## Aditya N Bhatt

231057015

DL Assignment-1

05-01-2024

github link: https://github.com/adityab24840/Deep-Learning/blob/Assignments/Assignments/Softmax_Classifier.ipynb

```
1 ## Load libraries
2 import numpy as np
3 import sys
4 import matplotlib.pyplot as plt
5 import matplotlib.cm as cm
6 plt.style.use('dark_background')
7 %matplotlib inline
```

```
1 import tensorflow as tf
```

```
1 tf.__version__
```

```
'2.15.0'
```

```
1 # Generate artificial data with 5 samples, 4 features per sample
2 # and 3 output classes
3 num_samples = 5 # number of samples
4 num_features = 4 # number of features (a.k.a. dimensionality)
5 num_labels = 3 # number of output labels
6 # Data matrix (each column = single sample)
7 X = np.random.choice(np.arange(3, 10), size = (num_features, num_samples), replace = True)
8 # Class labels
9 y = np.random.choice([0, 1, 2], size = num_samples, replace = True)
10 print(X)
11 print('------')
12 print(y)
13 print('------')
14 # One-hot encode class labels
15 y = tf.keras.utils.to_categorical(y)
16 print(y)
```

```
[[5 7 7 7 9]
 [6 9 3 5 3]
 [6 4 7 7 4]
 [8 6 6 3 3]]
------
[2 2 0 2 0]
------
[[0. 0. 1.]
 [0. 0. 1.]
 [1. 0. 0.]
 [0. 0. 1.]
 [1. 0. 0.]]
```

A generic layer class with forward and backward methods

```
1   class Layer:
2     def __init__(self):
3       self.input = None
4       self.output = None
5
6     def forward(self, input):
7       pass
8
9     def backward(self, output_gradient, learning_rate):
10       pass
```

The softmax classifier steps for a generic sample $\mathbf{x}$ with (one-hot encoded) true label $\mathbf{y}$ (3 possible categories) using a randomly initialized weights matrix (with bias abosrbed as its last last column):

1. Calculate raw scores vector for a generic sample $\mathbf{x}$ (bias feature added):

$$\mathbf{z} = \mathbf{W}\mathbf{x}.$$

2. Calculate softmax probabilities (that is, softmax-activate the raw scores)

$$\mathbf{a} = \mathrm{softmax}(\mathbf{z}) \Rightarrow \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \mathrm{softmax}\left(\begin{bmatrix} z_0 \\ z_1 \\ z_2 \end{bmatrix}\right) = \begin{bmatrix} \frac{e^{z_0}}{e^{z_0}+e^{z_1}+e^{z_2}} \\ \frac{e^{z_1}}{e^{z_0}+e^{z_1}+e^{z_2}} \\ \frac{e^{z_2}}{e^{z_0}+e^{z_1}+e^{z_2}} \end{bmatrix}$$

3. Softmax loss for this sample is (where output label $y$ is not yet one-hot encoded)

$$L = -\log([a]_y)$$
$$= -\log\Big([\mathrm{softmax}(\mathbf{z})]_y\Big)$$
$$= -\log\Big([\mathrm{softmax}(\mathbf{Wx})]_y\Big).$$

4. Predicted probability vector that the sample belongs to each one of the output categories is given a new name

$$\hat{\mathbf{y}} = \mathbf{a}.$$

5. One-hot encoding the output label

$$\underbrace{y \to \mathbf{y}}$$
$$\text{e.g. } 2 \to \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

results in the following representation for the softmax loss for the sample which is also referred to as the categorical crossentropy (CCE) loss:

$$L = L\left(\mathbf{y}, \hat{\mathbf{y}}\right) = \sum_{k=0}^{2} -y_k \log(\hat{y}_k).$$

6. Calculate the gradient of the loss for the sample w.r.t. weights by following the computation graph from top to bottom (that is, backward):

$$L$$
$$\downarrow$$
$$\hat{\mathbf{y}} = \mathbf{a}$$
$$\downarrow$$
$$\mathbf{z}$$
$$\downarrow$$
$$\mathbf{W}$$

$$\Rightarrow \nabla_{\mathbf{W}}(L) = \nabla_{\mathbf{W}}(\mathbf{z}) \times \nabla_{\mathbf{z}}(\mathbf{a}) \times \nabla_{\mathbf{a}}(L)$$
$$= \underbrace{\nabla_{\mathbf{W}}(\mathbf{z})}_{\text{first term}} \times \underbrace{\nabla_{\mathbf{z}}(\mathbf{a})}_{\text{second to last term}} \times \underbrace{\nabla_{\hat{\mathbf{y}}}(L)}_{\text{last term}}.$$

7. Now focus on the last term $\nabla_{\hat{\mathbf{y}}}(L)$:

$$\nabla_{\hat{\mathbf{y}}}(L) = \begin{bmatrix} \nabla_{\hat{y}_0}(L) \\ \nabla_{\hat{y}_1}(L) \\ \nabla_{\hat{y}_2}(L) \end{bmatrix} = \begin{bmatrix} -y_0/\hat{y}_0 \\ -y_1/\hat{y}_2 \\ -y_0/\hat{y}_2 . \end{bmatrix}$$

8. Now focus on the second to last term $\nabla_{\mathbf{z}}(\mathbf{a})$:

$$\nabla_{\mathbf{z}}(\mathbf{a}) = \nabla_{\mathbf{z}}\left(\begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix}\right)$$
$$= \begin{bmatrix} \nabla_{\mathbf{z}}(a_0) & \nabla_{\mathbf{z}}(a_1) & \nabla_{\mathbf{z}}(a_2) \end{bmatrix}$$
$$= \begin{bmatrix} \nabla_{z_0}(a_0) & \nabla_{z_0}(a_1) & \nabla_{z_0}(a_2) \\ \nabla_{z_1}(a_0) & \nabla_{z_1}(a_1) & \nabla_{z_1}(a_2) \\ \nabla_{z_2}(a_0) & \nabla_{z_2}(a_1) & \nabla_{z_2}(a_2) \end{bmatrix}$$
$$= \begin{bmatrix} a_0(1-a_0) & -a_1 a_0 & -a_2 a_0 \\ -a_0 a_1 & a_1(1-a_1) & -a_1 a_1 \\ -a_0 a_2 & -a_1 a_2 & a_2(1-a_2) \end{bmatrix}.$$

9. On Monday, we will focus on the first term to complete the gradient calculation using the computation graph.

```
1 ## Softmax activation class
2 class Softmax(Layer):
3   def forward(self, input):
4     self.output = np.array(tf.nn.softmax(input))
5
6   def backward(self, output_gradient, learning_rate):
7     return(np.dot((np.identity(np.size(self.output))-self.output.T) * self.output, output_gradient))
```

```
1 ## Define the loss function and its gradient
2 def cce(y, yhat):
3   return(-np.sum(y*np.log(yhat)))
4
5 def cce_gradient(y, yhat):
6   return(-y/yhat)
7
8 # TensorFlow in-built function for categorical crossentropy loss
9 cce = tf.keras.losses.CategoricalCrossentropy()
10 cce
```

```
<keras.src.losses.CategoricalCrossentropy at 0x7ac0abf573a0>
```

```
1 ## Train the 0-layer neural network using batch training with batch size = 1
2
3 # Steps: run over each sample, calculate loss, gradient of loss,
4 # and update weights.
5
6 # Step-1: add the bias feature to all the samples
7 # Step-2: initialize the entries of the weights matrix randomly
8 # Step-3: create softmax layer object softmax
9
10 # Step-4: run over each sample
11 for i in range(X.shape[1]):
12   # Step-5: forward step
13   # (a) Raw scores z = Wx = np.dot(W, x[:, i])
14   # (b) Softmax activation: softmax.forward(z)
15   # (c) Calculate cce loss for sample: cce(y[i, :], softmax.output)
16   # (d) Print cce loss
17
18   # Step-6: backward step
19   # (a) Calculate the gradient of the sample loss w.r.t. input of the
20   # softmax layer: softmax.backward(output_gradient = cce_gradient(y[i, :], softmax.output))
21   # (d) Print gradient
```

```
  File "<ipython-input-8-7ceb0d13911d>", line 21
    # (d) Print gradient
                        ^
SyntaxError: incomplete input
```

SUGGEST FIX

1. Calculate raw scores vector for a generic sample $\mathbf{x}$ (bias feature added):
$$\mathbf{z} = \mathbf{W}\mathbf{x}.$$

```
1   x = tf.constant([[1, 2, 3], [4, 5, 6]])
```

```
1   W = tf.Variable(tf.random.normal([3, 2]))
```

```
1   import numpy as np
2   z = np.dot(W, x)
3   z
```

```
array([[ -9.08701134, -11.99186206, -14.89671278],
       [ -6.39343089,  -7.85449612,  -9.31556135],
       [  4.78612366,   6.11788398,   7.4496443 ]])
```

2. Calculate softmax probabilities (that is, softmax-activate the raw scores)

$$\mathbf{a} = \mathrm{softmax}(\mathbf{z}) \Rightarrow \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \mathrm{softmax}\left( \begin{bmatrix} z_0 \\ z_1 \\ z_2 \end{bmatrix} \right) = \begin{bmatrix} \frac{e^{z_0}}{e^{z_0}+e^{z_1}+e^{z_2}} \\ \frac{e^{z_1}}{e^{z_0}+e^{z_1}+e^{z_2}} \\ \frac{e^{z_2}}{e^{z_0}+e^{z_1}+e^{z_2}} \end{bmatrix}$$

```
1 z = tf.matmul(W, tf.cast(x, tf.float32))
2 z
```

```
<tf.Tensor: shape=(3, 3), dtype=float32, numpy=
array([[ -9.087011 , -11.991862 , -14.896712 ],
       [ -6.3934307,  -7.8544965,  -9.315561 ],
       [  4.7861238,   6.117884 ,   7.449644 ]], dtype=float32)>
```

```
1 softmax = lambda z: np.exp(z) / np.sum(np.exp(z))
2 a = softmax(z)
3 a
```

```
array([[4.9335377e-08, 2.7014548e-09, 1.4792358e-10],
       [7.2939986e-07, 1.6921265e-07, 3.9255497e-08],
       [5.2261811e-02, 1.9795233e-01, 7.4978489e-01]], dtype=float32)
```

3. Softmax loss for this sample is (where output label $y$ is not yet one-hot encoded)

$$L = -\log([a]_y)$$
$$= -\log\Big([\text{softmax}(\mathbf{z})]_y\Big)$$
$$= -\log\Big([\text{softmax}(\mathbf{Wx})]_y\Big).$$

```
1 def softmax_loss(y, a):
2     return -np.log(a[y])
3
4 y = 1
5 a = softmax(z)
6 loss = softmax_loss(y, a)
7 print(loss)
8
```

```
[14.131043 15.59211  17.053175]
```

4. Predicted probability vector that the sample belongs to each one of the output categories is given a new name

$$\hat{\mathbf{y}} = \mathbf{a}.$$

```
1 y_hat = a
2 y_hat
```

```
array([[4.9335377e-08, 2.7014548e-09, 1.4792358e-10],
       [7.2939986e-07, 1.6921265e-07, 3.9255497e-08],
       [5.2261811e-02, 1.9795233e-01, 7.4978489e-01]], dtype=float32)
```

5. One-hot encoding the output label

$$\underbrace{y \to \mathbf{y}}$$
$$\text{e.g. } 2 \to \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

results in the following representation for the softmax loss for the sample which is also referred to as the categorical crossentropy (CCE) loss:

$$L = L\left(\mathbf{y}, \hat{\mathbf{y}}\right) = \sum_{k=0}^{2} -y_k \log(\hat{y}_k).$$

```
1 y_encoded = tf.keras.utils.to_categorical(y, num_labels)
2 y_encoded
```

```
array([0., 1., 0.], dtype=float32)
```

6. Calculate the gradient of the loss for the sample w.r.t. weights by following the computation graph from top to bottom (that is, backward):

$$L$$
$$\downarrow$$
$$\hat{\mathbf{y}} = \mathbf{a}$$
$$\downarrow$$
$$\mathbf{z}$$
$$\downarrow$$
$$\mathbf{W}$$
$$\Rightarrow \nabla_{\mathbf{W}}(L) = \nabla_{\mathbf{W}}(\mathbf{z}) \times \nabla_{\mathbf{z}}(\mathbf{a}) \times \nabla_{\mathbf{a}}(L)$$
$$= \underbrace{\nabla_{\mathbf{W}}(\mathbf{z})}_{\text{first term}} \times \underbrace{\nabla_{\mathbf{z}}(\mathbf{a})}_{\text{second to last term}} \times \underbrace{\nabla_{\hat{\mathbf{y}}}(L)}_{\text{last term}}.$$

```
1 import numpy as np
2
3 def softmax(z):
4     exp_z = np.exp(z)
5     return exp_z / np.sum(exp_z)
6
7 def softmax_loss(y, a):
8     return -np.log(a[y])
9
10 def gradient_W(X, y, a):
11     # Compute the gradient of the loss with respect to z
12     grad_z = a
13     grad_z[y] -= 1
14
15     # Compute the gradient of z with respect to W
16     grad_W = np.outer(grad_z, X)
17
18     return grad_W
19
20 x_0 = np.random.rand()
21 x_1 = np.random.rand()
22 x_2 = np.random.rand()
23 x_3 = np.random.rand()
24
25 X = np.array([x_0, x_1, x_2, x_3])
26
27 print(X)
28 print('------')
29 grad_W = gradient_W(X, y, a)
30 print('------')
31 print(grad_W)
```

```
    [0.95981135 0.52359943 0.68329462 0.59092701]
    ------
    ------
    [[ 4.73526547e-08  2.58319749e-08  3.37105973e-08  2.91536066e-08]
     [ 2.59288703e-09  1.41448020e-09  1.84588955e-09  1.59636263e-09]
     [ 1.41978728e-10  7.74526996e-11  1.01075383e-10  8.74120365e-11]
     [-9.59810668e-01 -5.23599052e-01 -6.83294129e-01 -5.90926587e-01]
     [-9.59811182e-01 -5.23599332e-01 -6.83294495e-01 -5.90926904e-01]
     [-9.59811297e-01 -5.23599395e-01 -6.83294577e-01 -5.90926974e-01]
     [ 5.01614793e-02  2.73642541e-02  3.57102140e-02  3.08829155e-02]
     [ 1.89996894e-01  1.03647726e-01  1.35259762e-01  1.16975378e-01]
     [ 7.19652048e-01  3.92586936e-01  5.12323977e-01  4.43068141e-01]]
```

7. Now focus on the last term $\nabla_{\hat{\mathbf{y}}}(L)$:

$$\nabla_{\hat{\mathbf{y}}}(L) = \begin{bmatrix} \nabla_{\hat{y}_0}(L) \\ \nabla_{\hat{y}_1}(L) \\ \nabla_{\hat{y}_2}(L) \end{bmatrix} = \begin{bmatrix} -y_0/\hat{y}_0 \\ -y_1/\hat{y}_2 \\ -y_0/\hat{y}_2 \cdot \end{bmatrix}$$

```
1 def gradient_y(y, y_hat):
2     grad_y = -y / y_hat
3     return grad_y
4
5 grad_y = gradient_y(y_encoded, y_hat)
6 grad_y
7
```

```
    array([[-0.0000000e+00, -3.7017091e+08, -0.0000000e+00],
           [ 0.0000000e+00,  1.0000002e+00,  0.0000000e+00],
           [-0.0000000e+00, -5.0517211e+00, -0.0000000e+00]], dtype=float32)
```

8. Now focus on the second to last term $\nabla_{\mathbf{z}}(\mathbf{a})$:

$$\begin{aligned}
\nabla_{\mathbf{z}}(\mathbf{a}) &= \nabla_{\mathbf{z}}\left(\begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix}\right) \\
&= \begin{bmatrix} \nabla_{\mathbf{z}}(a_0) & \nabla_{\mathbf{z}}(a_1) & \nabla_{\mathbf{z}}(a_2) \end{bmatrix} \\
&= \begin{bmatrix} \nabla_{z_0}(a_0) & \nabla_{z_0}(a_1) & \nabla_{z_0}(a_2) \\ \nabla_{z_1}(a_0) & \nabla_{z_1}(a_1) & \nabla_{z_1}(a_2) \\ \nabla_{z_2}(a_0) & \nabla_{z_2}(a_1) & \nabla_{z_2}(a_2) \end{bmatrix} \\
&= \begin{bmatrix} a_0(1-a_0) & -a_1 a_0 & -a_2 a_0 \\ -a_0 a_1 & a_1(1-a_1) & -a_1 a_1 \\ -a_0 a_2 & -a_1 a_2 & a_2(1-a_2) \end{bmatrix}.
\end{aligned}$$

```
1 import numpy as np
2
3 def gradient_z(a):
4     diag_a = np.diag(a)
5     return diag_a - np.outer(a, a)
6
7 a_0 = np.random.rand()
8 a_1 = np.random.rand()
9 a_2 = np.random.rand()
10
11 a = np.array([a_0, a_1, a_2])
12 grad_z = gradient_z(a)
13 grad_z
14
```

```
    array([[ 0.15138129, -0.04468256, -0.09806931],
           [-0.04468256,  0.18254319, -0.12671096],
           [-0.09806931, -0.12671096,  0.2492516 ]])
```

Assignment 1

```
1 ## Step-1: add the bias feature to all the samples
2
3 X = np.reshape(X, (X.shape[0], 1))
4 X_bias = np.hstack((np.ones((X.shape[0], 1)), X))
5
```

```
1 import numpy as np
2
3 # Step-2: initialize the entries of the weights matrix randomly
4 initialize_weights = lambda num_features, num_labels: np.random.randn(num_labels, num_features)
5
6 num_features = 4
7 num_labels = 3
8 W = initialize_weights(num_features, num_labels)
9 print(W)
10
```

```
    [[-0.13677493 -0.89663344 -0.4914765  -0.30098478]
     [ 1.99771577 -0.32930935 -1.2273706   0.53561762]
     [ 0.40719042  0.71531368  1.0341364  -1.16897045]]
```

```
1   # Step-3: create softmax layer object softmax
2   class SoftmaxLayer:
3       def __init__(self):
4           pass
5
6       def forward(self, z):
7           exp_z = np.exp(z)
8           self.output = exp_z / np.sum(exp_z)
9
10      def backward(self, y):
11          return self.output - y
12
13  softmax = SoftmaxLayer()
14
```

```
1   # Step-4: run over each sample
2
3   def cross_entropy_loss(y_true, y_pred):
4     return -np.sum(y_true * np.log(y_pred))
5
6   # Start the loop over each sample
7   for i in range(X.shape[1]):
8     # Calculate the raw scores by multiplying the weights with the sample
9     z = np.dot(W, X[:, i])
10
11    # Apply the softmax activation function
12    softmax.forward(z)
13
14    # Calculate the cross-entropy loss for the sample
15    def cross_entropy_loss(y_true, y_pred):
16      return -np.sum(y_true * np.log(y_pred))
17    # Print the loss for the sample
18    print(f"Loss for sample {i}: {loss}")
19
20    # Calculate the gradient of the loss with respect to the input of the softmax layer
21    grad_z = softmax.backward(y)
22
23    # Print the gradient for the sample
```

```
23     # Print the gradient for the sample
24     print(f"Gradient for sample {i}: {grad_z}")
```

```
Loss for sample 0: [14.131043 15.59211  17.053175]
Gradient for sample 0: [-0.94447549 -0.42514257 -0.63038194]
```