

Week 1 Assignment — Regularization

1. L2 Regularization

$$J = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(a^{LL}(i)) + (1-y^{(i)}) \log(1-a^{LL}(i)))$$

↓ regularization

$$J_{\text{regularized}} = \underbrace{-\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(a^{LL}(i)) + (1-y^{(i)}) \log(1-a^{LL}(i)))}_{\text{cross-entropy cost}} + \underbrace{\frac{\lambda}{2m} \sum_{i,j} W_{k,j}^{LL2}}_{\text{L2 Regularization cost}}$$

GRADED FUNCTION: compute_cost_with_regularization

```
def compute_cost_with_regularization(A3, Y, parameters, lambd):
```

```
    """
```

```
    Implement the cost function with L2 regularization. See formula (2) above.
```

```
    Arguments:
```

```
    A3 -- post-activation, output of forward propagation, of shape (output size, number of example
```

```
s)
```

```
    Y -- "true" labels vector, of shape (output size, number of examples)
```

```
    parameters -- python dictionary containing parameters of the model
```

```
    Returns:
```

```
    cost - value of the regularized loss function (formula (2))
```

```
    """
```

```
    m = Y.shape[1]
```

```
    W1 = parameters["W1"]
```

```
    W2 = parameters["W2"]
```

```
    W3 = parameters["W3"]
```

```
    cross_entropy_cost = compute_cost(A3, Y) # This gives you the cross-entropy part of the cost
```

```
    ### START CODE HERE ### (approx. 1 line)
```

```
    L2_regularization_cost = (1. / m) * (lambd / 2) * (np.sum(np.square(W1)) + np.sum(np.square(W2))
```

```
+ np.sum(np.square(W3)))
```

```
    cost = cross_entropy_cost + L2_regularization_cost
```

```
    return cost
```

For backward propagation: add the regularization terms gradient

```
### START CODE HERE ### (approx. 1 line)
```

```
dw3 = 1./m * (np.dot(dz3, A2.T) + lambd * W3)
```

```
### END CODE HERE ###
```

```
db3 = 1./m * np.sum(dz3, axis=1, keepdims = True)
```

```
da2 = np.dot(W3.T, dz3)
```

```
dz2 = np.multiply(da2, np.int64(A2 > 0))
```

```
### START CODE HERE ### (approx. 1 line)
```

```
dw2 = 1./m * (np.dot(dz2, A1.T) + lambd * W2)
```

```
### END CODE HERE ###
```

```
db2 = 1./m * np.sum(dz2, axis=1, keepdims = True)
```

```
da1 = np.dot(W2.T, dz2)
```

```
dz1 = np.multiply(da1, np.int64(A1 > 0))
```

```
### START CODE HERE ### (approx. 1 line)
```

```
dw1 = 1./m * (np.dot(dz1, X.T) + lambd * W1)
```

```
### END CODE HERE ###
```

```
db1 = 1./m * np.sum(dz1, axis=1, keepdims = True)
```

$$\frac{d}{dw} \left(\frac{1}{2} \frac{\lambda}{m} W^2 \right) = \frac{\lambda}{m} W$$

What is L2-regularization actually doing?:

L2-regularization relies on the assumption that a model with small weights is simpler than a model with large weights. Thus, by penalizing the square values of the weights in the cost function you drive all the weights to smaller values. It becomes too costly for the cost to have large weights! This leads to a smoother model in which the output changes more slowly as the input changes.

- The cost computation:
 - A regularization term is added to the cost
- The backpropagation function:
 - There are extra terms in the gradients with respect to weight matrices
- Weights end up smaller ("weight decay"):
 - Weights are pushed to smaller values.

2. Dropout

randomly shuts down some neurons in each iteration
(1) forward propagation

1. In lecture, we discussed creating a variable $d^{[1]}$ with the same shape as $a^{[1]}$ using `np.random.rand()` to randomly get numbers between 0 and 1. Here, you will use a vectorized implementation, so create a random matrix $D^{[1]} = [d^{1} d^{[1](2)} \dots d^{[1](m)}]$ of the same dimension as $A^{[1]}$.
2. Set each entry of $D^{[1]}$ to be 0 with probability $(1 - \text{keep_prob})$ or 1 with probability (keep_prob) , by thresholding values in $D^{[1]}$ appropriately. Hint: to set all the entries of a matrix X to 0 (if entry is less than 0.5) or 1 (if entry is more than 0.5) you would do: $X = (X < 0.5)$. Note that 0 and 1 are respectively equivalent to False and True.
3. Set $A^{[1]}$ to $A^{[1]} * D^{[1]}$. (You are shutting down some neurons). You can think of $D^{[1]}$ as a mask, so that when it is multiplied with another matrix, it shuts down some of the values.
4. Divide $A^{[1]}$ by `keep_prob`. By doing this you are assuring that the result of the cost will still have the same expected value as without drop-out. (This technique is also called inverted dropout.)

```
# LINEAR -> RELU -> LINEAR -> RELU -> LINEAR -> SIGMOID
Z1 = np.dot(W1, X) + b1
A1 = relu(Z1)
### START CODE HERE ### (approx. 4 lines)           # Steps 1-4 below correspond to the Steps 1-4 described above.

1. D1 = np.random.rand(A1.shape[0], A1.shape[1])           # Step 1: initialize matrix D1 = np.random.rand(..., ...)
2. D1 = D1 < keep_prob           # Step 2: convert entries of D1 to 0 or 1 (using keep_prob as the threshold)
3. A1 = np.multiply(A1, D1)           # Step 3: shut down some neurons of A1
4. A1 /= keep_prob           # Step 4: scale the value of neurons that haven't been shut down
### END CODE HERE ###
Z2 = np.dot(W2, A1) + b2
A2 = relu(Z2)
### START CODE HERE ### (approx. 4 lines)
D2 = np.random.rand(A2.shape[0], A2.shape[1])           # Step 1: initialize matrix D2 = np.random.rand(..., ...)
D2 = D2 < keep_prob           # Step 2: convert entries of D2 to 0 or 1 (using keep_prob as the threshold)
A2 = np.multiply(A2, D2)           # Step 3: shut down some neurons of A2
A2 /= keep_prob           # Step 4: scale the value of neurons that haven't been shut down
### END CODE HERE ###
Z3 = np.dot(W3, A2) + b3
A3 = sigmoid(Z3)
```

(2) backward propagation

1. You had previously shut down some neurons during forward propagation, by applying a mask $D^{[1]}$ to $A1$. In backpropagation, you will have to shut down the same neurons, by reapplying the same mask $D^{[1]}$ to $da1$.
2. During forward propagation, you had divided $A1$ by `keep_prob`. In backpropagation, you'll therefore have to divide $da1$ by `keep_prob` again (the calculus interpretation is that if $A^{[1]}$ is scaled by `keep_prob`, then its derivative $da^{[1]}$ is also scaled by the same `keep_prob`).


```

dZ3 = A3 - Y
dW3 = 1./m * np.dot(dZ3, A2.T)
db3 = 1./m * np.sum(dZ3, axis=1, keepdims = True)
dA2 = np.dot(W3.T, dZ3)
### START CODE HERE ### (= 2 lines of code)
1. dA2 = np.multiply(dA2, D2) # Step 1: Apply mask D2 to shut down the same neurons
as during the forward propagation
2. dA2 /= keep_prob # Step 2: Scale the value of neurons that haven't been shut down
### END CODE HERE ###
dZ2 = np.multiply(dA2, np.int64(A2 > 0))
dW2 = 1./m * np.dot(dZ2, A1.T)
db2 = 1./m * np.sum(dZ2, axis=1, keepdims = True)

dA1 = np.dot(W2.T, dZ2)
### START CODE HERE ### (= 2 lines of code)
dA1 = np.multiply(dA1, D1) # Step 1: Apply mask D1 to shut down the same neurons a
s during the forward propagation
dA1 /= keep_prob # Step 2: Scale the value of neurons that haven't been shut down
### END CODE HERE ###
dZ1 = np.multiply(dA1, np.int64(A1 > 0))
dW1 = 1./m * np.dot(dZ1, X.T)
db1 = 1./m * np.sum(dZ1, axis=1, keepdims = True)

```

- Dropout is a regularization technique.
- You only use dropout during training. Don't use dropout (randomly eliminate nodes) during test time.
- Apply dropout both during forward and backward propagation.
- During training time, divide each dropout layer by keep_prob to keep the same expected value for the activations. For example, if keep_prob is 0.5, then we will on average shut down half the nodes, so the output will be scaled by 0.5 since only the remaining half are contributing to the solution. Dividing by 0.5 is equivalent to multiplying by 2. Hence, the output now has the same expected value. You can check that this works even when keep_prob is other values than 0.5.