

CSE446 HW3

Brian Kang[‡]

Conceptual Questions

A.1

- a. False. By projecting our data on a k -dim subspace, $k = \text{full rank}(X)$ maximum. Its is possible to not have full rank, then we definitely will lose information.
- b. True. The max margin using support vectors are constructed such that the distance from the decision boundary is minimized, that is, lowest generalization error among all linear classifiers.
- c. True. Bootstrap used sampling with replacement.
- d. False, should be the rows of V^T .
- e. False. Applying LASSO after feature reduction using PCA does not return coefficients that have meaning, so it is not interpretable for the linear model.
- f. True. With varying k , the complexity of the clustering model varies. So, with the right k , we can find good clusters.
- g. Decrease sigma to increase variance and decrease bias and increase the fit so we are not underfitting the training data. This allows us to reduce our error as long as we do not decrease so much that we end up overfitting.

*Collaborated with Cindy Wu

[‡]Referred to Stackoverflow, Kaggle, and various other websites that helped structuring my codes.

Kernels and Bootstrap

A.2

Proof.

$$\begin{aligned}\phi(x) \cdot \phi(x') &= \sum_{i=0}^{\infty} \left(\frac{1}{\sqrt{i!}} e^{-\frac{x^2}{2}} x^i \right) * \left(\frac{1}{\sqrt{i!}} e^{-\frac{x'^2}{2}} x'^i \right) \\ &= e^{-\frac{x^2 + x'^2}{2}} \sum_{i=0}^{\infty} \frac{1}{i!} (xx')^i \\ &= e^{-\frac{x^2 + x'^2}{2}} * e^{xx'} \text{ by Taylor expansion of } e^z \\ &= e^{-\frac{x^2 + 2xx' + x'^2}{2}} = e^{-\frac{(x+x')^2}{2}}\end{aligned}$$

□

A.3

Here is the code: Plots and answers are below the code!

```
1 # Brian Kang
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from itertools import compress, combinations
5
6
7 # generate random data
8 def generateData(n):
9     np.random.seed(446)
10    X = np.random.uniform(size=n)
11    e = np.random.randn(n)
12    Y = f(X) + e
13    return X, Y
14
15
16 # for easier plotting later
17 def f(X):
18     return 4 * np.sin(np.pi * X) * np.cos(6 * np.pi * X ** 2)
19
20
21 # first kernel function
22 def kernelPoly(X, z, d):
23     return np.power(1 + X * z, d)
24
25
26 # second kernel function
27 def kernelRBF(X, z, gamma):
28     return np.exp(-gamma * np.square(X - z))
29
30
31 def kernelRidge(X, Y, kernelFunc, hyperparam, lmbda):
32     Xi, Xj = np.meshgrid(X, X) # get 2D array of all points made from [X by X]
33     K = kernelFunc(Xi, Xj, hyperparam) # do kernel evaluation
34     # np.linalg.lstsq(a, b, rcond = 'warning'):
35     # Solves the equation a*x = b by computing a vector x that
36     # minimizes the squared Euclidean 2-norm of b-ax
37     # https://stackoverflow.com/questions/54753132/understanding-numpys-lstsq
38     alphahat, _, _, _ = np.linalg.lstsq(K.T @ K + lmbda * K, K.T @ Y, rcond=None)
39     # get alpha hat
40     # return a lambda function of getting array of fhat = sum of alphahat *
41     # kernel poly or rbf
42     return lambda xprime: np.array([np.sum(alphahat * kernelFunc(X, z,
43     hyperparam)) for z in xprime])
44
45
46 # do both LOOCV and 10-fold by setting k=1 or 10
47 def kCV(k, X, Y, kernelFunc, hyperparam, lmbda):
48     total = len(X)
49     index = np.arange(total).astype(int)
50     # mix up index and make n/k by k matrix, each row is then a fold
51     fold = np.random.permutation(index).reshape(int(total / k), k)
52     error = np.zeros(int(total / k))
```

```

50 # do ridge
51 for i, testindx in enumerate(fold):
52     trainindx = np.ones(total).astype(int)
53     trainindx[testindx] = 0
54     # get training set from X using compress() in np.array format
55     Xtrain = np.array([xtrain for xtrain in compress(X, trainindx)])
56     Ytrain = np.array([ytrain for ytrain in compress(Y, trainindx)])
57     # define predictor function
58     fhat = kernelRidge(Xtrain, Ytrain, kernelFunc, hyperparam, lambda)
59     # SSR/n
60     error[i] = np.sum(np.power(Y[testindx] - fhat(X[testindx]), 2)) / len(
testindx)
61     return np.mean(error) # mean error
62
63
64 def problemsABC(n, nfold, B, figureNumber):
65     X, Y = generateData(n)
66
67     # Part 3.A.poly
68     # make some lambda values
69     LMBDA = np.power(10, np.linspace(-2, 2, 10))
70     # make some hyperparameter d's in the Natural numbers
71     D = np.arange(0, 15)
72     savePolyResult = np.zeros([len(LMBDA) * len(D), 3])
73     index = 0
74     for d in D:
75         for lambda in LMBDA:
76             savePolyResult[index] = np.array([d, lambda, kCV(nfold, X, Y,
kernelPoly, d, lambda)])
77             index = index + 1
78     # get vars corresponding with smallest error
79     bestd, bestlmbda, _ = savePolyResult[np.argmin(savePolyResult[:, 2])]
80     print("Best d for poly kernel is " + str(bestd))
81     print("Best lambda for poly kernel is " + str(bestlmbda))
82
83     # Part 3.B & C.poly
84     xvals = np.linspace(0, 1, 100)
85     # bootstrap
86     index = np.arange(n).astype(int)
87     saveFhat = []
88     for i in range(B):
89         bootstrap = np.random.choice(index, n, replace=True)
90         saveFhat.append(kernelRidge(X[bootstrap], Y[bootstrap], kernelPoly,
bestd, bestlmbda))
91     # confidence intervals using percentiles
92     sortedFhat = np.sort(np.array([fhat(xvals) for fhat in saveFhat]), axis=0)
93     CIdown = sortedFhat[int(0.025 * B), :]
94     CIup = sortedFhat[int(0.975 * B), :]
95     # define predictor function
96     fhatPoly = kernelRidge(X, Y, kernelPoly, bestd, bestlmbda)
97     plt.figure(figureNumber)
98     ax = plt.gca()
99     plt.scatter(X, Y)
100     plt.plot(xvals, f(xvals), label="True f") # true f(x)
101     plt.plot(xvals, fhatPoly(xvals), label="Predicted f") # fhat(x)

```

```

102 ax.fill_between(xvals, CIdown, CIup, color='b', alpha=0.1)
103 plt.xlabel("X")
104 plt.ylabel("Y")
105 plt.legend()
106 plt.savefig("hw3_" + str(figureNumber) + ".png")
107
108 # Part 3.A.rbf
109 # make some lambda values
110 LMBDA = np.power(10, np.linspace(-2, 2, 10))
111 # make some hyperparameter gamma's
112 GAMMA = np.linspace(0, 50, 20)
113 saveRBFResult = np.zeros([len(LMBDA) * len(GAMMA), 3])
114 index = 0
115 for gamma in GAMMA:
116     for lmbda in LMBDA:
117         saveRBFResult[index] = np.array([gamma, lmbda, kCV(1, X, Y,
118 kernelRBF, gamma, lmbda)])
119         index = index + 1
120 # get vars corresponding with smallest error
121 bestgamma, bestlmbda, _ = saveRBFResult[np.argmin(saveRBFResult[:, 2])]
122 # if get gamma using heuristic method
123 # bestgamma = 1/np.median([np.power(xi[0] - xi[1], 2) for xi in combinations
124 (X, 2)])
125 print("Best gamma for rbf kernel is " + str(bestgamma))
126 print("Best lambda for rbf kernel is " + str(bestlmbda))
127
128 # Part 3.B & C.rbf
129 xvals = np.linspace(0, 1, 100)
130 # bootstrap
131 index = np.arange(n).astype(int)
132 saveFhat = []
133 for i in range(B):
134     bootstrap = np.random.choice(index, n, replace=True)
135     saveFhat.append(kernelRidge(X[bootstrap], Y[bootstrap], kernelRBF,
136 bestgamma, bestlmbda))
137 # confidence intervals using percentiles
138 sortedFhat = np.sort(np.array([fhat(xvals) for fhat in saveFhat]), axis=0)
139 CIdown = sortedFhat[int(0.025 * B), :]
140 CIup = sortedFhat[int(0.975 * B), :]
141 # define predictor function
142 fhatRBF = kernelRidge(X, Y, kernelRBF, bestgamma, bestlmbda)
143 plt.figure(figureNumber + 1)
144 ax = plt.gca()
145 plt.scatter(X, Y)
146 plt.plot(xvals, f(xvals), label="True f") # true f(x)
147 plt.plot(xvals, fhatRBF(xvals), label="Predicted f") # fhat(x)
148 ax.fill_between(xvals, CIdown, CIup, color='b', alpha=0.1)
149 plt.xlabel("X")
150 plt.ylabel("Y")
151 plt.legend()
152 plt.savefig("hw3_" + str(figureNumber + 1) + ".png")
153
154 return fhatPoly, fhatRBF # for part E

```

```

154 # problem 3.abc.1
155 problemsABC(n=30, nfold=1, B=300, figureNumber=1)
156 print()
157 # problem 3.abc.2 = 3.d
158 fhatPoly, fhatRBF = problemsABC(n=300, nfold=10, B=300, figureNumber=3)
159
160 # problem 3.e
161 m = 1000
162 B = 300
163 X, Y = generateData(m)
164 # bootstrap
165 index = np.arange(m).astype(int)
166 saveMu = []
167 for i in range(B):
168     bootstrap = np.random.choice(index, m, replace=True) # size m
169     error = np.mean(
170         np.power(Y[bootstrap] - fhatPoly(X[bootstrap]), 2) - np.power(Y[
171             bootstrap] - fhatRBF(X[bootstrap]), 2))
172     saveMu.append(error)
173 # confidence intervals using percentiles
174 sortedMu = np.sort(np.array(saveMu), axis=0)
175 CIdown = sortedMu[int(0.025 * B)]
176 CIup = sortedMu[int(0.975 * B)]
177 print()
178 print("The CI is: (" + str(CIdown) + ", " + str(CIup) + ")")
179 if (CIdown <= 0) and (CIup >= 0):
180     print(
181         "The confidence interval includes 0, so there exists statistically
182         significant\nevidence that one of f_rbf or f_poly is better at predicting Y
183         from X.")
184 else:
185     print(
186         "The confidence interval does NOT include 0, so there is NO
187         statistically significant\nevidence that one of f_rbf or f_poly is better at
188         predicting Y from X.")

```

a. Refer to figure 1.

```
Best d for poly kernel is 2.0  
Best lambda for poly kernel is 0.027825594022071243  
Best gamma for rbf kernel is 50.0  
Best lambda for rbf kernel is 0.0774263682681127
```

Figure 1: Parameter values for Poly and RBF Kernels

b. Refer to figure 2.

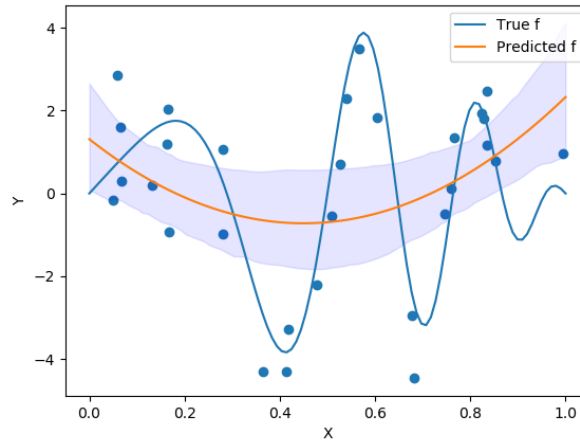


Figure 2: Poly Kernel v.s. True f

c. Refer to figure 3.

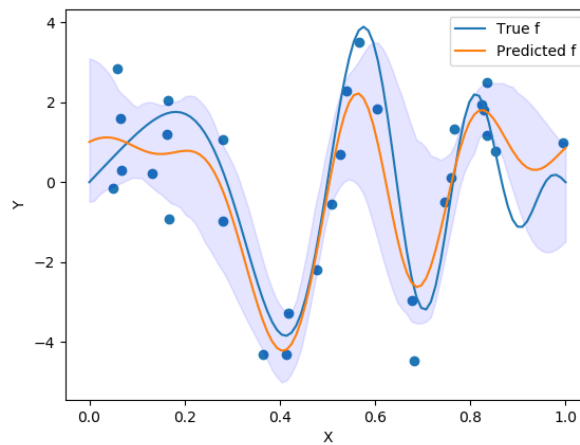


Figure 3: RBF Kernel v.s. True f

```
Best d for poly kernel is 13.0
Best lambda for poly kernel is 0.027825594022071243
Best gamma for rbf kernel is 50.0
Best lambda for rbf kernel is 0.01
```

Figure 4: Parameter values for Poly Kernel

d. Refer to figure 4 for parameter values.

Refer to figure 5 for poly kernel plot.

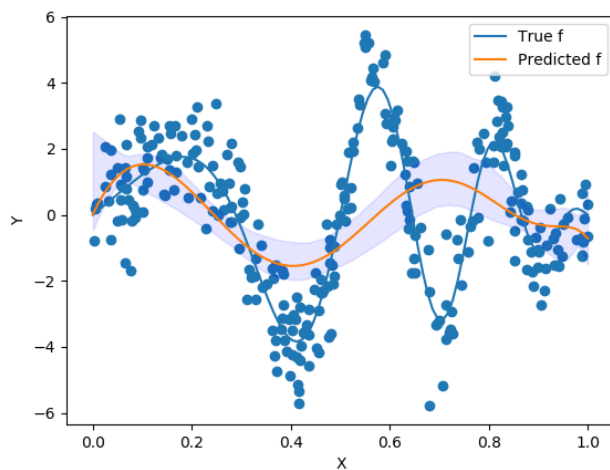


Figure 5: Poly Kernel v.s. True f

Refer to figure 6 for rbf kernel plot.

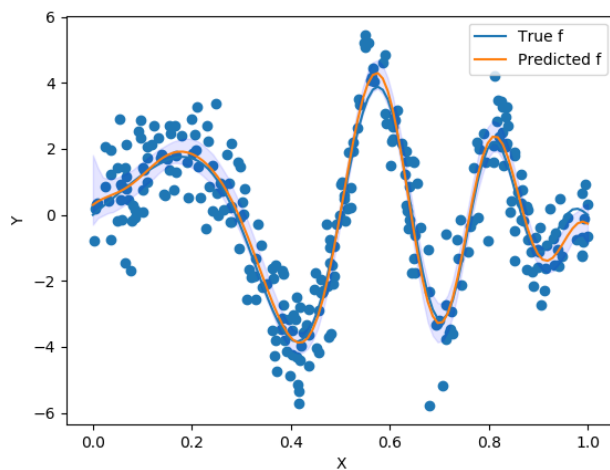


Figure 6: RBF Kernel v.s. True f

- e. Refer to figure 7 for the confidence interval and its interpretation.

TYPO Correction: "There **IS** evidence supporting the fact that one is better than the other because **0** is not in the confidence interval."

```
The CI is: (2.782856608988435, 3.4650986308114233)
```

```
The confidence interval does NOT include 0, so there is NO statistically significant  
evidence that one of f_rbf or f_poly is better at predicting Y from X.
```

Figure 7: Confidence Interval

k-means clustering

A.4

Here is the code: Plots and answers are below the code!

```
1 # Brian Kang
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from mnist import MNIST
5
6
7 def load_dataset():
8     mndata = MNIST('./data')
9     mndata.gz = True
10    X_train, labels_train = map(np.array, mndata.load_training())
11    X_test, labels_test = map(np.array, mndata.load_testing())
12    X_train = X_train / 255.0
13    X_test = X_test / 255.0
14    return X_train, labels_train, X_test, labels_test
15
16
17 # Part A
18 class kmeans:
19     def __init__(self, k):
20         self.k = k
21         self.objective = []
22         self.cluster = None
23         self.center = None
24
25     # lloyds algorithm
26     def train(self, X, maxIter):
27         # get k starting center points randomly chosen
28         n, _ = X.shape
29         np.random.seed(446)
30         startCenter = X[np.random.randint(0, n, size=self.k)]
31         # center_i = mu_i
32         center = np.copy(startCenter) # use the copy
33         dist = np.zeros((n, self.k)) # must have extra parentheses
34         for iter in range(maxIter):
35             for i in range(self.k):
36                 dist[:, i] = np.linalg.norm(X - center[i], axis=1) ** 2
37             # find centroids of each clusters
38             centroids = []
39             objective = 0
40             # referred to some stackoverflow page & lecture notes
41             for i in range(self.k):
42                 getCenter = X[np.argmin(dist, axis=1) == i] # point with min
43                 objective += np.sum(np.linalg.norm(getCenter - center[i], axis
44                 =1) ** 2)
45                 centroid = np.mean(getCenter, axis=0)
46                 centroids.append(centroid)
47                 center = np.copy(np.array(centroids)) # updated center
48                 self.objective.append(objective)
49             self.center = center
```

```

49         self.cluster = np.argmin(dist, axis=1)
50
51     def predict(self, X):
52         n, _ = X.shape
53         dist = np.zeros((n, self.k)) # must have extra parentheses
54         for i in range(self.k):
55             dist[:, i] = np.linalg.norm(X - self.center[i], axis=1) ** 2
56         pred = self.center[np.argmin(dist, axis=1)]
57         return pred
58
59
60 # Part B.1
61 x_train, y_train, x_test, y_test = load_dataset() # load in data
62
63 kmean1 = kmeans(10)
64 kmean1.train(x_train, 75)
65 plt.figure(1)
66 plt.plot(kmean1.objective)
67 plt.xlabel("Iteration")
68 plt.ylabel("Objective Value")
69 plt.savefig("hw3_9.png")
70
71 # Part B.2
72 # visualize the cluster centers
73 # k=10=2*5
74 plt.figure(2)
75 fig, axes = plt.subplots(2, 5)
76 for i, ax in enumerate(axes.flatten()):
77     ax.imshow(kmean1.center[i].reshape(28, 28), cmap='gray')
78     ax.set_title('Cluster Center {}'.format(i))
79 plt.tight_layout()
80 plt.savefig("hw3_10.png")
81
82 # Part C
83 ks = np.array([2, 4, 8, 16, 32, 64])
84 trainError = []
85 testError = []
86 for k in ks:
87     kmean1 = kmeans(k)
88     kmean1.train(x_train, 25)
89     predTrain = kmean1.predict(x_train)
90     predTest = kmean1.predict(x_test)
91     trainError.append(np.mean(np.linalg.norm(x_train - predTrain, axis=1) ** 2))
92     testError.append(np.mean(np.linalg.norm(x_test - predTest, axis=1) ** 2))
93
94 plt.figure(3)
95 plt.plot(ks, trainError, label='Train')
96 plt.plot(ks, testError, label='Test')
97 plt.xlabel("K")
98 plt.ylabel("Error")
99 plt.legend()
100 plt.savefig("hw3_11.png")

```

- a. Code for Lloyd's algorithm is above.
- b. Refer to figure below for objective function values.

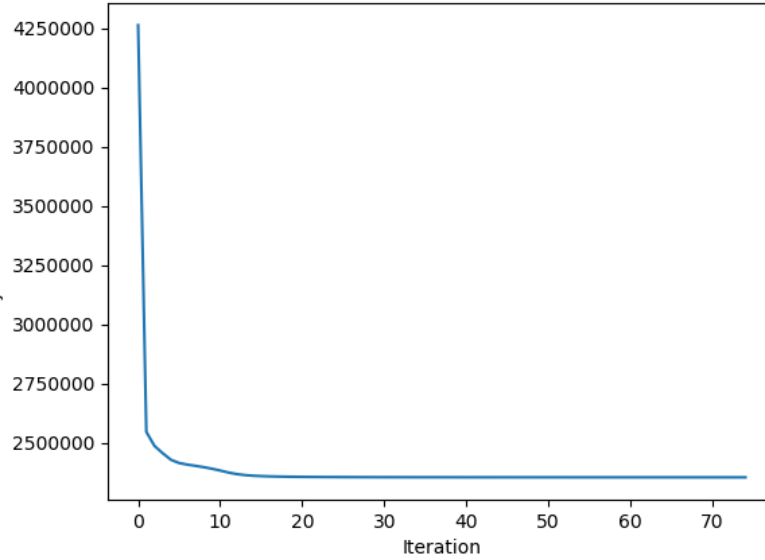


Figure 8: K-means Objective

Refer to figure below for Visualized Cluster Centers.

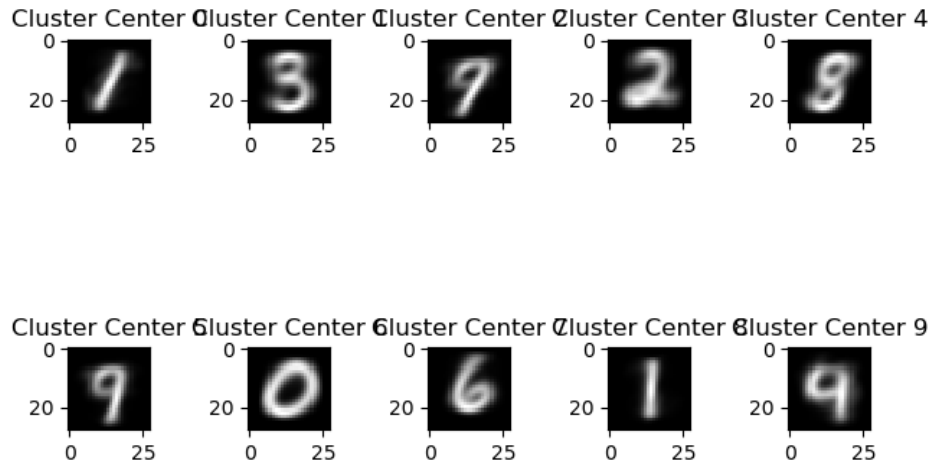


Figure 9: Visualized Cluster Centers

c. Refer to figure below for error values.

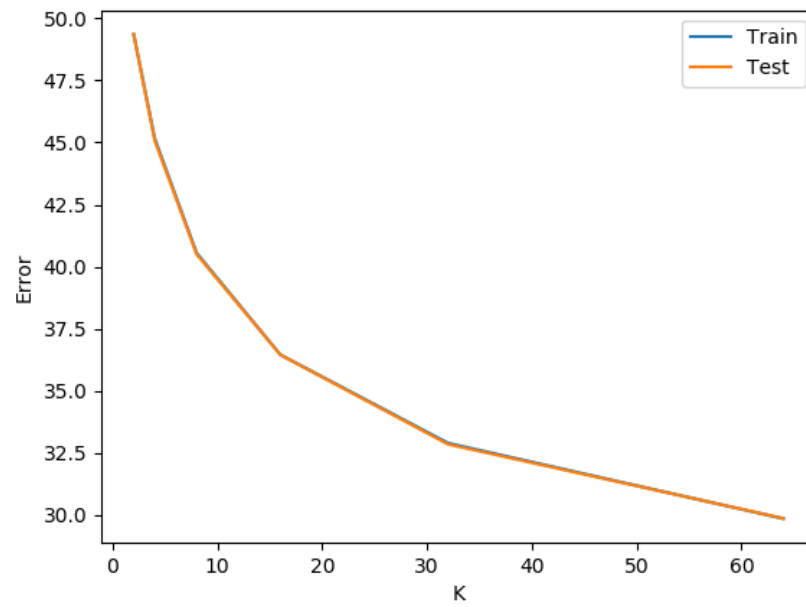


Figure 10: Model Errors per Varying K

Neural Networks for MNIST

A.5

Here is the code: Plots and answers are below the code!

```
1 # Brian Kang
2 import torch
3 import torch.nn as nn
4 import numpy as np
5 import matplotlib.pyplot as plt
6 from mnist import MNIST
7
8
9 def load_dataset():
10     mndata = MNIST('./data')
11     mndata.gz = True
12     X_train, labels_train = map(np.array, mndata.load_training())
13     X_test, labels_test = map(np.array, mndata.load_testing())
14     X_train = X_train / 255.0
15     X_test = X_test / 255.0
16     return X_train, labels_train, X_test, labels_test
17
18
19 class nnet(nn.Module):
20     def __init__(self, X, Y):
21         super(nnet, self).__init__()
22         self.X = X.T # 784 x n
23         self.Y = Y.T # 10 x n
24         d, n1 = self.X.shape # 784
25         k, n2 = self.Y.shape
26         h = 64 # number of nodes in the hidden layer
27         alpha = 1 / np.sqrt(d)
28         torch.manual_seed(446)
29         self.w0 = nn.Parameter(torch.FloatTensor(h, d).uniform_(-alpha, alpha),
requires_grad=True)
30         self.w1 = nn.Parameter(torch.FloatTensor(k, h).uniform_(-alpha, alpha),
requires_grad=True)
31         self.b0 = nn.Parameter(torch.FloatTensor(h, n1).uniform_(-alpha, alpha),
requires_grad=True)
32         self.b1 = nn.Parameter(torch.FloatTensor(k, n2).uniform_(-alpha, alpha),
requires_grad=True)
33         self.layer1 = torch.zeros(h, n1)
34         self.output = torch.zeros(k, n2)
35         self.relu = nn.ReLU() # we instantiate an instance of the ReLU module
36
37     def forward(self, X):
38         self.layer1 = self.relu(torch.add(torch.mm(self.w0, X.T), self.b0)) #
ReLU(W*X+b)
39         self.output = torch.add(torch.mm(self.w1, self.layer1), self.b1)
40         return self.output
41
42
43 # Part A
44 learningRate = 0.001
45 h = 64
```

```

46 d = 10
47
48 x_train, y_train, x_test, y_test = load_dataset() # load in data
49 num_class = 10 # this is k
50 onehot_label_train = np.eye(num_class)[y_train] # one-hot code the labels
51 onehot_label_test = np.eye(num_class)[y_test]
52
53 minibatch_size = 600
54 torch.manual_seed(446)
55 batch = np.random.randint(0, len(x_train), minibatch_size) # random index
    selection
56 model = nnet(np.array(x_train[batch, :]), np.array(onehot_label_train[batch, :])
    )
57
58 model.train()
59 optim = torch.optim.Adam(model.parameters(), lr=learningRate)
60 lossFunc = nn.CrossEntropyLoss()
61
62 i = 0
63 indx = []
64 lost = []
65 cont = True # start training
66 while cont:
67     optim.zero_grad() # zero all the gradients to get new gradients soon
68     outputs = model(torch.from_numpy(np.array(x_train[batch, :])).float())
69     label = [] # get label
70     for idx, vec in enumerate(np.array(onehot_label_train[batch, :])):
71         label = np.append(label, np.where(vec == 1)[0])
72     input = outputs.T
73     loss = lossFunc(input,
74                     torch.LongTensor(torch.from_numpy(np.array(label)).long()))
75     loss.backward() # computes derivatives of the loss with respect to W
76     optim.step() # make a step
77
78     i = i + 1
79     indx.append(i)
80     lost.append(loss.item())
81     if i % 1000 == 1:
82         learningRate = learningRate + 0.00025 # slowly increase learning rate
83     if i == 60000 and loss.item() >= 0.01:
84         i = 0
85     if loss.item() < 0.01:
86         cont = False # loss is small enough, stop training
87 # next random batch to train on
88 torch.manual_seed(446)
89 batch = np.random.randint(0, len(x_train), minibatch_size) # random index
    selection
90
91 # test loss
92 model.eval()
93 torch.manual_seed(446)
94 batch = np.random.randint(0, len(x_test), minibatch_size)
95 ii = 0
96 lostVal = []
97 for iter in range(0, i):

```

```

98     outputs = model(torch.from_numpy(np.array(x_test[batch, :])).float())
99     label = []
100     for idx, vec in enumerate(np.array(onehot_label_test[batch, :])):
101         label = np.append(label, np.where(vec == 1)[0])
102     input = outputs.T
103     loss = lossFunc(input,
104                     torch.LongTensor(torch.from_numpy(np.array(label)).long()))
105     ii = ii + 1
106     lostVal.append(loss.item())
107     torch.manual_seed(446)
108     batch = np.random.randint(0, len(x_test), minibatch_size)
109
110 # loss plot
111 plt.figure(1)
112 plt.plot(indx, lost, label="Train", linewidth=0.5)
113 plt.plot(indx, lostVal, label="Test", linewidth=0.5)
114 plt.xlabel("Iteration")
115 plt.ylabel("Cross Entropy Training Loss")
116 plt.suptitle("Minibatch of Size 600")
117 plt.legend()
118 plt.savefig("hw3_5.png")
119
120 # accuracy plot
121 acc = [1 - x for x in lost]
122 accVal = [1 - x for x in lostVal]
123 plt.figure(2)
124 plt.plot(indx, acc, label="Train", linewidth=0.5)
125 plt.plot(indx, accVal, label="Test", linewidth=0.5)
126 plt.xlabel("Iteration")
127 plt.ylabel("Accuracy")
128 plt.suptitle("Minibatch of Size 600")
129 plt.legend()
130 plt.ylim(bottom=0)
131 plt.savefig("hw3_6.png")
132
133
134 # Part B
135 class nnet2(nn.Module):
136     def __init__(self, X, Y):
137         super(nnet2, self).__init__()
138         self.X = X.T # 784 x n
139         self.Y = Y.T # 10 x n
140         d, n = self.X.shape # 784 x minibatch
141         k, _ = self.Y.shape
142         h1 = 32 # number of nodes in the hidden layer
143         h2 = 32 # second hidden layer
144         alpha = 1 / np.sqrt(d)
145         torch.manual_seed(446)
146         self.w0 = nn.Parameter(torch.FloatTensor(h1, d).uniform_(-alpha, alpha),
147                                requires_grad=True)
148         self.w1 = nn.Parameter(torch.FloatTensor(h2, h1).uniform_(-alpha, alpha),
149                                requires_grad=True)
150         self.w2 = nn.Parameter(torch.FloatTensor(k, h2).uniform_(-alpha, alpha),
151                                requires_grad=True)
152         self.b0 = nn.Parameter(torch.FloatTensor(h1, n).uniform_(-alpha, alpha),

```



```

requires_grad=True)
150     self.b1 = nn.Parameter(torch.FloatTensor(h2, n).uniform_(-alpha, alpha),
requires_grad=True)
151     self.b2 = nn.Parameter(torch.FloatTensor(k, n).uniform_(-alpha, alpha),
requires_grad=True)
152     self.layer1 = torch.zeros(h1, n)
153     self.layer2 = torch.zeros(h2, n)
154     self.output = torch.zeros(k, n)
155     self.relu = nn.ReLU() # we instantiate an instance of the ReLU module
156
157     def forward(self, X):
158         self.layer1 = self.relu(torch.add(torch.mm(self.w0, X.T), self.b0)) #
ReLU(W*X+b)
159         self.layer2 = self.relu(torch.add(torch.mm(self.w1, self.layer1), self.
b1))
160         self.output = torch.add(torch.mm(self.w2, self.layer2), self.b2)
161         return self.output
162
163
164 learningRate = 0.001
165 h1 = 32
166 h2 = 32
167 d = 10
168
169 minibatch_size = 600
170 torch.manual_seed(446)
171 batch = np.random.randint(0, len(x_train), minibatch_size) # random index
selection
172 model = nnet2(np.array(x_train[batch, :]), np.array(onehot_label_train[batch,
:])))
173
174 model.train()
175 optim = torch.optim.Adam(model.parameters(), lr=learningRate)
176 lossFunc = nn.CrossEntropyLoss()
177
178 i = 0
179 indx = []
180 lost = []
181 cont = True # start training
182 while cont:
183     optim.zero_grad() # zero all the gradients to get new gradients soon
184     outputs = model(torch.from_numpy(np.array(x_train[batch, :])).float())
185     label = [] # get label
186     for idx, vec in enumerate(np.array(onehot_label_train[batch, :])):
187         label = np.append(label, np.where(vec == 1)[0])
188     input = outputs.T
189     loss = lossFunc(input,
190                     torch.LongTensor(torch.from_numpy(np.array(label)).long()))
191     loss.backward() # computes derivatives of the loss with respect to W
192     optim.step() # make a step
193
194     i = i + 1
195     indx.append(i)
196     lost.append(loss.item())
197     if i % 1000 == 1:

```

```

198         learningRate = learningRate + 0.00025 # slowly increase learning rate
199     if i == 60000 and loss.item() >= 0.01:
200         i = 0
201     if loss.item() < 0.01:
202         cont = False # loss is small enough, stop training
203     # next random batch to train on
204     torch.manual_seed(446)
205     batch = np.random.randint(0, len(x_train), minibatch_size) # random index
        selection

206
207 # test loss
208 model.eval()
209 torch.manual_seed(446)
210 batch = np.random.randint(0, len(x_test), minibatch_size)
211 ii = 0
212 lostVal = []
213 for iter in range(0, i):
214     outputs = model(torch.from_numpy(np.array(x_test[batch, :])).float())
215     label = []
216     for idx, vec in enumerate(np.array(onehot_label_test[batch, :])):
217         label = np.append(label, np.where(vec == 1)[0])
218     input = outputs.T
219     loss = lossFunc(input,
220                     torch.LongTensor(torch.from_numpy(np.array(label)).long()))
221     ii = ii + 1
222     lostVal.append(loss.item())
223     torch.manual_seed(446)
224     batch = np.random.randint(0, len(x_test), minibatch_size)
225
226 # loss plot
227 plt.figure(3)
228 plt.plot(indx, lost, label="Train", linewidth=0.5)
229 plt.plot(indx, lostVal, label="Test", linewidth=0.5)
230 plt.xlabel("Iteration")
231 plt.ylabel("Cross Entropy Training Loss")
232 plt.suptitle("Minibatch of Size 600")
233 plt.legend()
234 plt.savefig("hw3_7.png")
235
236 # accuracy plot
237 acc = [1 - x for x in lost]
238 accVal = [1 - x for x in lostVal]
239 plt.figure(4)
240 plt.plot(indx, acc, label="Train", linewidth=0.5)
241 plt.plot(indx, accVal, label="Test", linewidth=0.5)
242 plt.xlabel("Iteration")
243 plt.ylabel("Accuracy")
244 plt.suptitle("Minibatch of Size 600")
245 plt.legend()
246 plt.ylim(bottom=0)
247 plt.savefig("hw3_8.png")

```

a. Refer to figure below for loss values.

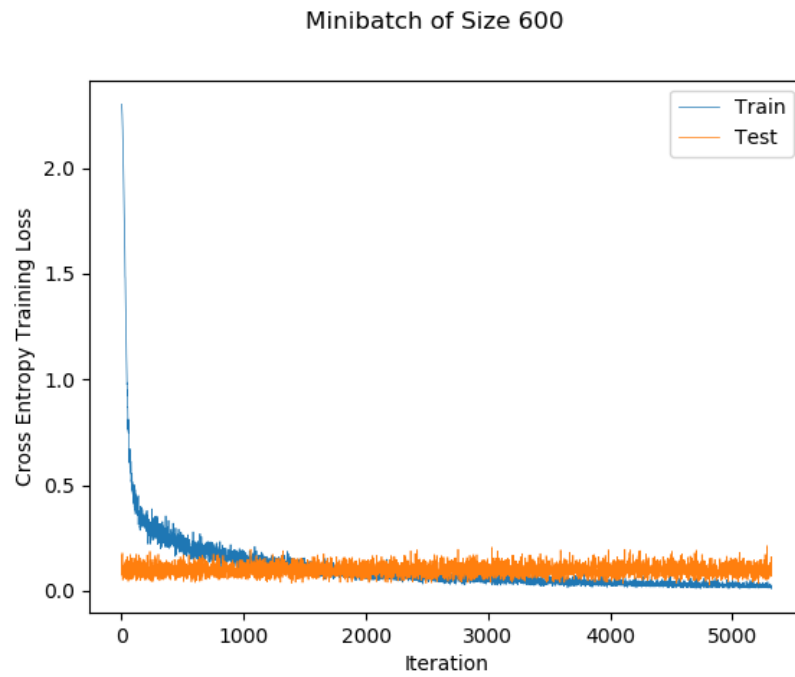


Figure 11: Loss: Wide and Shallow

Refer to figure below for accuracy values.

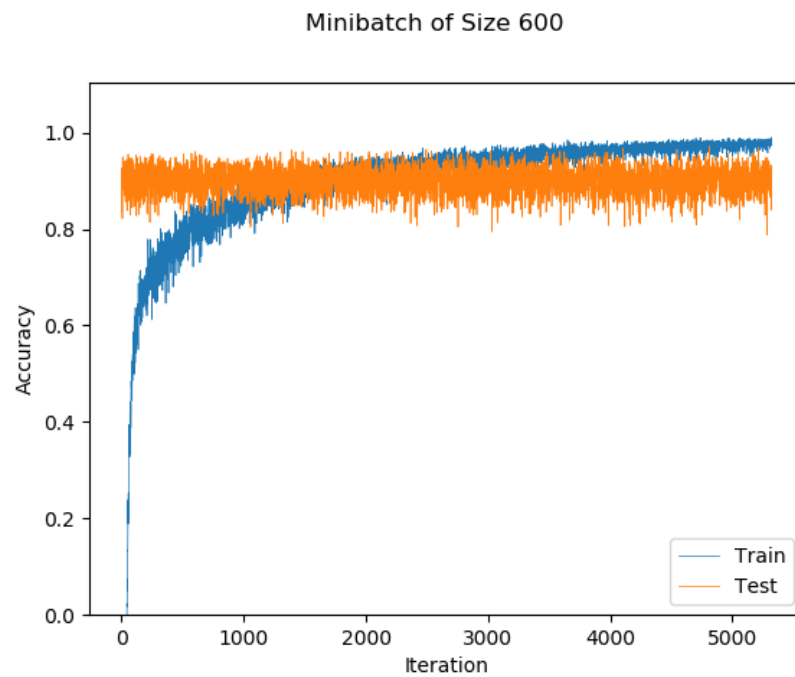


Figure 12: Accuracy: Wide and Shallow

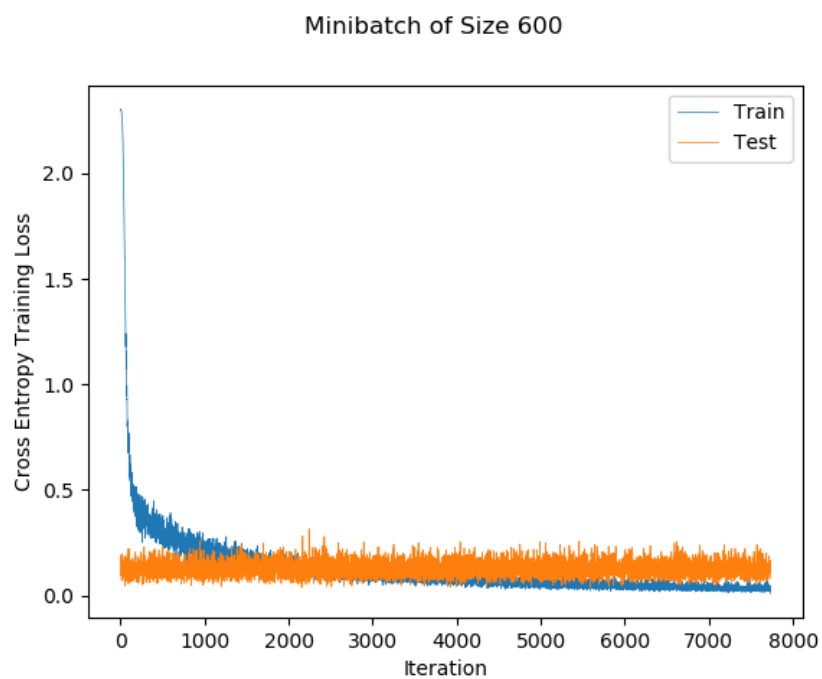


Figure 13: Loss: Narrow and Deep

- b. Refer to figure above for loss values.
Refer to figure below for accuracy values.

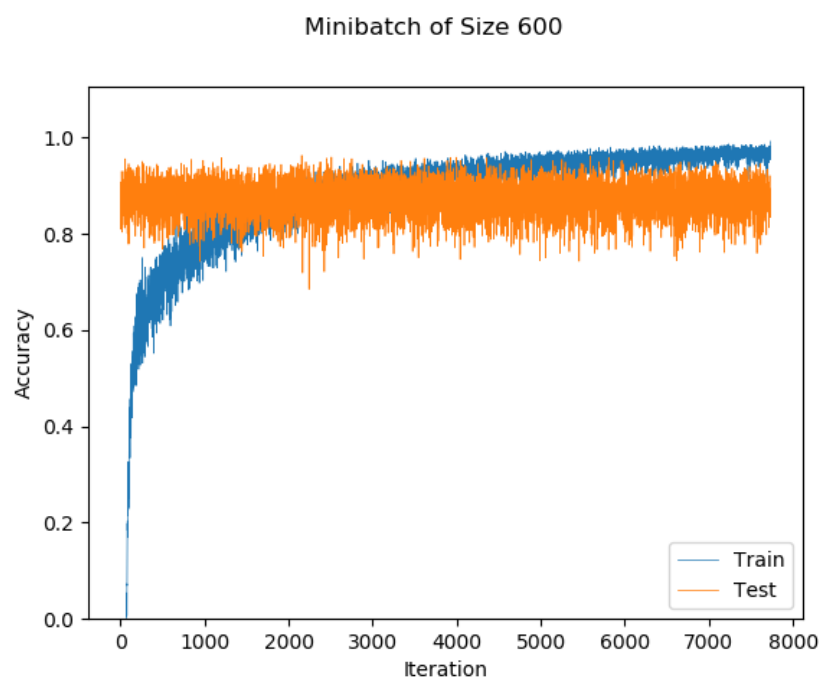


Figure 14: Accuracy: Narrow and Deep

- c. NOTE: Before answering I would like to first acknowledge the behavior of the tests results. I expected the loss to start high and then decrease (or accuracy start low then increase) like the training set at first. But we see that it hovers around a single value like noise behaves. Not sure if this is implementation error, but I reason this behavior that it happened because I use the trained weight and bias on the test set. Such results in predicting the label "at its best" everytime, so the variation can be explained by the randomness of each mini-batch iteration.

Now my answer to the question, it looks like general behavior of Wide&Shallow and Narrow&Deep are very similar. But after running the neural networks multiple times, it seemed like the wide&shallow one computed faster than the deeper one even if it was narrower. This can be correlated to the smaller iterations done for wide&shallow (5500-ish) versus narrow&deep (7800-ish). Also, it's observable that test narrow&deep has more variation around a value than the other. This might be related to the idea that wider and shallower neural nets are better at "memorization" while narrow and deep neural nets are better at "generalization." In a sense, wide neural nets use more of the data to "memorize" the behavior, while deep neural nets use each layer to "generally understand" a component.

PCA

A.6

Here is the code: Plots and answers are below the code!

```
1 # Brian Kang
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from mnist import MNIST
5
6
7 def load_dataset():
8     mndata = MNIST('./data')
9     mndata.gz = True
10    x_train, labels_train = map(np.array, mndata.load_training())
11    x_test, labels_test = map(np.array, mndata.load_testing())
12    x_train = x_train / 255.0
13    x_test = x_test / 255.0
14    return x_train, labels_train, x_test, labels_test
15
16
17 # Part A
18 x_train, y_train, x_test, y_test = load_dataset()
19 mu = np.mean(x_train, axis=0)
20 sigma = np.dot((x_train - mu).T, x_train - mu) / len(x_train) # covariance
21                        matrix
22 lmbda, V = np.linalg.eig(sigma)
23 # sort in decreasing order
24 sortIndx = np.argsort(lmbda)[::-1]
25 lmbda = lmbda[sortIndx].astype('float') # change type for later
26 V = V[:, sortIndx].astype('float')
27 indx = [0, 1, 9, 29, 49]
28 print("Eigenvalues:", lmbda[indx]) # get highest selected lmbda
29 sumlmbda = np.sum(lmbda)
30 print("Sum of Eigenvalues:", sumlmbda)
31
32 # Part C.1
33 trainError = []
34 testError = []
35 ks = np.arange(1, 101)
36 for k in ks:
37     predTrain = mu + np.dot(np.dot(x_train - mu, V[:, :k]), V[:, :k].T)
38     predTest = mu + np.dot(np.dot(x_test - mu, V[:, :k]), V[:, :k].T)
39     trainError.append(np.mean(np.linalg.norm(x_train - predTrain, axis=1) ** 2))
40     # mse
41     testError.append(np.mean(np.linalg.norm(x_test - predTest, axis=1) ** 2))
42
43 plt.figure(1)
44 plt.plot(ks, trainError, label="Train")
45 plt.plot(ks, testError, label="Test")
46 plt.xlabel("K")
47 plt.ylabel("Error")
48 plt.legend()
49 plt.savefig("hw3_13.png")
50
```

```

49 # Part C.2
50 value = []
51 for k in ks:
52     value.append(1 - np.sum(lmbda[:k]) / sumlmbda)
53
54 plt.figure(2)
55 plt.plot(ks, value)
56 plt.xlabel("K")
57 plt.ylabel("1-Sum(lambda 1...k) / Sum(lambda 1...d)")
58 plt.savefig("hw3_14.png")
59
60 # Part D
61 # visualize the PCA eigenvectors
62 # k=10=2*5
63 plt.figure(3)
64 fig, axes = plt.subplots(2, 5)
65 for i, ax in enumerate(axes.flatten()):
66     ax.imshow(V[:, i].reshape(28, 28), cmap='gray')
67     ax.set_title('PCA Eigen {}'.format(i))
68 plt.tight_layout()
69 plt.savefig("hw3_15.png")
70
71 # Part E
72 fignum = 16
73 for nums in [2, 6, 7]:
74     digit = x_train[y_train == nums][0] # arbitrary
75     ks = [5, 15, 40, 100]
76     # keep the brackets
77     reconstruct = [digit] + [mu + np.dot(np.dot(digit - mu, V[:, :k]), V[:, :k].
T) for k in ks]
78     fig, axes = plt.subplots(1, 5)
79     for i, ax in enumerate(axes.flatten()):
80         ax.imshow(reconstruct[i].reshape(28, 28), cmap='gray')
81         if i == 0:
82             ax.set_title("Real Num.")
83         else:
84             ax.set_title("k=" + str(ks[i - 1]))
85     plt.tight_layout()
86     plt.savefig("hw3_" + str(fignum) + ".png")
87     fignum = fignum + 1

```

```
Eigenvalues: [5.11678773 3.74132848 1.24272938 0.36425572 0.16970843]
Sum of Eigenvalues: 52.72503549512705
```

a.

b. For any k ,

$$x = \mu + (x - \mu)V_k V_k^T$$

V_k is the matrix of eigenvectors of the first k eigenvalues for any $k = 1, \dots, d$.

c. Refer to the plots below.

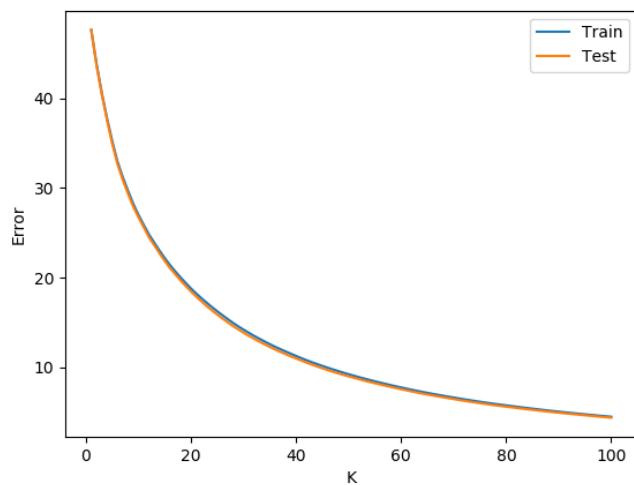


Figure 15: Reconstruction Error

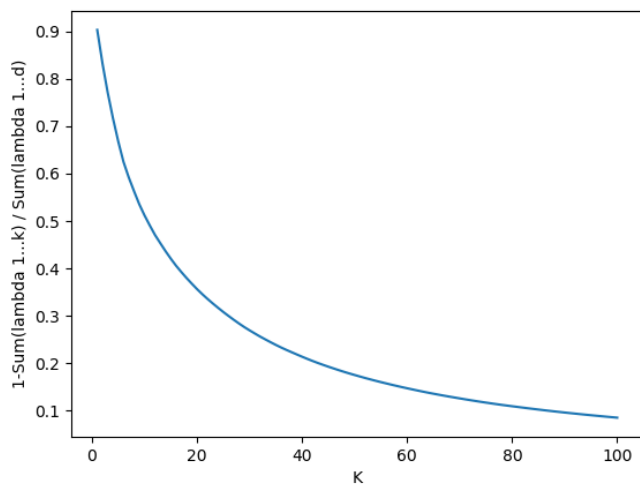
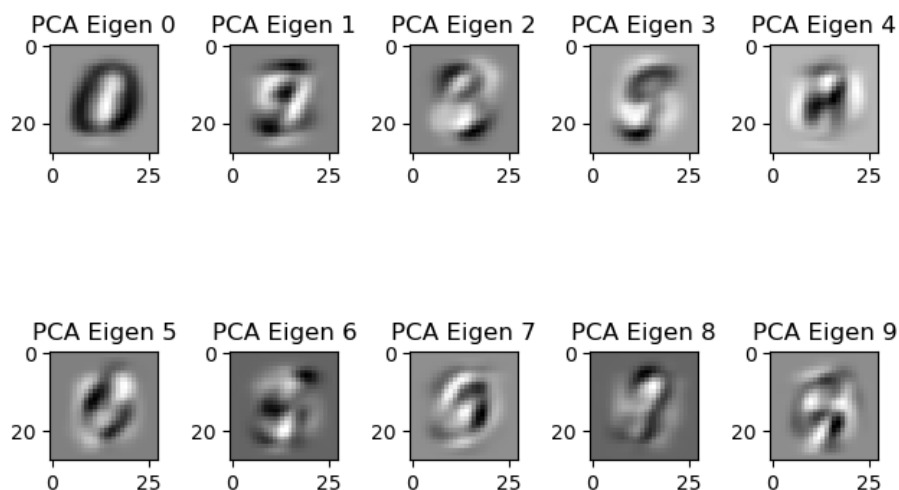


Figure 16: Plot of $1 - \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i}$

- d. Above is what the PCA directions captured. PCA is alternatively known in the math world as SVD, or Singular Value Decomposition. Vectors are decomposed into certain number of orthogonal vectors, which allows to get a more accurate picture as we multiply more of these orthogonal vectors (leading to higher dimensions). We are seeing this. PCA captured parts of the characteristics of the digits. Each plot captures a different characteristic, and when these are all multiplied to each other, it is likely that we will get a accurate representation of the digits.



- e. Below is presented the numbers 2, 6, 7 and what PCA (or SVD) captured with various k 's, i.e., various dimensions of orthogonal vectors. As said in (d), with the increase of k , we can see that each digit is appears more clear.

