINEAR→RELU] * ($L - 1$) → LINEAR → SIGMOID group.

Arguments:

AL – probability vector, output of the forward propagation (L_model_forward())

Y – true "label" vector (containing 0 if non-cat, 1 if cat)

caches – list of caches containing: every cache of linear_activation_forward() with "relu" (it's caches[1], for l in range($L-1$) i.e $l = 0 \dots L-2$) the cache of linear_activation_forward() with "sigmoid" (it's caches[$L - 1$])

Returns:

grads – A dictionary with the gradients

grads["dA" + str(l)] = ...

grads["dW" + str(l)] = ...

grads["db" + str(l)] = ...

```

def L_model_backward(AL, Y, caches):

    grads = {}
    L = len(caches) # the number of layers
    m = AL.shape[1]
    Y = Y.reshape(AL.shape) # after this line, Y is the same shape as AL

    #Initializing the backpropagation
    dAL = - (np.divide(Y, AL) - np.divide(1 - Y, 1 - AL))

    # Lth layer (SIGMOID -> LINEAR) gradients

    current_cache = caches[-1]
    dA_prev_temp, dW_temp, db_temp = linear_activation_backward(dAL, current_cache, "sigmoid")
    grads["dA" + str(L-1)] = dA_prev_temp
    grads["dW" + str(L)] = dW_temp
    grads["db" + str(L)] = db_temp

    # L-1 layer to 1 (RELU -> LINEAR) gradients
    for l in reversed(range(L-1)):
        current_cache = caches[l]
        dA_prev_temp, dW_temp, db_temp = linear_activation_backward(grads["dA" + str(l+1)], current_cache, "relu")
        grads["dA" + str(l)] = dA_prev_temp
        grads["dW" + str(l+1)] = dW_temp
        grads["db" + str(l+1)] = db_temp

    return grads

```

Finally, we reach the learning step. We upgrade our new $W^{[l]}$ and $b^{[l]}$ on each layer using the gradient descent. In other words, for every iteration, we get closer to the minima. The definition is given at the beginning of the subsection, and the code is given below:

Algorithm 1.6.9: Update Parameters

Updates parameters using gradient descent.

Arguments:

params – python dictionary containing your parameters

grads – python dictionary containing your gradients, output of L_model_backward

Returns:

parameters – python dictionary containing your updated parameters

parameters["W" + str(l)] = ...

parameters["b" + str(l)] = ...

```

def update_parameters(params, grads, learning_rate):

    parameters = copy.deepcopy(params)
    L = len(parameters) // 2 # number of layers in the neural network

    for l in range(L):
        parameters["W" + str(l+1)] = parameters["W" + str(l+1)] - learning_rate*grads["dW" + str(l+1)]
        parameters["b" + str(l+1)] = parameters["b" + str(l+1)] - learning_rate*grads["db" + str(l+1)]

    return parameters

```

Learning rate determines how big the step is per iteration as it goes down to the minima. It is a hyperparameter that may lead to divergence if the rate is too high, or taking long computational time if the rate is too low.

1.7 Putting Everything Together

We have now introduced all players of a neural network, let us recap and assemble the neural network.

Recall that a neural network consists of layers and nodes, as well as the connection between nodes of the previous layer and the current layer. Each node is a function that takes in data from the previous nodes and outputs new data using the following variables:

- The weight of previous nodes
- The bias of previous nodes
- A chosen activation function

To train a neural network consists of the following steps.

1. Collect and input training data which includes the expected output of the network for each datum.
2. Forward Propagation: Compute the actual output of the training data using the network.
3. Backwards Propagation: Optimize the values of the weights and biases of each node in every layer by performing gradient descent on the chosen cost function.

This completes a training cycle. It is often to repeat the above with multiple batches of training data.

2 Improving the Adaptability of the Prediction

Applied machine learning is highly iterative. This is to find out the best set of hyperparameters to train the parameters. In particular, it is important to find out

- The number of layers of the network
- The number of hidden units per layer
- The learning rate of the gradient descent
- The choice of activation functions for each node

through the use of the idea \rightarrow code \rightarrow experiment cycle. Moreover, intuition from one application area usually does not transfer to another. In this section, we try to control the data set as well as fine tuning the change in $W^{[l]}$ and $b^{[l]}$ to improve the prediction output of the learning process.

2.1 Hypervariables

Definition 2.1.1: Hypervariables

The hypervariables of a program refers to variables that controls the output of other variables which further affects the output of the program.

The following are all hypervariables of a neural network.

- The number of layers of the network
- The number of hidden units per layer
- The learning rate of the gradient descent
- The choice of activation functions for each node

They control what W and b will look like for each node. For different values of hypervariables, the optimization of the cost function will be different.

2.2 Data Sets

Suppose we have n training data available. Traditionally, we split the training data into the training set, the dev set and the test set using either the 60%/20%/20% split or the 70%/30% split ignoring the test set. However, as we gather more and more data, it is useful to put in more data into training since we will still have a large amount of data for testing.

While previously data sets maybe of size nearing 10,000, we now have data sets of size 1,000,000. Modern deep learning appliers take the 98%/1%/1% split.

It is important to make sure the distribution of the train and test set to be similar. For example, training the cat picture identifier only using pet owner photos, while testing the program using security camera cat photos may result in the program being highly inaccurate. It is therefore important both the training data and test data come from the same distribution.

2.3 Bias and Variance

High bias is when the program unable to give correct predictions for most of the data set. It is also called underfitting. High variance is when the program is trained too specifically for the training data so that when given other data sets, the program becomes inaccurate. It is also called overfitting.

For example, suppose we are training a program to identify photos between cats and dogs. If the train set error is at 1% while the dev set error is at 10%. This indicates that the program is too specific to the training set, hence the high dev set error. In this case the prediction has high variance. Now if both the

training set error and the dev set error are both 15%, this means that the program in general is just not performing well. Hence this program is high bias.

An important thing to know is that programs can be simultaneously high bias and high variance. This is when the training set error is high all the while dev set error is even higher.

All this is assuming that human error is almost 0%. If the data with manually provide to train is wrong, it is hard for the program to get right. This is also called Bayes' error in some literature.

In general, there are different solutions to whether the prediction is high bias or high variance. In the case that the prediction is high bias, we can use one of the following methods:

- Increasing the size of the network
- Allow the program to train longer (try different hypervariables)
- Test different Optimization Algorithms
- Change the network structure

with increasing the size of the network being almost always efficient. On the other hand, high variance can be solved using the following methods

- Collecting more training data
- Regularization
- Change the network structure

with collecting more training data and regularization being most useful.

One thing to note is the Bias-variance trade off. Usually, different methods while useful to solving one problem, may increase magnitude of the problem. For example, methods that reduced bias may increase variance. In modern day, this problem is not so much apparent. In particular, increasing the network size helps decreasing bias, fixes variance given we perform regularization. Collecting more training data will solve the problem of high variance while fixing the bias.

2.4 Regularization

Recall that in general, the core of deep learning is to optimize the function

$$J(W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(A^{[L](i)}, Y^{(i)})$$

Regularization refers to the fact that we add a regularization term to the cost function.

Definition 2.4.1: The Regularization Term

Suppose that a neural network with 1 layer has n inputs. Regularization refers to the fact that we choose one of the following regularization terms and add it to the cost function:

- L_2 -regularization:

$$\frac{\lambda}{2m} \sum_{k=1}^n w_k^2 = \frac{\lambda}{2m} \|w\|_2^2$$

(most common)

- L_1 -regularization:

$$\frac{\lambda}{m} \sum_{k=1}^n w_k = \frac{\lambda}{m} \|w\|_1$$

where λ is the regularization parameter. It is a hyperparameter.

In a neural network with L layers, L_2 regularization is given by

$$\frac{\lambda}{2m} \sum_{l=1}^L \|W^{[l]}\|_F^2 = \frac{\lambda}{2m} \sum_{l=1}^L \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{[l]})^2$$

The back propagation, depending on the cost function, also changes by adding the term

$$\frac{\lambda}{m} W^{[l]}$$

for layer l which is just the derivative of the L_2 -regularization term. The effect becomes apparent when we see that the gradient descent step becomes

$$W_{n+1}^{[l]} = W_n^{[l]} - \frac{\alpha \lambda}{m} W_n^{[l]} - \alpha \cdot (\text{normal back prop term})$$

Notice here that W_n decreases by a small amount due to the introduction of the regularization term. Moreover, the larger the hyperparameter β , the larger the decrease in W_n becomes. In the gradient descent step, the importance of the weight decreases, and decays over time. This is why regularization is also sometimes called weight decay.

Since $W^{[l]}$ becomes smaller, this means that some of the nodes in layer l becomes less important. Their contribution becomes very small so that the dependency on the some of the nodes becomes less. Recalling that high variance, means that the prediction fits the training data too well. By reducing node dependency, we can penalize nodes that the program become way too dependent on. Thus allowing the prediction to fit less of the training data.

Let us use tanh to illustrate. As λ increases, the gradient descent step changes W in to closer to 0. This means that Z follows W and decreases to close to 0. As seen in the graph of tanh, at near 0 the derivative is almost linear. This means that we are punishing too complicate fittings that are non-linear, which prevents overfitting.

There is also the notion of dropout regularization, in which each node is assigned a probability to be disabled, meaning that their contributions is 0. As such it is also a type of regularization, by decentralizing nodes that are overly dependent from the training. We demonstrate this using an implementation of inverted dropout.

We define a variable which is the probability for a node to be kept in the layer. We then standardize it by dividing the output $A^{[l]}$ with the probability so that the expected value remains the same. Otherwise it will be hard for the program to train the weights and bias. During making predictions at test time, we don't want any dropout since we don't want randomization in testing the data. Intuitively, nodes should not rely on any node in the previous layer. So by randomly disabling the nodes, we make sure that the weight are spread out and any strong dependency is penalized and so is lost.

Studies have shown that dropout has an extremely similar effect to that of L_2 regularization.

Furthermore, we can also vary the probability that keeps the node. In particular, for layers with less node, the probability of keeping the node should be higher so that these already small amount of nodes keep getting trained in a more constant manner. Conversely, layers with more nodes or layers in which overfitting may occur, the probability to keep nodes should be lower. We can also drop out inputs although it is an uncommon practise.

The down side of this is that J is now less well defined since it now randomly depends on different sets of nodes. It becomes hard to debug the cost function itself this way. Optimizing the cost function itself should be seen as a completely different task. (Orthogonalization)

Finally, there are also other regularization methods such as data augmentation, to flip, rotate crop and zoom in or out of images so that more data can be used to train it. Though this is not as good since the images are fundamentally the same, but it is a cheap way to obtain more data. Early stopping is a technique that as we are carrying out the test using the dev set, once we see that the cost function is not minimized or is not minimized the same way as the minimization done in the training set, we halt the process. Early stopping may not be a good way to regularize the program because it contradicts orthogonality. Orthogonality essentially means to do one task at a time / change one variable at a time. We will see more of this later. But early stopping is contrary to orthogonality since we are stopping the optimization of J as well as stopping the prevention of overfitting but halting early. This makes calculating a good λ for regularization harder.

3 Improving the Speed of the Algorithm

3.1 Normalization

Normalization is a process that allows training to be more efficient. This is because we standardize all our data into constant mean and variance each time W and b for each layer is renewed. We do it through the following:

Definition 3.1.1: Normalization

For a neural network at layer l , we normalize the inputs by the following equations:

$$A_{\text{norm}}^{[l]} = \frac{A^{[l]} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

where μ is the mean of $A^{[l]}$ calculating over the training data and σ is the standard deviation. The ϵ term is extremely small to prevent σ from becoming 0 due to rounding.

For $A^{[l]} = \begin{pmatrix} | & & | \\ A^{[l](1)} & \dots & A^{[l](m)} \\ | & & | \end{pmatrix}$, the mean is calculated by

$$\mu = \frac{1}{m} \sum_{k=1}^m A^{[l](k)}$$

and the variance is calculated by

$$\sigma^2 = \frac{1}{m} \sum_{k=1}^m (A^{[l](k)} * A^{[l](k)})$$

where $*$ is element wise multiplication. Finally division in $\frac{A^{[l]} - \mu}{\sigma}$ is given by column wise division of the matrix $A^{[l]} - \mu$ with the column vector σ , where the column divides the column element wise.

This results in the variance in different inputs to be the same, which is 1.

In the case that J is not normalized, the graph of J may be elongated or in weird shapes so that gradient descent becomes difficult. By normalizing J , gradient descent becomes easier, we can take larger steps for each step in the gradient descent. Indeed say x_1 is a value that ranges between 1 and 1000 and x_2 is a value that ranges between 0 and 1. Then W_1 and W_2 becomes very different.

3.2 Weight Initialization

Sometimes derivatives get huge or tiny. In the case that the activation function is linear, the weight matrix becomes diagonal so that inputs grow or shrink exponentially. Initializing the weights may solve this problem. Recall that $Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$. If we want the size of $Z^{[l]}$ to be appropriate, for larger n we need smaller w_{ij} in general. So we initialize $W^{[l]}$ so that each element is between 0 and $\text{Var}(w_{ij}) = \sqrt{\frac{1}{n^{[l-1]}}}$. This allows the elements of the matrix to have a variance of $\frac{1}{n^{[l-1]}}$.

Other weight initialization equations include the Xavier Initialization:

$$\tanh\left(\sqrt{\frac{1}{n^{[l-1]}}}\right)$$

or

$$\sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$$

. This could also be a hypervariable.

3.3 Minibatch

One way to lower the computation time is to process the training data in a minibatch instead of one batch, as matrix multiplication becomes computationally costly for large matrices. We break the training data $A^{[0]}$ into

$$A^{\{1\}[0]}, \dots, A^{\{m/s\}[0]}$$

where s is each size of the minibatch.

3.4 Optimization of the Gradient Descent

3.5 Learning Rates

4 Convolutional Neural Networks

4.1 Matrix Convolutions

Definition 4.1.1: Matrix Convolutions

Definition 4.1.2: Padding

Definition 4.1.3: Strided Convolutions

4.2 Convolutions in Neural Networks

A neural network consists of layers of nodes organized in a single file, and hence we use matrix and vectors to manipulate the data. A convolutional neural network arranges the nodes in 3D and furthermore, changes the size each layer through convolutions, hence the name.

Definition 4.2.1: Convolutional Layers

Let $A^{[l-1]}$ denote the output of the $l - 1$ layer. Suppose that $A^{[l-1]}$ is of size $m \times n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$, where m is the size of the training data. We will consider the case $m = 1$ and generalize it accordingly. A convolutional layer consists of the following data:

- $n_C^{[l]}$ amount of nodes with a weight and a bias and a chosen activation function.
- Each weights are of size $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$, where $f^{[l]}$ is called the filter size
- Each biases is of size $n_H^{[l]} \times n_W^{[l]}$ where

$$n_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} \right\rfloor + 1 \quad \text{and} \quad n_W^{[l]} = \left\lfloor \frac{n_W^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} \right\rfloor + 1$$

(p, s explained below)

Both of which acts as variables of the neural network, such that the output of the layer is calculated as follows.

1. For each node with weight W , apply the convolution

$$A^{[l-1]} * W$$

with stride $s^{[l]}$ and padding size $p^{[l]}$ to obtain $W^{[l]}$ amount of $n_H^{[l]} \times n_W^{[l]}$ matrices.

2. For each node, add in the corresponding biases to the above matrices
3. Apply the new matrices to their respective activation functions.
4. Stack up the matrices to form the output which has dimension $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

For general training sizes m , all matrices have an extra dimension which is of size m .

Definition 4.2.2: Convolutional Neural Networks

Convolutional neural networks is a neural network that consists of convolutional layers.

4.3 Types of Convolutions

5 Deep Neural Networks

6 Recurrent Neural Networks

7 Large Language Models

Large language models refer to large, general purpose language models that can be pre-trained and then fine tuned for specific purposes. They are trained to solve common language problems such as text classification, questions answering, document summarization and text generation. They can then be tailored to solve problems in specific fields of interests.

Large here means that there is a large training data set and a large number of parameters. It is also general purpose in the sense that human language is crucial for a society to function. They are usually pre-trained with a smaller data set. They usually require minimal training data to tailor it to a specific situation such as retail. A recent example of such a large language model is PaLM, created by Google.

With LLM development, there is no need for machine learning expertise nor any training examples. This is because there is no need to train a model and instead people are more focused in thinking about prompt design. This is compared to traditional machine learning development in which machine learning expertise is needed and one needs to train a model and think about computational costs and thinking about minimizing the cost function.

Question answering is a subfield of natural language processing (NLP) that deals with the task of answering questions given in a language naturally. They are trained in a large number of texts and are able to answer a wide range of questions. There is no need for specific domain knowledge.

Prompt design is the process of creating a prompt tailored to the specific task that the system is being asked to perform. Prompt engineering is the process of creating a prompt that is designed to improve performance. The latter may require specific domain knowledge.

There are three main types of large language models. Generic language models predict the next word / token based on the language in the given training data. This involves completing sentences. Instruction tuned models are trained to predict a response to the instructions given in the input. This involves answering full questions. Dialogue tuned models are trained to have a dialogue predicting a response. They are a special case of instruction tuned models. They work best with paraphrasing / humane type responses.

Tuning is the process of adapting a model to a new domain or a set of custom use cases by training the model on new training data. Fine tuning is expensive and not efficient in most cases. Therefore there are more specific types of tuning such as parameter efficient tuning methods (PETM). This is a method for tuning an LLM on your own custom data without duplicating the model. Another type is prompt tuning which is one of the easiest parameter-efficient tuning methods. Generative AI and Vertex AI studio allows one to quickly explore and customize different AI. They come with a library of pre trained models, a tool for fine tuning models, a tool for deploying models and a community forum. Vertex AI also requires little to none coding experiences. It is useful for creating Chat-bots, digital assistants, custom search engines, knowledge bases, training applications and more. PaLM API allows one to test and customize Google created AI models. They are equipped with a graphical user interface. A transformer model consists of an encoding component and a decoding component. The encoding component takes in an input and the decoding component solves the task at hand.

7.1 Data Preprocessing

Large language models take in sentences / paragraphs as the total input. In order to understand the meaning behind each sentence, we must first break down the sentence / paragraph in a meaningful

way. Namely, we break down each paragraphs into sentences and then each sentence into words.

Given an input P that is a string (eg a paragraph or a sentence), we use an array X to store all the words of P . However we do not store the contents of P as strings. Instead, we establish a dictionary, which is an array that stores some amount of vocabulary, indexed by $I \subseteq \mathbb{N}$. For example, one may have

$$\begin{bmatrix} 1 & 2 & \cdots & 10000 \\ a & aaron & \cdots & zulu \end{bmatrix}$$

Suppose that the dictionary has length N . Each word is looked up in the dictionary and is returned the index of the word in the vocabulary. It is then represented as a one-hot vector with position the index. This means that if word i in paragraph P has index k_i in the dictionary, then we obtain a vector with 1 at position k_i and 0 everywhere else. We combine the results into a single matrix X , of size $N \times d$. Each column in X is precisely the one-hot vector for the words.

7.2 Attention

Since each output of a training set is independent of the of the other outputs, we can parallelize the process. In this case, we first impose that the maximum character of each input x_i a value d . If the training set has size m , then the input $X \in M_{m \times d}(\mathbb{R})$.

Transformers in fact are more complicated than self attention. For one, different words in a sentence can relate to each word in multiple ways instead of one. This leads to multihead attention.

In multihead attention, we first choose the number of heads h , and then provide with each head the weight matrices for query W_i^Q , for key W_i^K and for value W_i^V . Individually, we compute the attention for each head

$$\text{head}_i = \text{Selfattention}(Q, K, V)$$

and then we compute the multihead attention A by projecting it down to an appropriate dimension using an orthogonal matrix W^O after we put all the self attention matrices into a direct sum:

$$A = \left(\bigoplus_{i=1}^h \text{head}_i \right) W^O$$

7.3 Transformer Blocks

We will need a few more components to assemble the transformer.

Definition 7.3.1: Feed Forward Layer

Definition 7.3.2: Residual Connection

Residual Connection inputs two sequences of the same dimensions and sum it element wise.

The point of residual connections is that we can give higher level layers direct access to lower level information, without going through the intermediate layers.

Definition 7.3.3: Layer Normalize

Definition 7.3.4: Transformer Blocks

A transformer block consists of a sequence of operations.

- 0) The sequence of inputs are provided to the transformer.
- 1) Multihead Attention Layer
- 2) Residual Connection: The output of 1) with 0)
- 3) Layer Normalize
- 4) Feed Forward Layer
- 5) Residual Connection: The output of 4) with the output of 3)
- 6) Layer Normalize