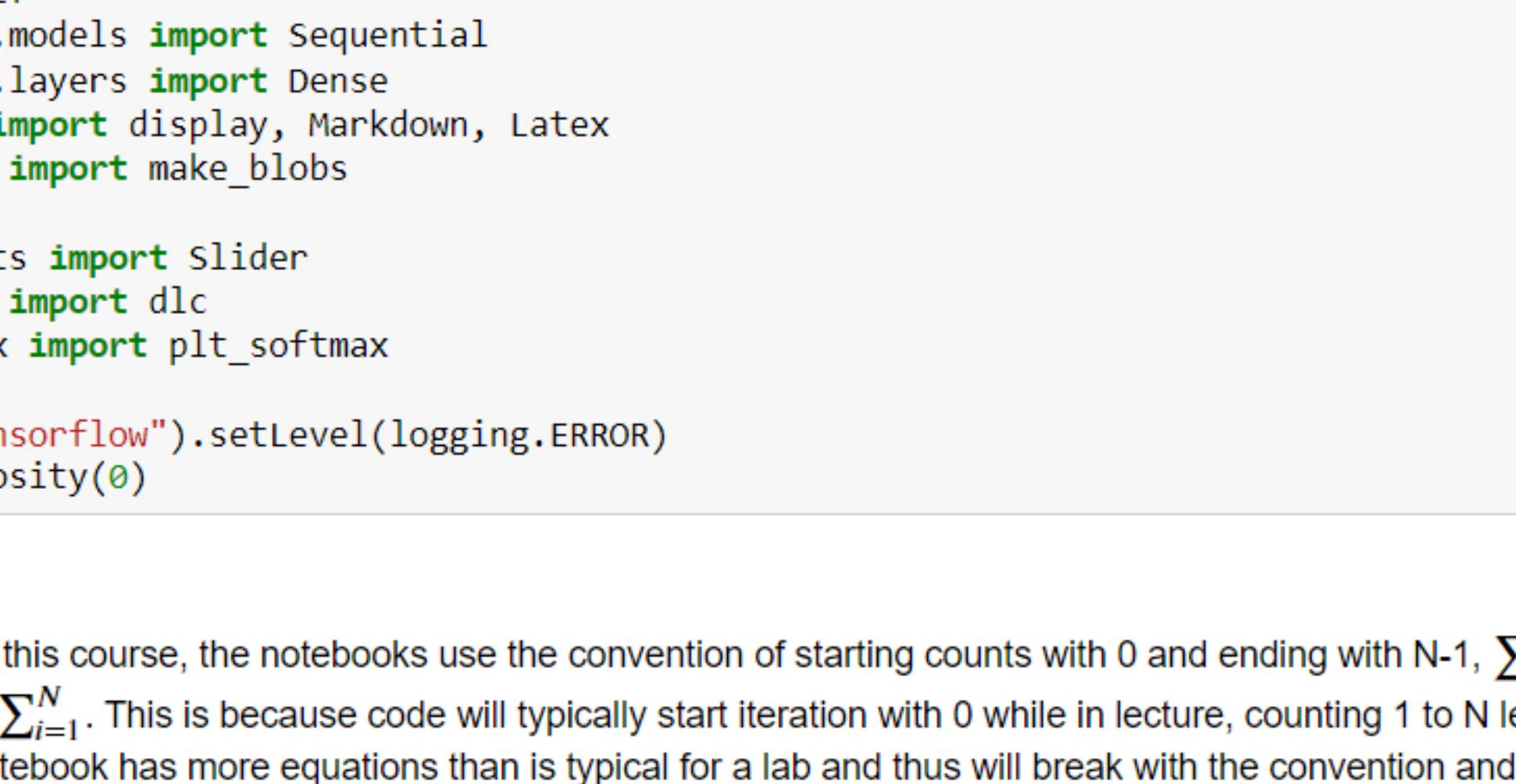


Optional Lab - Softmax Function

In this lab, we will explore the softmax function. This function is used in both Softmax Regression and in Neural Networks when solving Multiclass Classification problems.

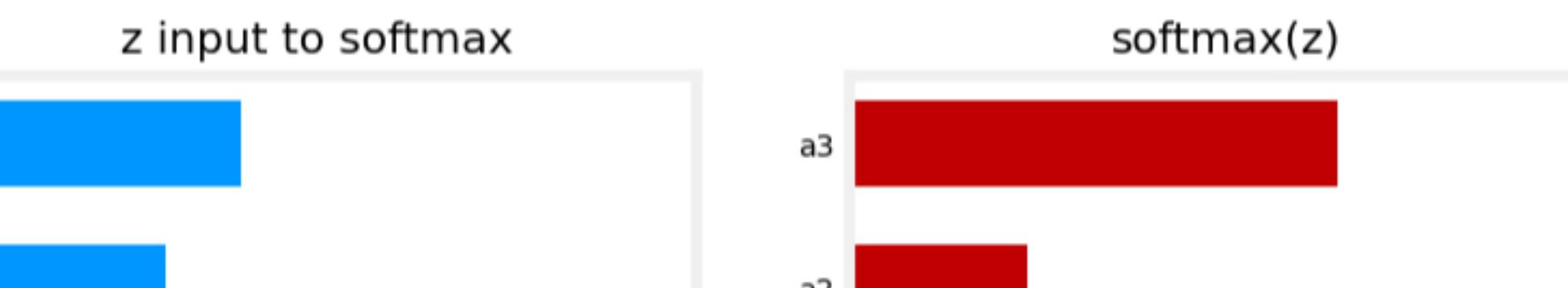


```
In [1]: import numpy as np
import matplotlib.pyplot as plt
plt.style.use('/deeplearning/mplstyle')
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
from IPython.display import display, Markdown, Latex
from sklearn.datasets import make_blobs
import matplotlib.widgets
from ipywidgets import Slider
from lab_utils_common import dlc
from lab_utils_softmax import softmax
import logging
logging.getLogger("tensorflow").setLevel(logging.ERROR)
tf.autograph.set_verbosity(0)
```

Note: Normally, in this course, the notebooks use the convention of starting counts with 0 and ending with $N-1$, while lectures start with 1 and end with N . This is because code will typically start iteration with 0 while in lecture, counting 1 to N leads to cleaner, more succinct equations. This notebook has more equations than is typical for a lab and thus will break with the convention and will count 1 to N .

Softmax Function

In both softmax regression and neural networks with Softmax outputs, N outputs are generated and one output is selected as the predicted category. In both cases a vector \mathbf{z} is generated by a linear function which is applied to a softmax function. The softmax function converts \mathbf{z} into a probability distribution as described below. After applying softmax, each output will be between 0 and 1 and the outputs will add to 1, so that they can be interpreted as probabilities. The larger inputs will correspond to larger output probabilities.



The softmax function can be written:

$$a_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}} \quad (1)$$

The output \mathbf{a} is a vector of length N , so for softmax regression, you could also write:

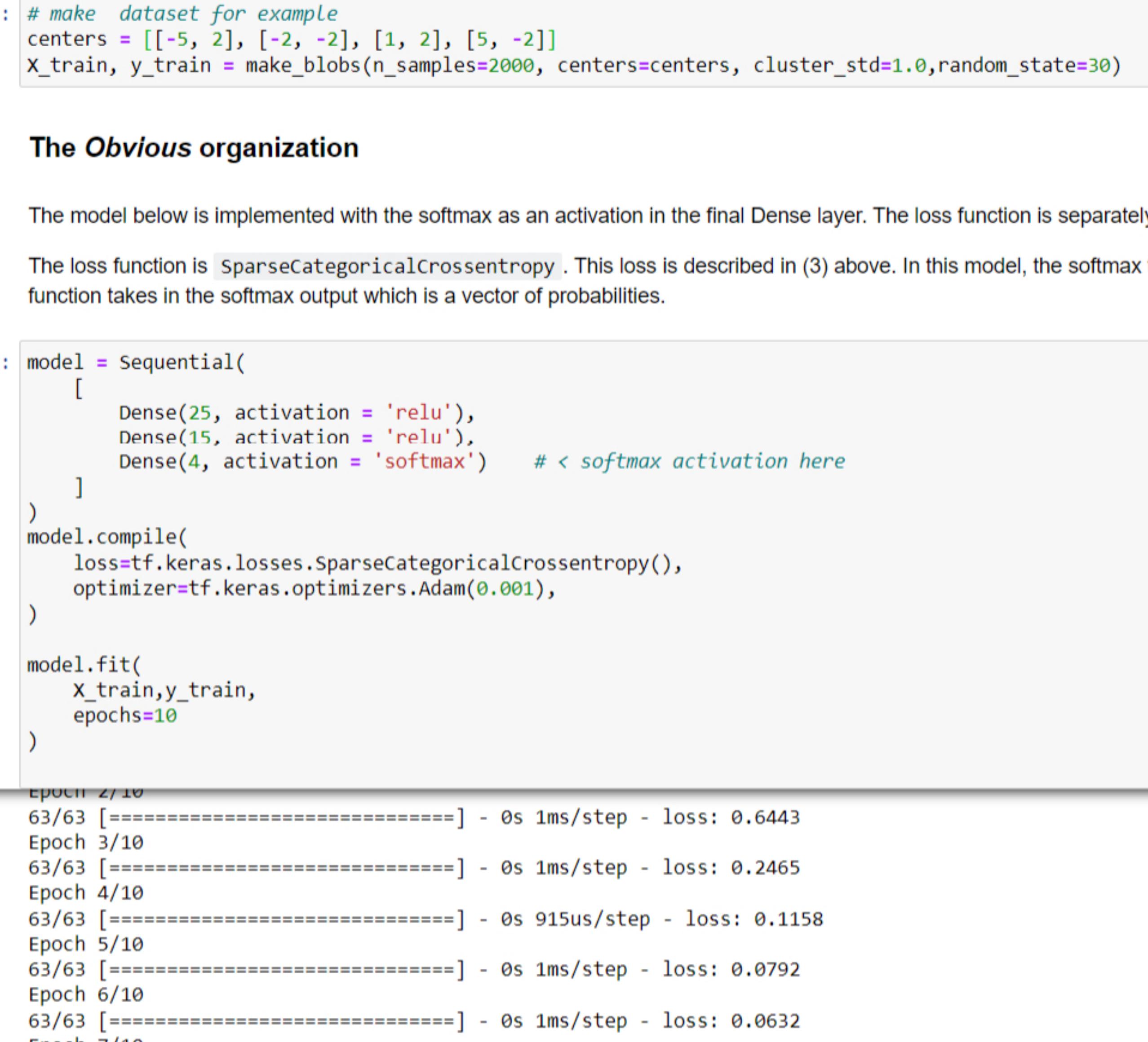
$$\mathbf{a}(x) = \begin{bmatrix} P(y=1|\mathbf{x}; \mathbf{w}, b) \\ \vdots \\ P(y=N|\mathbf{x}; \mathbf{w}, b) \end{bmatrix} = \frac{1}{\sum_{k=1}^N e^{z_k}} \begin{bmatrix} e^{z_1} \\ \vdots \\ e^{z_N} \end{bmatrix} \quad (2)$$

Which shows the output is a vector of probabilities. The first entry is the probability the input is the first category given the input \mathbf{x} and parameters \mathbf{w} and b . Let's create a NumPy implementation:

```
In [2]: def my_softmax(z):
    ez = np.exp(z)           #element-wise exponential
    sm = ez/ez.sum()
    return(sm)
```

Below, vary the values of the \mathbf{z} inputs using the sliders.

```
In [3]: plt.close('all')
plt_softmax(my_softmax)
```



As you are varying the values of the \mathbf{z} 's above, there are a few things to note:

- the exponential in the numerator of the softmax magnifies small differences in the values
- the output values sum to one
- the softmax spans all of the outputs. A change in z_0 for example will change the values of $a_0 - a_3$. Compare this to other activations such as ReLU or Sigmoid which have a single input and single output.

Cost

Logistic regression Softmax regression

$$\text{Logistic regression: } a_i = g(x) = \frac{1}{1+e^{-x}} = P(y=1|x) \quad a_i = \frac{e^{z_i}}{e^{z_1} + \dots + e^{z_N}} = P(y=1|\mathbf{x})$$

$$a_{\bar{y}} = 1 - a_y = P(y=0|x) \quad a_{\bar{y}} = \frac{e^{z_{\bar{y}}}}{e^{z_1} + \dots + e^{z_N}} = P(y=0|\mathbf{x})$$

$$\text{Loss: } L(\mathbf{y}, \mathbf{a}) = -\log(a_y) + (1 - y)\log(a_{\bar{y}}) \quad L(\mathbf{y}, \mathbf{a}) = -\log(a_y) + (1 - y)\log(a_{\bar{y}})$$

$$J(\mathbf{w}, b) = \text{average loss} \quad J(\mathbf{w}, b) = \text{average loss}$$

Figure 1 was created using the softmax function from the lab_utils_softmax module.

The loss function associated with Softmax, the cross-entropy loss, is:

$$L(\mathbf{a}, \mathbf{y}) = \begin{cases} -\log(a_i), & \text{if } y = i. \\ \vdots \\ -\log(a_N), & \text{if } y = N \end{cases} \quad (3)$$

Where y is the target category for this example and \mathbf{a} is the output of a softmax function. In particular, the values in \mathbf{a} are probabilities that sum to one.

Recall: In this course, Loss is for one example while Cost covers all examples.

Note in (3) above, only the line that corresponds to the target contributes to the loss, other lines are zero. To write the cost equation we need an 'indicator function' that will be 1 when the index matches the target and zero otherwise.

$$1_{\{y == n\}} = \begin{cases} 1, & \text{if } y == n. \\ 0, & \text{otherwise.} \end{cases}$$

Now the cost is:

$$J(\mathbf{w}, b) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{j=1}^N 1_{\{j^{(i)} == j\}} \log \frac{e^{z_j^{(i)}}}{\sum_{k=1}^N e^{z_k^{(i)}}} \right] \quad (4)$$

Where m is the number of examples, N is the number of outputs. This is the average of all the losses.

Tensorflow

This lab will discuss two ways of implementing the softmax, cross-entropy loss in Tensorflow, the 'obvious' method and the 'preferred' method. The former is the most straightforward while the latter is more numerically stable.

Let's start by creating a dataset to train a multiclass classification model.

```
In [4]: # make dataset for example
centers = [[-5, 2], [-2, -2], [1, 2], [5, -2]]
X_train, y_train = make_blobs(n_samples=2000, centers=centers, cluster_std=1.0, random_state=30)
```

The Obvious organization

The model below is implemented with the softmax as an activation in the final Dense layer. The loss function is separately specified in the `compile` directive.

The loss function is `SparseCategoricalCrossentropy`. This loss is described in (3) above. In this model, the softmax takes place in the last layer. The loss function takes in the softmax output which is a vector of probabilities.

```
In [5]: model = Sequential([
    Dense(25, activation = 'relu'),
    Dense(15, activation = 'relu'),
    Dense(4, activation = 'softmax')    # < softmax activation here
])
model.compile(
    loss=tf.keras.losses.SparseCategoricalCrossentropy(),
    optimizer=tf.keras.optimizers.Adam(0.001),
)
model.fit(
    X_train,y_train,
    epochs=10
)
```

```
Epoch 0/10
63/63 [=====] - 0s 1ms/step - loss: 0.6443
Epoch 3/10
63/63 [=====] - 0s 1ms/step - loss: 0.2465
Epoch 4/10
63/63 [=====] - 0s 915us/step - loss: 0.1158
Epoch 5/10
63/63 [=====] - 0s 1ms/step - loss: 0.0792
Epoch 6/10
63/63 [=====] - 0s 1ms/step - loss: 0.0632
Epoch 7/10
63/63 [=====] - 0s 1ms/step - loss: 0.0543
Epoch 8/10
63/63 [=====] - 0s 932us/step - loss: 0.0482
Epoch 9/10
63/63 [=====] - 0s 1ms/step - loss: 0.0425
Epoch 10/10
63/63 [=====] - 0s 1ms/step - loss: 0.0385
63/63 [=====] - 0s 1ms/step - loss: 0.0385
```

```
Out[5]: <keras.callbacks.History at 0x7f02264cc690>
```

Because the softmax is integrated into the output layer, the output is a vector of probabilities.

```
In [6]: p_nonpreferred = model.predict(X_train)
print(p_nonpreferred[:2])
print("largest value", np.max(p_nonpreferred), "smallest value", np.min(p_nonpreferred))
```

```
[3.15e-03 3.8e-03 9.83e-01 8.62e-02]
[9.92e-01 8.15e-03 2.64e-05 2.33e-05]
largest value 0.9999964 smallest value 3.4137022e-09
```

Preferred

Recall from lecture, more stable and accurate results can be obtained if the softmax and loss are combined during training. This is enabled by the 'preferred' organization shown here.

More numerically accurate implementation of softmax

Softmax regression

$$a_i = \frac{e^{z_i}}{e^{z_1} + \dots + e^{z_N}} = P(y=i|\mathbf{x})$$

$$a_{\bar{y}} = \frac{e^{z_{\bar{y}}}}{e^{z_1} + \dots + e^{z_N}} = P(y=\bar{y}|\mathbf{x})$$

$$L(\mathbf{y}, \mathbf{a}) = -\log(a_y) + (1 - y)\log(a_{\bar{y}})$$

$$L(\mathbf{y}, \mathbf{a}) = -\log(a_y) + (1 - y)\log(a_{\bar{y}})$$

$$J(\mathbf{w}, b) = \text{average loss}$$

$$J(\mathbf{w}, b) = \text{average loss}$$

Figure 1 was created using the softmax function from the lab_utils_softmax module.

The loss function associated with Softmax, the cross-entropy loss, is:

$$L(\mathbf{a}, \mathbf{y}) = \begin{cases} -\log(a_i), & \text{if } y = i. \\ \vdots \\ -\log(a_N), & \text{if } y = N \end{cases} \quad (3)$$

Where y is the target category for this example and \mathbf{a} is the output of a softmax function. In particular, the values in \mathbf{a} are probabilities that sum to one.

Recall: In this course, Loss is for one example while Cost covers all examples.

Note in (3) above, only the line that corresponds to the target contributes to the loss, other lines are zero. To write the cost equation we need an 'indicator function' that will be 1 when the index matches the target and zero otherwise.

$$1_{\{y == n\}} = \begin{cases} 1, & \text{if } y == n. \\ 0, & \text{otherwise.} \end{cases}$$

Now the cost is:

$$J(\mathbf{w}, b) = -\frac{1}{m} \left[\sum_{i=1}^m 1_{\{y^{(i)} == y\}} \log \frac{e^{z_y^{(i)}}}{\sum_{k=1}^N e^{z_k^{(i)}}} \right] \quad (4)$$

Where m is the number of examples, N is the number of outputs. This is the average of all the losses.

Tensorflow

This lab will discuss two ways of implementing the softmax, cross-entropy loss in Tensorflow, the 'obvious' method and the 'preferred' method. The former is the most straightforward while the latter is more numerically stable.

Let's start by creating a dataset to train a multiclass classification model.

```
In [4]: # make dataset for example
centers = [[-5, 2], [-2, -2], [1, 2], [5, -2]]
X_train, y_train = make_blobs(n_samples=2000, centers=centers, cluster_std=1.0, random_state=30)
```

The Obvious organization

The model below is implemented with the softmax as an activation in the final Dense layer. The loss function is separately specified in the `compile` directive.

The loss function is `SparseCategoricalCrossentropy`. This loss is described in (3) above. In this model, the softmax takes place in the last layer. The loss function takes in the softmax output which is a vector of probabilities.

```
In [5]: model = Sequential([
    Dense(25, activation = 'relu'),
    Dense(15, activation = 'relu'),
    Dense(4, activation = 'softmax')    # < softmax activation here
])
model.compile(
    loss=tf.keras.losses.SparseCategoricalCrossentropy(),
    optimizer=tf.keras.optimizers.Adam(0.001),
)
model.fit(
    X_train,y_train,
    epochs=10
)
```

```
Epoch 0/10
63/63 [=====] - 0s 1ms/step - loss: 0.6443
Epoch 3/10
63/63 [=====] - 0s 1ms/step - loss: 0.2465
Epoch 4/10
63/63 [=====] - 0s 915us/step - loss: 0.1158
Epoch 5/10
63/63 [=====] - 0s 1ms/step - loss: 0.0792
Epoch 6/10
63/63 [=====] - 0s 1ms/step - loss: 0.0632
Epoch 7/10
63/63 [=====] - 0s 1ms/step - loss: 0.0543
Epoch 8/10
63/63 [=====] - 0s 932us/step - loss: 0.0482
Epoch 9/10
63/63 [=====] - 0s 1ms/step - loss: 0.0425
Epoch 10/10
63/63 [=====] - 0s 1ms/step - loss: 0.0385
63/63 [=====] - 0s 1ms/step - loss: 0.0385
```

```
Out[5]: <keras.callbacks.History at 0x7f02264cc690>
```

Because the softmax is integrated into the output layer, the output is a vector of probabilities.

```
In [6]: p_nonpreferred = model.predict(X_train)
print(p_nonpreferred[:2])
print("largest value", np.max(p_nonpreferred), "smallest value", np.min(p_nonpreferred))
```

```
[3.15e-03 3.8
```