	<pre>import scripts.io as io import scripts.util as util import scripts.NN as nn import random import numpy as np import matplotlib.pyplot as plt from sklearn.metrics import auc</pre>
	Part 1: Autoencoder implementation  1) Implement a feed-forward, neural network with standard sigmoidal units. Your implementation should accept a vector as input, be able to adjust network weights by backpropagation, and report the average loss per epoch. It should also allow for variation in the size of input layer, hidden layers, and output layer. We expect that you will be able to produce fast enough code to be of use in the learning task at hand.  a. To confirm that your implementation functions correctly, demonstrate its ability to correctly solve the 8 x 3 x 8 autoencoder task. Specifically, set up an autoencoder network consisting of an input layer (8 nodes), a hidden layer (3 nodes), and an output layer (8 nodes), all with sigmoidal units.  # Model Training #
[2]:	<pre>n_train = 1000 n_test = 1000  # Define network dimensions dim = [8, 3, 8]  # Define training data x = util.gen_label_array((dim[0], n_train))</pre>
	<pre># Define the number of epochs and learning rate for model training n_epochs = 2000 learning_rate = 2e-3  # Create and fit model net =nn.NeuralNetwork(dim, 'logistic') net.fit(x, x, n_epochs, learning_rate)  # Plot loss over range of epochs plt.plot(net.loss_list)</pre>
Out[2]:	plt.title('Loss Over Epoch') plt.xlabel('Epoch') plt.ylabel('Loss (Sum of Squares)')  Text(0, 0.5, 'Loss (Sum of Squares)')  Loss Over Epoch
	2.0 - (sale of the control of the co
In [3]:	# Model Testing #  # Define testing data and labels  x_ = util.gen_label_array((dim[0], n_test))  # Make prediction on testing data  yfin = net.predict(x_)
In [4]:	<pre># Let's visually investigate the model accuracy fig, (ax1, ax2) = plt.subplots(2, 1, figsize=(8, 3.8)) ax1.imshow(x_[:,0:50], cmap='gray') # Plot the input data ax1.set_title('Input Label Data') ax2.imshow(yfin[:,0:50], cmap='gray') # Plot the output data ax2.set_title('Output Label Data')</pre>
	ax2.set_xlabel('Label Index') fig.text(0.08, 0.5, 'Label Value', va='center', rotation='vertical')  Text(0.08, 0.5, 'Label Value')  Input Label Data  O  S  Text(0.08, 0.5, 'Label Value')
	9 0 10 20 30 40 Output Label Data  5 0 10 20 30 40 Label Index
In [5]:	<pre># Determine accuracy of model  count_correct = []  # Loop through testing results for i in range(0, n_test):     call = yfin[:,i:i+1].argmax() # Find index of maximum value in output layer     label = x_[:,i:i+1].argmax()     if call == label: # Compare output layer to known class         count_correct.append(1)     else:</pre>
	count_correct.append(0)  # Calculate Accuracy acc = np.sum(count_correct)/len(count_correct) * 100 print(f'Model Accuracy: {acc:.2f}%')  Model Accuracy: 100.00%  As expected, the autoencoder has perfect accuracy. Visualizing the model results shows the noise present in the model output, which is expected given that we compress the input data from eight to three in the neural network.
	Part 2: Adapt for classification, and develop training regime  2) Set up a procedure to encode DNA sequences (Rap1 binding sites) as input for your neural network. Consider how your encoding strategy may influence your network predictions.  a. Describe your process of encoding your training DNA sequences into input vectors in detail. Include a description of how you think the representation might affect your network's predictions.  # Import DNA sequences
In [7]:	<pre>pos = io.txt2str('data/rap1-lieb-positives.txt') neg = io.fa2str('data/yeast-upstream-1k-negative.fa')  # One-hot encode the sequences  pos_array = np.hstack(util.one_hot_dna(pos, 17)) neg_array = np.hstack(util.one_hot_dna(neg, 1000))  Individual DNA sequence strings are read in like so:</pre>
	'ACATCCGTGCACCTCCG'  The sequence is then one-hot encoded. This encoding scheme converts the the nucleotide categories: 'A', 'T', 'G', and 'C', into numeric dummy variables: [0,0,0,1], [0,0,1,0], [0,1,0,0], and [1,0,0,0], which can be input into the neural network.  Following the encoding step, you have the option of either concatenating your encoded vectors into a 1D array of size (n, 4). I chose the (n4, 1) arrangement as this requires less neurons in the network than the other vector arrangement.
In [9]: Out[9]:	An encoded example of the sequence above is shown below.  pos_array[:,0]  array([1., 0., 0., 0., 0., 0., 0., 1., 1., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 1., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0
	3) Design a training regime that will use both positive and negative training data to train your predictive model. Note that if you use a naive scheme here, which overweights the negative data, your system will probably not converge (it will just call everything a non-Rap1 site).  a. Describe your training regime. How was your training regime designed so as to prevent the negative training data from overwhelming the positive training data?  # Split data into training and testing data  # Split positive data pos_train_array, pos_test_array = util.train_test_split(pos_array, 80)
	# As there are so many negative sequence examples, let's take only a portion # of the negative sequence data. Let's also trim the 1kb sequences down # to 17bp. Both of these selection will be done randomly. samp_len = 300 tot_len = neg_array.shape[1]  n = random.randint(0, tot_len-samp_len) neg_array_trim = neg_array[:,n:n+samp_len] neg_array_sample = util.sample_array(neg_array_trim, 17, 4)
	neg_train_array, neg_test_array = util.train_test_split(neg_array_sample, 180)  For my training regime I randomly split the positive and negative sets into training and testing sets. For the positive example this resulted in ~60% training ~40% testing (80, 57). For the negative set, I randomly selected 17bp sequences from the given 1kb sequences, and then used only a small portion ~10% of the overall testing data in order to avoid overfitting my model on negative training data. The negative training data was then split to ~60% training ~40% testing (180, 120).  # Create training and testing arrays  train_array = np.concatenate((pos_train_array, neg_train_array), axis=1)  test_array = np.concatenate((pos_test_array, neg_test_array), axis=1)
In [12]:	<pre>test_array = np.concatenate((pos_test_array, neg_test_array), axis=1)  train_shuffle = np.random.permutation(train_array.shape[1]) test_shuffle = np.random.permutation(test_array.shape[1])  train_array = train_array[:, train_shuffle] test_array = test_array[:, test_shuffle]  # Generate training and testing label arrays  trp, trn = pos_train_array.shape[1], neg_train_array.shape[1]</pre>
	trp, trn = pos_train_array.shape[1], neg_train_array.shape[1] tep, ten = pos_test_array.shape[1], neg_test_array.shape[1]  train_labels = np.array([[1]*trp + [0]*trn]) test_labels = np.array([[1]*tep + [0]*ten])  train_labels = train_labels[:, train_shuffle] test_labels = test_labels[:, test_shuffle]  4) Modify your implementation to take as input positive and negative examples of Rap1 binding sites (using your encoding from Q2) and produce an output probability between [0 - 1.0] indicating classification as a binding site (1.0) or not (0.0). Select a network architecture, and train your network using the training regime you described in Q3 on all the data.
	<ul> <li>a. Provide an example of the input and output for one true positive sequence and one true negative sequence.</li> <li>b. Describe your network architecture, and the results of your training. How did your network perform in terms of minimizing error?</li> <li>c. What was your stop criterion for convergence in your learned parameters? How did you decide this?</li> <li># Model Training #</li> </ul>
	<pre># Define training data and labels x = train_array y = train_labels  # Define the number of epochs and learning rate for model training n_epochs = 3500 learning_rate = 2e-3  # Create and fit model dim = [68, 10, 1] net = nn.NeuralNetwork(dim, 'logistic')</pre>
Out[13]:	<pre>net = nn.NeuralNetwork(dim, 'logistic') net.fit(x, y, n_epochs, learning_rate)  # Plot loss over range of epochs plt.plot(net.loss_list) plt.title('Loss Over Epoch') plt.xlabel('Epoch') plt.ylabel('Loss (Sum of Squares)')</pre> Text(0, 0.5, 'Loss (Sum of Squares)')
	0.16
In [14]:	# Model Testing #
In [15]:	<pre>print('\n True Positive Input Array:', x_[:,0],</pre>
	'\n True Positive Output:', yfin[:,0][0],     '\n\n True Negative Input Array:', x_[:,-1],     '\n True Positive Input Array: [0. 1. 0. 0. 0. 0. 1. 0. 1. 0. 0. 1. 0. 0. 0. 1. 0. 0. 0. 1. 0. 0.  True Positive Input Array: [0. 1. 0. 0. 0. 0. 1. 0. 1. 0. 0. 1. 0. 0. 0. 1. 0. 0.  1. 0. 0. 0. 1. 0. 0. 0. 1. 0. 0. 0. 1. 0. 0. 0. 1. 0. 0. 0. 1. 0. 0.  1. 0. 0. 0. 1. 0. 0. 0. 1. 0. 0. 0. 1. 0. 0. 0. 0. 0. 1.]  True Positive Output: 0.01249520372785716  True Negative Input Array: [0. 1. 0. 0. 0. 0. 1. 0. 1. 0. 0. 0. 0. 1. 0. 0. 0. 0. 1. 0. 0. 0.  0. 0. 0. 1. 1. 0. 0. 0. 0. 0. 1. 0. 1. 0. 0. 0. 1. 0. 0. 0. 0.  1. 0. 0. 0. 0. 1. 1. 0. 0. 0. 0. 0. 0. 1. 0. 1. 0. 0. 0. 0. 1. 0. 0. 0.  1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.
In [16]:	<pre>0. 1. 0. 0. 1. 0. 0. 1. 0. 0. 0. 1. 0. 0. 0. 1. 0. 0. 0. 1. 0. 0.] True Negative Output: 0.0003604144330925875  # Determine accuracy of model  tot = yshape[1]  res = np.around(yfin) mis = np.sum(np.abs(y res))  # Calculate accuracy</pre>
	acc = ((tot-mis)/tot) * 100 print(f'Model Accuracy: {acc:.2f}%')  Model Accuracy: 98.31%  b. The network has three layer with the following dimmensions:  68:10:1  Error miminimalization during training can be seen above, under the model training cell. Over 2000 epochs with a learning rate of .002, loss decreased approximately exponentially from ~0.12 to ~0.0015. I belive
	this is very good loss minimalization.  My stop criteria was a loss of <= 0.0001, as such a low loss indicates a well trainined model, and further training may result in overfitting. This value can be changed as needed for different learning tasks.  Part 3: Cross-validation  5) Evaluate your model's classification performance via k-fold cross validation.  a. How can you use k-fold cross validation to determine your model's performance?
	<ul> <li>b. Given the size of your dataset, positive and negative examples, how would you select a value for k?</li> <li>c. Using the selected value of k, determine a relevant metric of performance for each fold. Describe how your model performed under cross validation.</li> <li>a. k-fold cross validation is based on hold out validation, except that the data is cycled through in order to test on all values within the set one time. This allows for a less biased measurement of model accuracy, and helps to minimize potential overfitting.</li> <li>b. When selecting your value of k, you must balance the amount of training data you need for good model performance, with the amount of data you leave for testing within each fold. For our dataset of ~430 data points, a k of 6 allows for training on ~350 datapoints and testing on ~80 datapoints for each fold. I believe this strikes a good balance of not undertraining the model, while have a decent fidelety of your</li> </ul>
In [17]:	<pre>model metric of choice.  # Let's recombine our arrays from earlier so they can be split k-times input_array = np.concatenate((train_array, test_array), axis=1) output_array = np.concatenate((train_labels, test_labels), axis=1)  # K-fold cross validation setup k = 8</pre>
In [19]:	<pre>shuffle = np.random.permutation(input_array.shape[1]) split_shuffle = list(util.split(shuffle, k))  # Model training and testing  n_epochs = 2000 learning_rate = 2e-3 dim = [68, 10, 1]</pre>
	<pre>auc_list = []  for n, test in enumerate(split_shuffle):     train_list = split_shuffle.copy()     del train_list[n]     train = np.hstack(train_list)  x = input_array[:, train] y = output_array[:, train]</pre>
	<pre>x_ = input_array[:, test] y_ = output_array[:, test]  net = nn.NeuralNetwork(dim, 'logistic') net.fit(x, y, n_epochs, learning_rate)  # Make prediction on testing data yfin = net.predict(x_)  pred_list = util.pred_gen(yfin[0])  test_fixe = util.pred_gen(yfin[0])</pre>
	<pre>tpr, fpr = util.pr_calc(y_[0], pred_list) auc_list.append(auc(fpr, tpr))  plt.plot(fpr, tpr, lw=2) plt.plot([0, 1], [0, 1], color='C7', lw=2, linestyle=':') plt.title('Local Alignment ROC') plt.xlabel('False Positive Rate') plt.ylabel('True Positive Rate')  plt.figure(figsize=[5,5]) plt.show()</pre>
	<pre>print(f'Mean AUROC = {np.mean(auc_list):.3f}')  Local Alignment ROC  10 0.8</pre>
	0.0
	<figure 0="" 360x360="" axes="" size="" with=""> Mean AUROC = 1.000 c. I chose AUROC as my model metric. This metric allows me to measure model performance over all possible binary classification threshold values for each fold. More specifically, I can measure true positive rates over false positive rates for the given binary threshold values.</figure>
	My mean AUROC was 1.0, as all folds had perfect accuracy. This is very good model performance, and demonstrates a high classification accuracy.  Part 4: Extension
	Part 4: Extension  To extend my network, I added five other activation functions to test the performance of. The activation functions I tested were logistic, relu, arctan, tanh, and softplus. I measured model performance with these different activation functions using AUROC.  I attempted to ascribe appropriate learning rates for each activation function, but may not have been sucessful.  The logistic, acrtan, and tanh activation functions worked well overall, while the relu and softplus implementations were very similiar to random. While this may point to a bug in the implementation of those functions, it may speak to an inability of those functions to perform our models classification task. I believe it is more likely that the activation functions are not suited for our over model architecture/parameters.
	Part 4: Extension  To extend my network, I added five other activation functions to test the performance of. The activation functions I tested were logistic, relu, arctan, tanh, and softplus. I measured model performance with these different activation functions using AUROC.  I attempted to ascribe appropriate learning rates for each activation function, but may not have been sucessful.  The logistic, acrtan, and tanh activation functions worked well overall, while the relu and softplus implementations were very similiar to random. While this may point to a bug in the implementation of those functions, it may speak to an inability of those functions to perform our models classification task. I believe it is more likely that the activation functions are not suited for our over model architecture/parameters.
	Part 4: Extension  To extend my network, I added five other activation functions to test the performance of. The activation functions I tested were logistic, relu, arctan, tanh, and softplus. I measured model performance with these different activation functions using AUROC.  I attempted to ascribe appropriate learning rates for each activation function, but may not have been sucessful.  The logistic, acrtan, and tanh activation functions worked well overall, while the relu and softplus implementations were very similiar to random. While this may point to a bug in the implementation of those functions, it may speak to an inability of those functions to perform our models classification task. I believe it is more likely that the activation functions are not suited for our over model architecture/parameters.  # Define training data and labels x = train_array y = train_labels  # Define testing data and labels x_ = test_array y_ = test_labels  # Network parameters n_epochs = 3000
	Part 4: Extension  To extend my nework, I added five other activation functions to test the performance of. The activation functions it tested were logistic, relu, arctan, tanh, and softplus. I measured model performance with these different activation functions using AUROC.  I attempted to ascribe appropriate learning rates for each activation function, but may not have been successful.  The logistic, acrtan, and tanh activation functions worked well overall, while the relu and softplus implementations were very similar to random. While this may point to a bug in the implementation of those functions, it may speak to an inability of those functions to perform our models classification task. I believe it is more likely that the activation functions are not suited for our over model architecture/parameters.  # Define training data and labels  * Test a rarray  * Test in labels  # Network parameters  ** Opening testing data and labels  * Lest array  * Lest albels  # Network parameters  ** Opening testing data and labels  * Lest array  * Lest albels  # Network parameters  ** Opening testing data and labels  * Lest array  * Lest a labels  # Network parameters  ** Opening testing data and labels  * Lest array  * Lest a labels  # Network parameters  ** Opening testing data and labels  * Lest array  * Lest a labels  # Network parameters  ** Opening testing data and labels  * Lest array  * Lest a labels  # Network parameters  ** Opening testing data and labels  * Lest a raray  * Lest a labels  # Network parameters  ** Opening testing data and labels  * Lest a raray  * Lest a labels  # Network parameters  ** Opening testing test
	Part 4: Extension  To extend my network. I added the other activation functions to test the performance of. The activation functions I tested were logistic, relu, arctan, tanh, and sortplus. I measured model performance with these different activation functions using AUROC.  Lattempted to ascribe appropriate learning rates for each activation function, but may not have been successful.  The logistic, actron, and tanh activation functions worked well overall, while the real and sortplus implementations were very similar to random. While this may point to a bug in the implementation of those functions, it may speak to an inability of those functions to perform our models classification task. I believe it is more likely that the activation functions are not suited for our ower model architectural parameters.  # perform extending data and labels  * a lettain, arrain, labels  # Motora? parameters  # in a labels  # with a label parameters  # in a labels  # with a label parameters  # in a labels  # with a label parameters  # in a labels  # with a label parameters  # in a labels  # with a label parameters  # in a labels  # with a label parameters  # in a labels  # with a label parameters  # in a labels  # with a label parameters  # in a labels  # with a label parameters  # in a labels  # with a label parameters  # in a labels  # with a label parameters  # in a labels  # with a label parameters  # in a labels  # with a labels  # with a label parameters  # in a labels  # with a label parameters  # in a labels  # with a label parameters  # in a label
	Part 4: Extension  To extend my motivant, I added fire other activation functions to test the performance of. The activation functions the set were beginner activation functions and ALPOCC.  Interrupted to asserbe agropmance learning rates for each activation function. but may not have been successful.  The logistic action, and under activation functions worked well recent while the reduced sufficient requirements are not suited for our over model architecture parameters. If a finite produced in the set of the set for our over model architecture parameters. If a finite requirement is a finite parameter of these functions the period of these functions that activation functions are not suited for our over model architecture parameters. If a finite result is a finite requirement of these functions are not suited for our over model architecture parameters.  If a finite requirement is a finite requirement of these functions are not suited for our over model architecture parameters.  If a finite requirement is a finite requirement of these functions are not suited for our over model architecture parameters.  If a finite requirement is a finite requirement of the set of the se
In [20]:	Part 4: Extension  To oxed my network, I added five other advacion functions to test the performance of. The activation functions is resident were logistic, resp., arcsan, cach, and softplus. I measured model performance with threse different enablation bandlans using AUROC.  The logistic, cancer, and test activation functions worked well overall, while the resil and softplus implementations were very similar to random. While this may point to a bug in the implementation of mose functions, through you can a insolid you those functions to perform our models classification task. I believe it is more likely that the advactor functions are not suited for our over model architecture/parameters.  ### Destina continuity data and liabels  ### a continuity data and liabels  #### a continuity data and liabels  ### a continuity data and liabels  ### a con
In [20]:	Part 4: Extension To extend by reduced, Lacted the cate advision histories to test the performance of. The activation biscories flower to surgicular forms activation flower to surgicular forms activation flower to surgicular forms activation flower to the surgicular forms activation for the surgicular flower to make the surgicular flower flowe
In [20]:	Part 4: Extension  Its section by retards, tabled by cultimated and testing the performance of the activation function because over small an anathor, while the recy produce and a compared to the performance of the activation function in a compared to a market of the performance of the activation and the activation and a compared to a market of the activation and an activation in a compared to a market of the activation and an activation in a compared to a compared to control while the received and adjusts in present a compared to a compared t
In [20]:	Part 4: Extension  to device up record instant to our a constant is core at the time of the constant is constant in the constant of the constant in the consta
In [20]: In [21]:	Part 4: Extension  To extend you could listed the office addition hands to be not be performed of The periods introduce belonity, dual parts, period and income modified formation with the county of
In [20]: In [21]: In [22]:	Part 4: Extension  To extend you could listed the office additional values to be all a personne of Tito person function. Excellent to the acceptance completely and acceptance to the personne of the personne
In [20]:  In [21]:  In [22]:	Part 4 - Extension  The standard and an administration and administrat
In [20]: In [21]: In [22]: In [24]:	Part 4: Extension  To contrary out: Listed to an outdood shower shower dependence of the shower legisle, usus on an adold in house the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout the contrary of the shower legisle was an another throughout throug
In [21]: In [22]: In [23]: In [25]:	Part La Capacida  Part La Capa
In [21]: In [22]: In [23]: In [25]:	Fig. 1. Example 1. Exa
In [21]: In [22]: In [23]: In [25]:	Reviols (Description of the Content
In [21]: In [22]: In [23]: In [25]:	Part I From Son To Continue to
In [21]: In [22]: In [23]: In [25]:	Ratic Contact And Anthonis Contact Anthonis Contac
In [20]:  In [21]:  In [23]:  In [24]:	Ratic Contact And Anticon
In [21]: In [22]: In [23]: In [25]:	Radic Contact And Anticon Contact Anticon Cont