

Biostatistics (STA5195) Lab 3 Report

Iris Dataset Analysis

Christopher Rutherford

September 18, 2020

Abstract

The Iris dataset is a commonly used dataset used for demonstrating machine/deep learning algorithms. It contains 150 data points, 50 each for 3 species of Iris: Iris setosa, Iris virginica, and Iris versicolor. It includes the measurements, in cm, of each flower's sepal length and width along with the petal length and width. We use a neural network to predict which subspecies of Iris a flower belongs to. We will also use the **wandb** (**w**eights **a**nd **b**iases) service to keep track of our neural network's metrics, hyperparameters, CPU/GPU utilization, and more.

1 Introduction

The Iris dataset is one of the most commonly used datasets for the demonstration of classification algorithms in machine learning. The measurements were collected by Edgar Anderson, and the dataset introduced by Ronald Fisher in 1936 to present an example of linear discriminant analysis [1, 2]. R includes the dataset in the base package, and the scikit-learn library includes it for Python. By today's "big data" standards, it is a relatively simple dataset containing 150 samples of three species of the Iris plant. There are 50 samples each of Iris setosa, Iris virginica, and Iris setosa. Each sample contains a measurement for the plant's sepal length, sepal width, petal length, and petal width, along with the species of Iris it belongs to.

2 Theory

With the independent variables all being numeric and the dependent variable being categorical, there is lots of flexibility with which model we can choose for our analysis. We will demonstrate how a neural network can be used to predict which species of Iris a flower belongs to given the four measurements. Figure 1 provides an overview of how this neural network can be visualized. We provide the four variables as an input to the model (yellow circles). This data is fed to the "hidden" layers (blue circles), which perform the necessary calculations. This information is fed to the output layer (orange circles) to determine which species of Iris the flower belongs to.

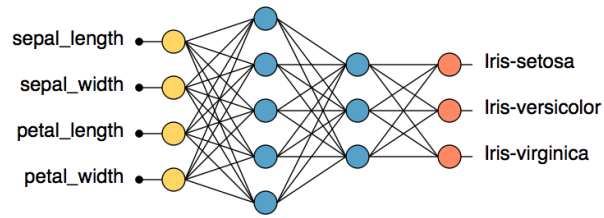


Figure 1: Overview of neural network for Iris species prediction

3 Procedures

As usual, we begin by importing the necessary libraries. This package is extremely convenient due to the inclusion of common libraries such as pandas, numpy, matplotlib, and more.

```
1 !pip install git+https://github.com/williamwardhahn/mpcr
2 from mpcr import *
```

Once this is done importing and installing, we can mount our google drive to our colab notebook:

```
1 drive.mount('/content/drive')
```

The next code block assumes the file is named `Iris.csv` and is saved in a folder named `Data`. Adjust this as necessary.

```
1 dataset = pd.read_csv('/content/drive/My Drive/Data/Iris.csv')
2 dataset
```

The printout of the dataset shows us our columns and the head/tail of our rows:

	Id	SepallengthCm	SepalwidthCm	PetallengthCm	PetalwidthCm	Species
	0	1	5.1	3.5	1.4	0.2 Iris-setosa
	1	2	4.9	3.0	1.4	0.2 Iris-setosa
	2	3	4.7	3.2	1.3	0.2 Iris-setosa
	3	4	4.6	3.1	1.5	0.2 Iris-setosa
	4	5	5.0	3.6	1.4	0.2 Iris-setosa

	145	146	6.7	3.0	5.2	2.3 Iris-virginica
	146	147	6.3	2.5	5.0	1.9 Iris-virginica
	147	148	6.5	3.0	5.2	2.0 Iris-virginica
	148	149	6.2	3.4	5.4	2.3 Iris-virginica
	149	150	5.9	3.0	5.1	1.8 Iris-virginica

150 rows x 6 columns

Figure 2: Iris dataset overview

We are interested in using the four measurements in order to predict the species, but first, we create some visualizations to see if we can detect some possible patterns in our data.

```
1 g = sns.pairplot(dataset) #create a histogram for each variable and
    a scatterplot for each pair of variables
```

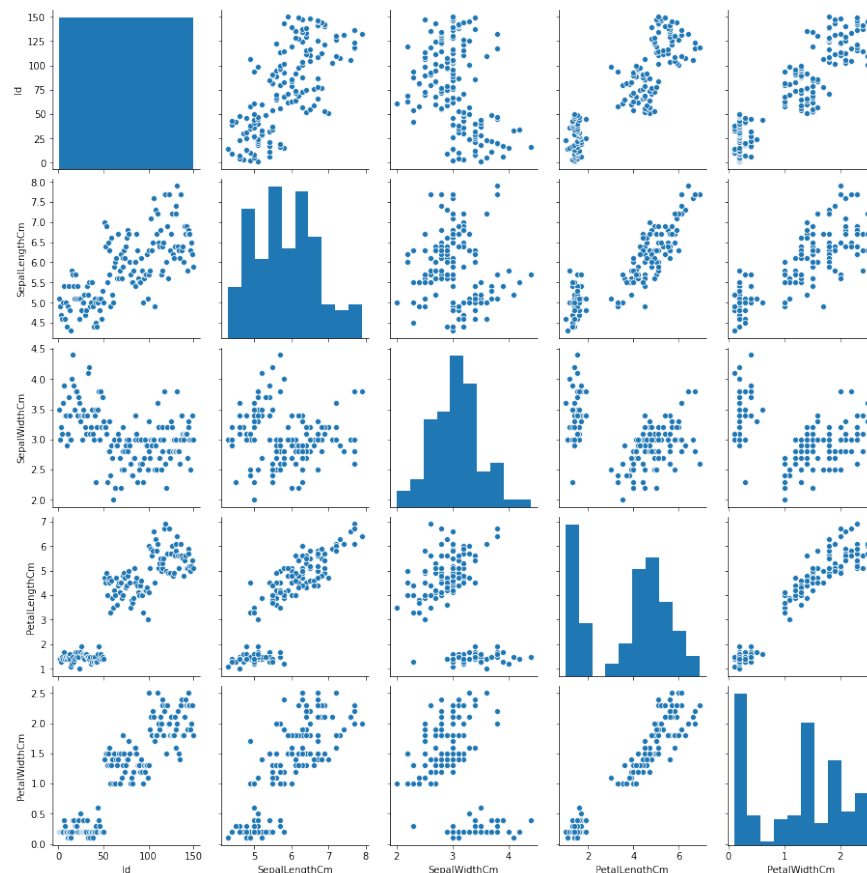


Figure 3: Histograms and scatterplots for each variable

Along the diagonal, we have histograms for each variable, and off of the diagonal, we have scatterplots between each pair of variables. Some of the plots show some very evident clusters of data. We can color each data point based on what species it belongs to to get a better look at the data:

```
1 g = sns.PairGrid(dataset, hue="Species") #same as above, but color
    each dot based on its species
2 g = g.map_diag(plt.hist) #histograms along the diagonal
3 g = g.map_offdiag(plt.scatter) #scatterplots off the diagonal
4 g = g.add_legend()
```

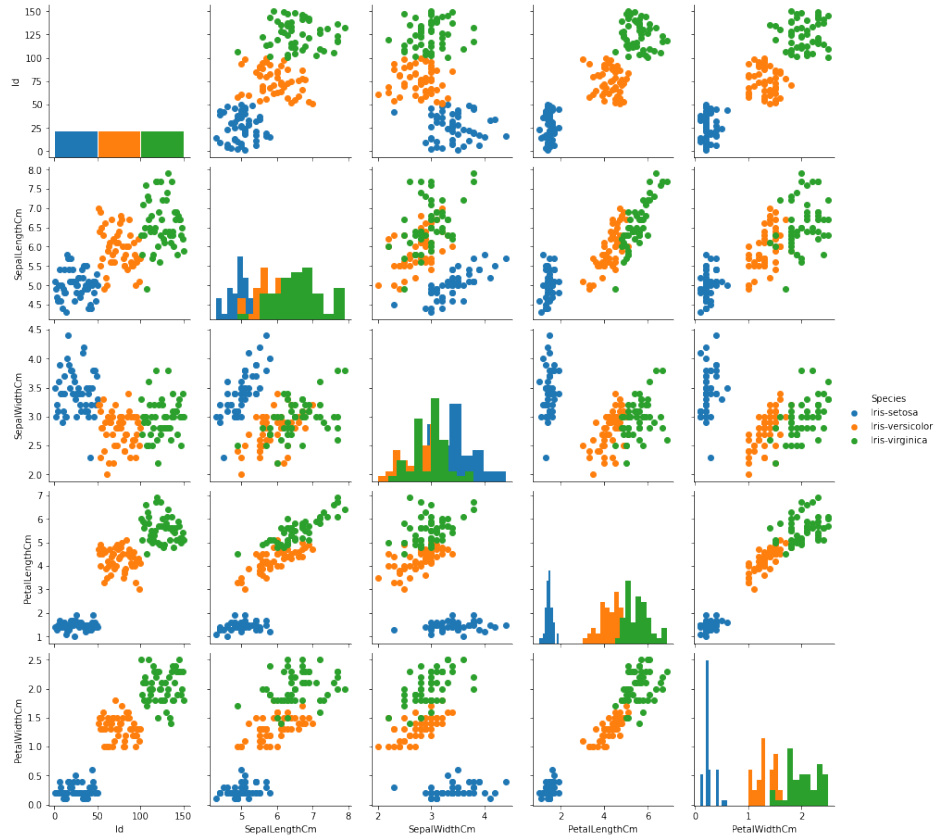


Figure 4: Histograms and scatterplots for each variable, colored by species

Providing a color for each species makes the clusters much more apparent in most of the scatterplots.

4 Analysis

Now we begin preparing our data for our model.

```
1 target_data = dataset[['Species']]
2 input_data = dataset.drop(['Id', 'Species'], axis=1)
```

The `target_data` variable stores the Y variable, which is the data for which species/class each flowers belongs to. `input_data` stores our X variables, or the ones used to predict the class. Since `input_data` is currently a DataFrame, we need to convert it to a numpy array for our model.

```
1 input_data = np.array(input_data)
```

We can create a visualization of this array:

```
1 plt.imshow(input_data[0:10,:])
```

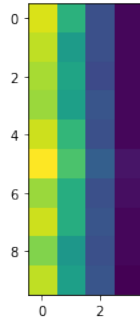


Figure 5: Heatmap of values of our input array

This heatmap indicates that the first column (sepal length) has the largest values (yellow) in the array, and the last column (petal width) has the smallest values (purple). Now, we one hot encode our target variable to make it easier for our network to work with:

```
1 target_data = pd.get_dummies(target_data.Species)
2 target_data
```

	Iris-setosa	Iris-versicolor	Iris-virginica
0	1	0	0
1	1	0	0
2	1	0	0
3	1	0	0
4	1	0	0
...
145	0	0	1
146	0	0	1
147	0	0	1
148	0	0	1
149	0	0	1

150 rows x 3 columns

Figure 6: One-hot encoded output variables

Instead of using one `species` column with labels, each species gets its own column. Data points with a 1 belong to that species, and 0 if they do not. As such, each data point will have a 1 in only one of the three columns. We can then convert this back into a single array where Iris setosa flowers will have a value of 0, Iris versicolor 1, and Iris virginica 2.

```
1 _, target_data = np.where(target_data==1)
2 target_data
```

Output: array([0, 0, 0, ..., 1, 1, 1, ..., 2, 2, 2,])

Now that we have our input and output data formatted for the model, we can create our train/test splits.

```

1 r = np.random.permutation(input_data.shape[0])
2 cut = int(0.8*len(r))
3 input_data.shape

```

Output: (150, 4)

The above code generates a random permutation of the integers 1, 2, ..., 150. Since we are using 80% of the data for training the model, we set the cutoff to be 80% of the length of `r`. This leaves us with 120 data points for training and 30 for testing. We create our splits:

```

1 X = input_data[r[:cut],:]
2 X_test = input_data[r[cut:],:]
3 Y = target_data[r[:cut]]
4 Y_test = target_data[r[cut:]]

```

Our `X` and `Y` variables both contain 120 random data points. `X` has the first 120 measurements and `Y` has the first 120 corresponding species. `X_test` and `Y_test` contain the remaining 30 data points.

We now define the softmax function, which will convert our output values to probabilities to determine which species a flower most likely belongs to:

```

1 def softmax(x):
2     s1 = torch.exp(x - torch.max(x,1)[0][:,None])
3     s = s1 / s1.sum(1)[: ,None]
4     return s

```

Mathematically, the softmax function is defined as follows:

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \quad (1)$$

where $\sigma(\mathbf{z})_i \in (0, 1)$ is the array (or vector) of probabilities of the i th element belonging to each of the output classes.

Since applying the softmax function to an array turns it into a probability distribution [3], we can use cross entropy as our loss function.

```

1 def cross_entropy(outputs, labels):
2     return -torch.sum(softmax(outputs).log()[range(outputs.size()
3     [0]), labels.long()])/outputs.size()[0]

```

The following block of code generates random normal variables with a mean of 0 and standard deviation of 0.1, but only values within 2 standard deviations of the mean are generated.

```

1 def randn_trunc(s): #Truncated Normal Random Numbers
2     mu = 0
3     sigma = 0.1
4     R = stats.truncnorm((-2*sigma - mu) / sigma, (2*sigma - mu) /
5     sigma, loc=mu, scale=sigma)
6     return R.rvs(s)

```

This allows us to generate random normal variables without getting extremely large/small values.

Now, we define a function to calculate the model's accuracy:

```
1 def acc(out, y):
2     with torch.no_grad():
3         return (torch.sum(torch.max(out, 1)[1] == y).item())/y.shape
4         [0]
```

This calculates how many flowers had their species labeled correctly, and divided that count by the total number of flowers (150).

The following functions tell Python to perform the computations on the GPU instead of the CPU, which greatly reduces model training time.

```
1 def GPU(data):
2     return torch.tensor(data, requires_grad=True, dtype=torch.float,
3         device=torch.device('cuda'))
4
5 def GPU_data(data):
6     return torch.tensor(data, requires_grad=False, dtype=torch.float,
7         device=torch.device('cuda'))
```

This function tells the model whether or not is currently training the model or evaluating it on the test set and gives it a random amount of data from the training/test set. The subset of data starts at the random integer `r` and is of length `b`, which is the batch size.

```
1 def get_batch(mode):
2     b = c.b
3     if mode == "train":
4         r = np.random.randint(X.shape[0]-b)
5         x = X[r:r+b,:]
6         y = Y[r:r+b]
7     elif mode == "test":
8         r = np.random.randint(X_test.shape[0]-b)
9         x = X_test[r:r+b,:]
10        y = Y_test[r:r+b]
11    return x,y
```

The following function tells us how much to adjust our gradient descent based on the learning rate, `c.h`.

```
1 def gradient_step(w):
2     for j in range(len(w)):
3         w[j].data = w[j].data - c.h*w[j].grad.data
4         w[j].grad.data.zero_()
```

An easy way to view our model's performance is by plotting the accuracy for both the training and testing sets:

```
1 def make_plots():
2     acc_train = acc(model(x,w),y)
3     xt,yt = get_batch('test')
4     acc_test = acc(model(xt,w),yt)
5     wb.log({"acc_train": acc_train, "acc_test": acc_test})
```

The following commands will tell Python to place our input and output data on the GPU by converting them to Tensor objects.

```
1 X = GPU_data(X)
2 Y = GPU_data(Y)
3 X_test = GPU_data(X_test)
4 Y_test = GPU_data(Y_test)
```

We define the ReLU layer, which is just $\max(0, x)$. That is, negative values are mapped to zero, while positive values remain the same.

```
1 def relu(x):
2     return x * (x > 0)

1 def model(x, w):
2     for j in range(len(w)):
3         x = relu(matmul(x, w[j]))
4     return x
```

```
1 wb.init(project="Iris");
2 c = wb.config
3 c.h = 0.05
4 c.b = 20
5 c.layers = 3
6 c.epochs = 2500
7 c.f_n = [4, 16, 16, 3]
8
9 w = [ GPU(randn_trunc((c.f_n[i], c.f_n[i+1]))) for i in range(c.
    layers) ]
10
11 for i in range(c.epochs):
12     x, y = get_batch('train')
13     loss = cross_entropy(softmax(model(x, w)), y)
14     loss.backward()
15     gradient_step(w)
16     if (i+1) % 1 == 0:
17         make_plots()
```

```
1 acc(model(X, w), Y)
```

Output: 0.30833333333333335

```
1 acc(model(X_test, w), Y_test)
```

Output: 0.46666666666666667

```
1 for i in range(len(w)):
2     plt.imshow(w[i].cpu().detach().numpy())
3     plt.show()
```

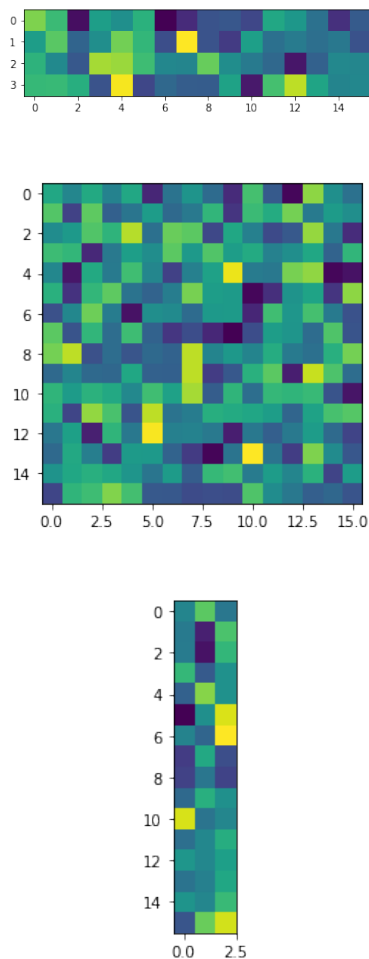



Figure 7: Visualization of weights

5 Conclusions

Overall, our model did not seem to be able to perform very well on this dataset. The model was only able to correctly label about 31% of the training data and 46.7% of the testing data, or 14 of the 30 testing flowers.

References

- [1] *Iris (plant)*, available at [https://en.wikipedia.org/wiki/Iris_\(plant\)](https://en.wikipedia.org/wiki/Iris_(plant))
- [2] *Iris flower data set*, available at https://en.wikipedia.org/wiki/Iris_flower_data_set
- [3] Ramsundar, B., Eastman, P., Walters, P., Pande, V. (2019). *Deep learning for the life sciences: Applying deep learning to genomics, microscopy, drug discovery and more* (1st ed.). Sebastopol, CA: O'Reilly Media.