

CSE546 Machine Learning HW3

Bobby Deng — 1663039 — dengy7

16 May 2020

A.1

a. [2 points] True or False: Given a data matrix $X \in R^{n \times d}$ where d is much smaller than n , if we project our data onto a k dimensional subspace using PCA where $k = \text{rank}(X)$, our projection will have 0 reconstruction error (we find a perfect representation of our data, with no information loss).

Solution:

True: If k is $\text{rank}(X)$, it means we have captured all variation from the data, and it can perfectly represent the data.

b. [2 points] True or False: The maximum margin decision boundaries that support vector machines construct have the lowest generalization error among all linear classifiers.

Solution:

False, SVM gives us a nice result, but it is not necessarily the lowest. It really depends on use cases and datasets. For example, the underlining true model need to be considered as well.

c. [2 points] True or False: An individual observation x_i can occur multiple times in a single bootstrap sample from a dataset X , even if x_i only occurs once in X .

Solution:

True: Bootstrap is done by draw sample with replacement.

d. [2 points] True or False: Suppose that the SVD of a square $n \times n$ matrix X is USV^T , where S is a diagonal $n \times n$ matrix. Then the rows of V are equal to the eigenvectors of $X^T X$.

Solution:

False, the columns are eigenvectors.

e. [2 points] True or False: Performing PCA to reduce the feature dimensionality and then applying the Lasso results in an interpretable linear model.

Solution:

False: When we use PCA, the data lost the interpretability, the data are reconstructed. Lasso can not make this model interpret-able.

f. [2 points] True or False: choosing k to minimize the k-means objective (see Equation (1) below) is a good way to find meaningful clusters.

Solution:

False, the bigger k we used, the less distance we will have, and this does not mean that our model gets better and better.

g. [2 points] Say you trained an SVM classifier with an RBF kernel $K(u, v) = \exp(-\frac{\|u-v\|_2^2}{2\sigma^2})$. It seems to under-fit the training set: should you increase or decrease σ ?

Solution:

We should decrease σ to increase some complexity of our model. When σ is big, the boundary tends to be pretty smooth, but when σ getting smaller, the boundary gets more strict.

A.2

Kernels and the Bootstrap

A2. [5 points] Suppose that our inputs x are one-dimensional and that our feature map is infinite-dimensional: $\phi(x)$ is a vector whose i th component is

$$\frac{1}{\sqrt{i!}} e^{-\frac{x^2}{2}} x^i$$

for all nonnegative integers i . (Thus, ϕ is an infinite-dimensional vector.) Show that $K(x, x') = e^{-\frac{(x-x')^2}{2}}$ is a kernel function for this feature map, i.e.,

$$\phi(x) \cdot \phi(x') = e^{-\frac{(x-x')^2}{2}}.$$

Hint: Use the Taylor expansion of e^z . (This is the one dimensional version of the Gaussian (RBF) kernel).

Solution:

From the question, we know that:

$$K(x, x') = \phi(x) \cdot \phi(x')$$

We also can get that:

$$\begin{aligned}\phi(x) &= \frac{1}{\sqrt{i!}} e^{-\frac{x^2}{2}} x^i \\ \phi(x') &= \frac{1}{\sqrt{i!}} e^{-\frac{x'^2}{2}} x'^i = \frac{1}{\sqrt{i!}} \frac{1}{\sqrt{e}}\end{aligned}$$

Then do the dot product:

$$\begin{aligned}K(x, x') &= \phi(x) \cdot \phi(x') \\ \phi(x) \cdot \phi(x') &= \sum_{i=0}^{\infty} \frac{1}{\sqrt{i!}} e^{-\frac{x^2}{2}} x^i \frac{1}{\sqrt{i!}} \frac{1}{\sqrt{e}} \\ \phi(x) \cdot \phi(x') &= \sum_{i=0}^{\infty} \frac{1}{i!} e^{-\frac{x^2}{2}} x^i \frac{1}{\sqrt{e}} \\ \phi(x) \cdot \phi(x') &= \sum_{i=0}^{\infty} \frac{1}{i!} e^{-\frac{x^2+1}{2}} x^i \\ \phi(x) \cdot \phi(x') &= e^{-\frac{x^2+1}{2}} \sum_{i=0}^{\infty} \frac{x^i}{i!}\end{aligned}$$

After doing Taylor Expansion:

$$\begin{aligned}\phi(x) \cdot \phi(x') &= e^{-\frac{x^2+1}{2}} e^x \\ \phi(x) \cdot \phi(x') &= e^{-\frac{x^2-2x+1}{2}} \\ \phi(x) \cdot \phi(x') &= e^{-\frac{(x-1)^2}{2}}\end{aligned}$$

Here we know that $x' = 1$, then:

$$\phi(x) \cdot \phi(x') = e^{-\frac{(x-x')^2}{2}}$$

So,

$$K(x, x') = e^{-\frac{(x-x')^2}{2}}$$

A.3

a.

A3. This problem will get you familiar with kernel ridge regression using the polynomial and RBF kernels. First, let's generate some data. Let $n = 30$ and $f_*(x) = 4 \sin(\pi x) \cos(6\pi x^2)$. For $i = 1, \dots, n$ let each x_i be drawn uniformly at random on $[0, 1]$ and $y_i = f_*(x_i) + \epsilon_i$ where $\epsilon_i \sim \mathcal{N}(0, 1)$.

For any function f , the true error and the train error are respectively defined as

$$\mathcal{E}_{true}(f) = E_{XY}[(f(X) - Y)^2], \quad \hat{\mathcal{E}}_{train}(f) = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2.$$

Using kernel ridge regression, construct a predictor

$$\hat{\alpha} = \arg \min_{\alpha} \|K\alpha - y\|^2 + \lambda \alpha^T K \alpha, \quad \hat{f}(x) = \sum_{i=1}^n \hat{\alpha}_i k(x_i, x)$$

where $K_{i,j} = k(x_i, x_j)$ is a kernel evaluation and λ is the regularization constant. Include any code you use for your experiments in your submission.

a. [10 points] Using leave-one-out cross validation, find a good λ and hyperparameter settings for the following kernels:

- $k_{poly}(x, z) = (1 + x^T z)^d$ where $d \in \mathbb{N}$ is a hyperparameter,
- $k_{rbf}(x, z) = \exp(-\gamma \|x - z\|^2)$ where $\gamma > 0$ is a hyperparameter¹.

Report the values of d , γ , and the λ values for both kernels.

Solution:

Best lamb for Poly kernel: 0.48828125.

Best d for Poly kernel: 44.

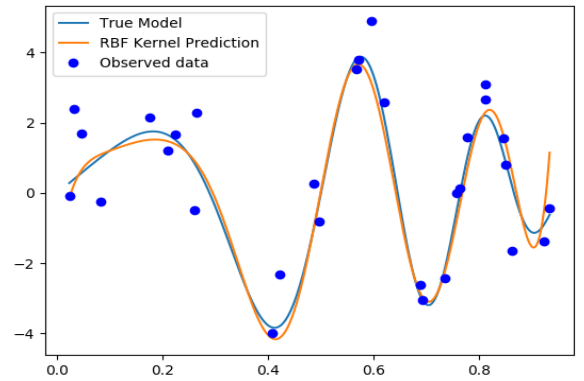
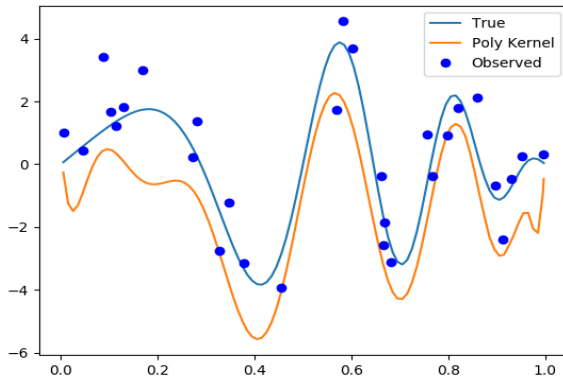
Best gamma for RBF kernel is : 8.175399541327897

Best Lambda for RBF kernel is : 9.313225746154785e-07

b.

b. [10 points] Let $\hat{f}_{poly}(x)$ and $\hat{f}_{rbf}(x)$ be the functions learned using the hyperparameters you found in part a. For a single plot per function $\hat{f} \in \{\hat{f}_{poly}(x), \hat{f}_{rbf}(x)\}$, plot the original data $\{(x_i, y_i)\}_{i=1}^n$, the true $f(x)$, and $\hat{f}(x)$ (i.e., define a fine grid on $[0, 1]$ to plot the functions).

Solution:



1

C.

- c. [5 points] We wish to build bootstrap percentile confidence intervals for $\hat{f}_{poly}(x)$ and $\hat{f}_{rbf}(x)$ for all $x \in [0, 1]$ from part b.² Use the non-parametric bootstrap with $B = 300$ bootstrap iterations to find 5% and 95% percentiles at each point x on a fine grid over $[0, 1]$.

Specifically, for each bootstrap sample $b \in \{1, \dots, B\}$, draw uniformly at random with replacement n samples from $\{(x_i, y_i)\}_{i=1}^n$, train an \hat{f}_b using the b th resampled dataset, compute $\hat{f}_b(x)$ for each x in your fine grid; let the 5th percentile at point x be the largest value ν such that $\frac{1}{B} \sum_{b=1}^B \mathbf{1}\{\hat{f}_b(x) \leq \nu\} \leq .05$, define the 95% analogously.

Plot the 5 and 95 percentile curves on the plots from part b.

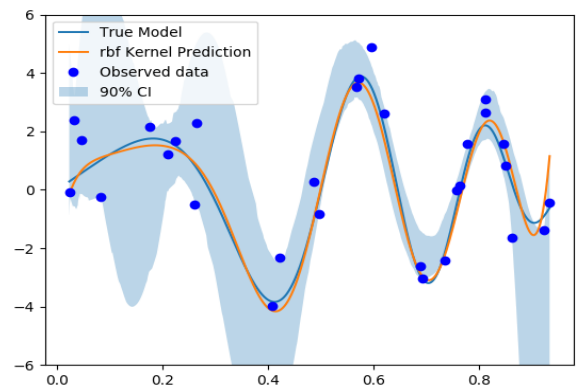
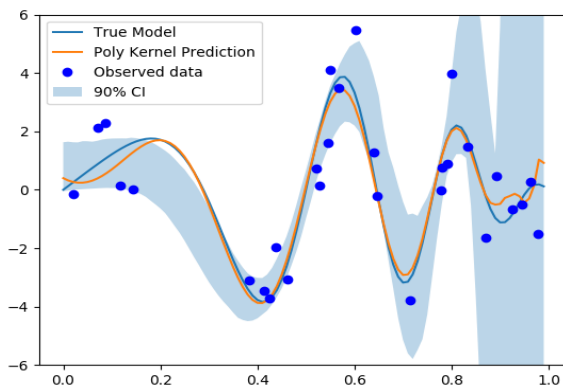
¹Given a dataset $x_1, \dots, x_n \in \mathbb{R}^d$, a heuristic for choosing a range of γ in the right ballpark is the inverse of the median of all $\binom{n}{2}$ squared distances $\|x_i - x_j\|_2^2$.

²See Hastie, Tibshirani, Friedman Ch. 8.2 for a review of the bootstrap procedure.

Solution:

In order to train a kernel model for each re-sampled bootstrap sample, we need to compute a new alpha and K each time, according to:

$$\hat{f}(x) = \sum_{i=1}^n \hat{\alpha}_i k(x_i, x)$$



Code for question a,b,c,

```
1 # b1 can also be used as part a's code
2 # A3.b1
3 import numpy as np
4 import matplotlib.pyplot as plt
```

```

5
6 n = 30
7 # np.random.seed(1)
8 x = np.random.uniform(0,1,n)
9 x_mean = np.mean(x)
10 x_sd = np.std(x)
11 # x = (x-x_mean) # x after standardization
12
13 y = 4*np.sin(np.pi*x)*np.cos(6*np.pi*(x**2)) + np.random.standard_normal(n)
14 # y = (y - x_mean) / x_sd
15 def k_poly(x, z, d):
16     a = x @ z.T
17     k = (1 + x @ z.T)**d
18     return k
19
20 error_validation_list = []
21 lamb = 500
22 lamb_list = []
23 d_list = []
24 for lamb in list(500 * (1/2)**(np.arange(0,20))):
25     for d in list(range(0, 51)):
26         error_validation = 0
27         print("Lam: ", lamb, ", d: ", d)
28         for i in range(n):
29             x_train = np.append(x[0:i], x[i+1:n])
30             y_train = np.append(y[0:i], y[i+1:n])
31             x_validation = x[i]
32             y_validation = y[i]
33             K = k_poly(x_train[:, np.newaxis], x_train[:, np.newaxis], d)
34             alpha = np.linalg.pinv(K + lamb) @ y_train
35             # in predicted y formula
36             k_xi_x = (1 + x_validation * x_train[np.newaxis, :]) ** d # use this when polynomial kernel
37             # k_xi_x = np.exp(-gamma*np.linalg.norm(x_validation - x_train[np.newaxis, :], 2))
38             y_predicted = alpha @ k_xi_x.T
39             error_validation += (y_predicted - y_validation).T @ (y_predicted - y_validation)
40             # error_validation = error_validation[0][0]
41         error_validation /= n
42         print("error_validation: ", error_validation)
43         error_validation_list.append(error_validation)
44         lamb_list.append(lamb)
45         d_list.append(d)
46
47 min_error = min(error_validation_list)
48 index_bootstrap_sample_min_error = error_validation_list.index(min(error_validation_list))
49 lamb_best_poly = lamb_list[index_bootstrap_sample_min_error]
50 d_best = d_list[index_bootstrap_sample_min_error]
51 print("Best lamb: ", lamb_best_poly, ", Best d: ", d_best)
52
53 # lamb_best_poly = 0.48828125
54 d_best = 30
55 # plots the comparison
56 # np.random.seed(1)
57 x_fine = np.array(list(np.arange(min(x),max(x), 0.01)))
58 n = len(x_fine)
59 y_fine_true = 4*np.sin(np.pi*x_fine)*np.cos(6*np.pi*(x_fine**2))
60 y_fine_grid = y_fine_true + np.random.standard_normal(n)
61 f_poly_predicted = []
62 for xi in x_fine:

```

```

63 K = k_poly(x_fine[:, np.newaxis], x_fine[:, np.newaxis], d_best)
64 alpha = np.linalg.pinv(K + lamb_best_poly) @ y_fine_grid
65 k_xi_x = (1 + xi * x_fine[np.newaxis, :]) ** d_best # use this when polynomial kernel
66 y_predicted = alpha @ k_xi_x.T
67 f_poly_predicted.append(y_predicted)
68
69 plt.plot(x_fine, y_fine_true, label='True')
70 plt.plot(x_fine, f_poly_predicted, label='Poly Kernel')
71 plt.plot(x, y, 'bo', label='Observed')
72 plt.xlabel("X")
73 plt.ylabel("Y")
74 plt.legend()
75 plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A3b_1_test.png")
76 plt.show()
77
78
79 # A3.c1
80 B = 300
81 n = 30
82 n_fine = len(x_fine)
83 # np.random.seed(0)
84 bootstrap_predicted_poly_matrix = []
85
86 for j in range(B):
87     index_bootstrap_sample = np.random.choice(n, n)
88     x_training = x[index_bootstrap_sample]
89     y_training = y[index_bootstrap_sample]
90     K = k_poly(x_training[:, np.newaxis], x_training[:, np.newaxis], d_best)
91     alpha = np.linalg.solve((K + lamb_best_poly * np.eye(n, n)), y_training)
92     y_predicted_bootstrap_ploy = []
93     for xi in x_fine:
94         y_predicted_bootstrap_ploy.append(np.sum((1 + xi * x_training[np.newaxis, :]) ** d_best @ alpha))
95     bootstrap_predicted_poly_matrix.append(y_predicted_bootstrap_ploy)
96 bootstrap_predicted_poly_matrix = np.array(bootstrap_predicted_poly_matrix)
97
98 percent_5_list_poly = []
99 percent_95_list_poly = []
100 for i in range(n_fine):
101     sorted_xi_from_300_B_sample = np.sort(bootstrap_predicted_poly_matrix[:, i])
102     x_percentile_5 = sorted_xi_from_300_B_sample[int(B * 0.05)]
103     x_percentile_95 = sorted_xi_from_300_B_sample[int(B * 0.95)]
104     percent_5_list_poly.append(x_percentile_5)
105     percent_95_list_poly.append(x_percentile_95)
106
107 plt.plot(x_fine, y_fine_true, label = 'True Model')
108 plt.plot(x_fine, f_poly_predicted, label = 'Poly Kernel Prediction')
109 plt.plot(x, y, 'bo', label = 'Observed data')
110 plt.fill_between(x_fine, percent_5_list_poly, percent_95_list_poly, alpha=0.3, label="90% CI")
111 plt.ylim(-6, 6)
112 plt.legend()
113 plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A3c_1_test.png")
114 plt.show()
115
116
117 #####
118 # A3.b2
119
120 def k_rbf(x, z, gamma):

```

```

121     return np.exp(-gamma*(x-z)*(x-z))
122
123 n = 30
124 # np.random.seed(0)
125 # x = np.random.rand(n)
126 x = np.random.uniform(0,1,n)
127 y_true = 4*np.sin(np.pi*x)*np.cos(6*np.pi*(x**2))
128 y = y_true + np.random.randn(n)
129 error_validation_list = []
130 lamb_list = []
131 gamma_list = []
132 d_list = []
133
134 lamb = 1
135 for lamb in list(500 * (1/2)**(np.arange(0,30))):
136     for gamma in list(50 * (1/1.1)**(np.arange(0,30))):
137         print("Lam: ", lamb, ", gamma: ", gamma)
138         error_validation = 0
139         for i in range(n):
140             x_train = np.append(x[0:i], x[i+1:n])
141             y_train = np.append(y[0:i], y[i+1:n])
142             x_validation = x[i]
143             y_validation = y[i]
144             K = k_rbf(x_train[:,np.newaxis],x_train[np.newaxis,:], gamma)
145             alpha = np.linalg.pinv(K + lamb) @ y_train
146             k_xi_x = np.exp(-gamma*(x_validation-x_train[np.newaxis,:])**2)
147             error_validation += (k_xi_x@alpha - y_validation).T@(k_xi_x@alpha - y_validation)
148         error_validation_list.append(error_validation)
149         print("error_validation: ", error_validation)
150         lamb_list.append(lamb)
151         gamma_list.append(gamma)
152
153 min_error = min(error_validation_list)
154 index_bootstrap_sample_min_error = error_validation_list.index(min_error)
155 lamb_best_rbf = lamb_list[index_bootstrap_sample_min_error]
156 gamma_best = gamma_list[index_bootstrap_sample_min_error]
157 print('Best gamma for RBF kernel is : ', gamma_best)
158 print('Best Lambda for RBF kernel is : ', lamb_best_rbf)
159
160
161 gamma_best= 10.175399541327897
162 lamb_best_rbf= 9.313225746154785e-07
163 # np.random.seed(10)
164
165 x_fine = np.arange(min(x),max(x),0.001)
166 n = len(x_fine)
167 y_fine_true = 4*np.sin(np.pi*x_fine)*np.cos(6*np.pi*(x_fine**2))
168 y_fine_grid = y_fine_true + np.random.standard_normal(n)
169
170 f_rbf_predicted = []
171 K_rbf = k_rbf(x_fine[:,np.newaxis],x_fine[np.newaxis,:], gamma_best)
172 alpha = np.linalg.solve((K_rbf + lamb_best_rbf*np.eye(n, n)), y_fine_grid)
173 for xi in x_fine:
174     f_rbf_predicted.append(np.sum(alpha * np.exp(-gamma_best*(xi-x_fine)**2)))
175
176 plt.plot(x_fine, y_fine_true, label = 'True Model')
177 plt.plot(x_fine, f_rbf_predicted, label = 'RBF Kernel Prediction')
178 plt.plot(x, y, 'bo', label = 'Observed data')

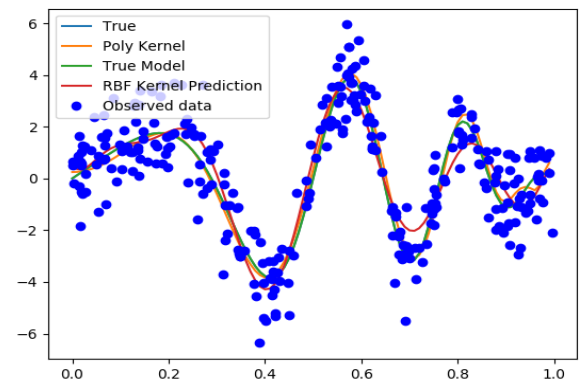
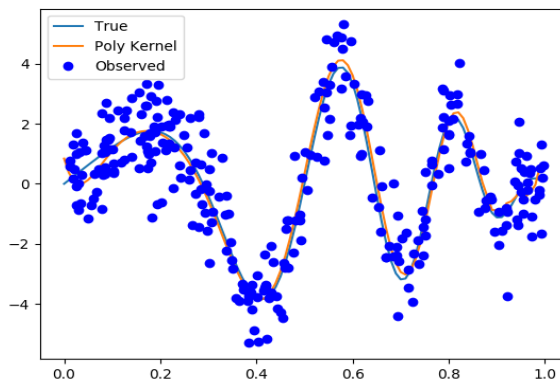
```

```

179 plt.legend()
180 plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A3b_2.png")
181 plt.show()
182
183 # A3.c2
184 B = 300
185 n=30
186 n_fine = len(x_fine)
187 # np.random.seed(0)
188 bootstrap_predicted_rbf_matrix = []
189 # user x, y from previous
190 for j in range(B):
191     index_bootstrap_sample = np.random.choice(n,n)
192     x_training = x[index_bootstrap_sample]
193     y_training = y[index_bootstrap_sample]
194     K_rbf = k_rbf(x_training[:, np.newaxis], x_training[np.newaxis, :], gamma_best)
195     alpha = np.linalg.solve((K_rbf + lamb_best_rbf * np.eye(n, n)), y_training)
196
197     y_predicted_bootstrap_rbf = []
198     for xi in x_fine:
199         y_predicted_bootstrap_rbf.append(np.sum(alpha * np.exp(-gamma_best*(xi-x_training)**2)))
200 bootstrap_predicted_rbf_matrix.append(y_predicted_bootstrap_rbf)
201 bootstrap_predicted_rbf_matrix = np.array(bootstrap_predicted_rbf_matrix)
202
203 percent_5_list_rbf = []
204 percent_95_list_rbf = []
205 for i in range(n_fine):
206     sorted_xi_from_300_B_sample = np.sort(bootstrap_predicted_rbf_matrix[:, i])
207     x_percentile_5 = sorted_xi_from_300_B_sample[int(B * 0.05)]
208     x_percentile_95 = sorted_xi_from_300_B_sample[int(B * 0.95)]
209     percent_5_list_rbf.append(x_percentile_5)
210     percent_95_list_rbf.append(x_percentile_95)
211
212
213 plt.plot(x_fine, y_fine_true, label = 'True Model')
214 plt.plot(x_fine, f_rbf_predicted, label = 'rbf Kernel Prediction')
215 plt.plot(x, y, 'bo', label = 'Observed data')
216 plt.fill_between(x_fine, percent_5_list_rbf, percent_95_list_rbf, alpha=0.3, label="90% CI")
217 plt.ylim(-6, 6)
218 plt.legend()
219 plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A3c_2_test.png")
220 plt.show()
221 #####

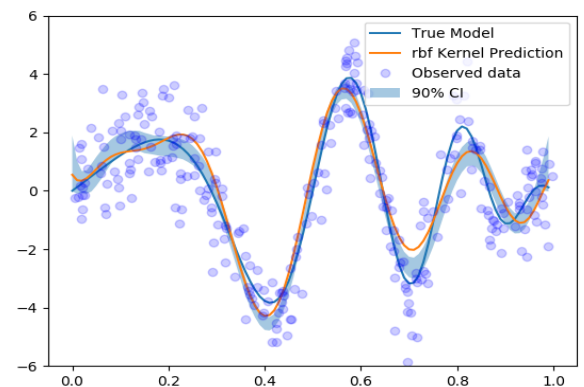
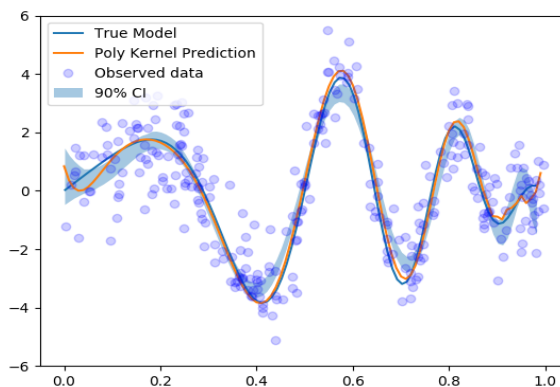
```


d.



Best lamb: 0.003814697265625 , Best d: 33

Best gamma for RBF kernel is : 8.992939495460687
Best Lambda for RBF kernel is : 1.862645149230957e-06



e.

5%: -0.14401336443720553

95%: -0.0253960255137911

In my case, 0 is not included in the confidence interval. And since the 5% and 95% are negative, so it is suggesting that with 90% confidence that RBF kernel is better than Polynomial kernel regression.

Code for question d,e,

```
1  # A3.d - 1 #####
2  import numpy as np
3  import matplotlib.pyplot as plt
4
5  n = 300
6  # np.random.seed(1)
7
8  x = np.random.uniform(0,1,n)
9  y = 4*np.sin(np.pi*x)*np.cos(6*np.pi*(x**2)) + np.random.standard_normal(n)
10
11 def k_poly(x, z, d):
12     a = x @ z.T
13     k = (1 + x @ z.T)**d
14     return k
15
```

```

16 error_validation_list = []
17 lamb = 500
18 cv_size = int(n/10)
19 lamb_list = []
20 d_list = []
21 for lamb in list(500 * (1/2)**(np.arange(0, 20))):
22     for d in list(range(0, 51)):
23         error_validation = 0
24         print("Lam: ", lamb, ", d: ", d)
25         for i in range(0, n, cv_size):
26             x_train = np.append(x[0:i], x[i+cv_size:n])
27             y_train = np.append(y[0:i], y[i+cv_size:n])
28             x_validation = x[i:i+cv_size]
29             y_validation = y[i:i+cv_size]
30             K = k_poly(x_train[:, np.newaxis], x_train[:, np.newaxis], d)
31             alpha = np.linalg.pinv(K + lamb) @ y_train
32             # in predicted y formula
33             k_xi_x = (1 + np.multiply(x_validation, np.repeat(x_train[np.newaxis, :],
34                 cv_size).reshape(270, cv_size))))*d
35             # y_predicted = alpha @ k_xi_x.T
36             y_predicted = alpha[np.newaxis, :] @ k_xi_x
37             error_validation += np.sum((y_predicted - y_validation).T @ (y_predicted - y_validation))
38             # error_validation = error_validation[0][0]
39             error_validation /= n
40             print("error_validation: ", error_validation)
41             error_validation_list.append(error_validation)
42             lamb_list.append(lamb)
43             d_list.append(d)
44
45 # min_error = min(error_validation_list)
46 index_bootstrap_sample_min_error = error_validation_list.index(min(error_validation_list))
47 lamb_best_poly = lamb_list[index_bootstrap_sample_min_error]
48 d_best = d_list[index_bootstrap_sample_min_error]
49 print("Best lamb: ", lamb_best_poly, ", Best d: ", d_best)
50
51 # Best lamb: 0.003814697265625 , Best d: 40
52
53 # plots the comparison
54 np.random.seed(1)
55 n = 100
56 x_fine = np.array(list(range(0, 100, 1))) / 100
57 y_fine_true = 4*np.sin(np.pi*x_fine)*np.cos(6*np.pi*(x_fine**2))
58 y_fine_grid = y_fine_true + np.random.standard_normal(n)
59 f_poly_predicted = []
60 for xi in x_fine:
61     K = k_poly(x_fine[:, np.newaxis], x_fine[:, np.newaxis], d_best)
62     alpha = np.linalg.pinv(K + lamb_best_poly) @ y_fine_grid
63     k_xi_x = (1 + xi * x_fine[np.newaxis, :]) ** d_best # use this when polynomial kernel
64     y_predicted = alpha @ k_xi_x.T
65
66     f_poly_predicted.append(y_predicted)
67
68 plt.plot(x_fine, y_fine_true, label='True')
69 plt.plot(x_fine, f_poly_predicted, label='Poly Kernel')
70 plt.plot(x, y, 'bo', label='Observed')
71 plt.legend()
72 plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A3d_1_test.png")
73 plt.show()

```

```

74
75 # ----- #
76
77 B = 300
78 n=300
79 m = 1000
80 x_fine = np.arange(min(x),max(x),0.01)
81 n_fine = len(x_fine)
82 # np.random.seed(10)
83 bootstrap_predicted_poly_matrix = []
84 x = np.random.uniform(0,1,n)
85 y_true_sample = 4*np.sin(np.pi*x)*np.cos(6*np.pi*(x**2))
86 y_observed = y_true_sample + np.random.randn(n)
87
88 for j in range(B):
89     index_bootstrap_sample = np.random.choice(n,n)
90     x_training = x[index_bootstrap_sample]
91     y_training = y_observed[index_bootstrap_sample]
92     K = k_poly(x_training[:,np.newaxis],x_training[:,np.newaxis], d_best)
93     alpha = np.linalg.solve((K + lamb_best_poly*np.eye(n, n)), y_training)
94     y_predicted_bootstrap_ploy = []
95     # for xi in np.array(list(range(0, 100, 1))) / 100:
96     # test
97     K = k_poly(x_training[:,np.newaxis],x_training[:,np.newaxis], d_best)
98     alpha = np.linalg.solve((K + lamb_best_poly*np.eye(n, n)), y_training)
99
100     for xi in x_fine:
101         y_predicted_bootstrap_ploy.append(np.sum((1+xi*x_training[np.newaxis,:]) ** d_best @ alpha))
102     bootstrap_predicted_poly_matrix.append(y_predicted_bootstrap_ploy)
103 bootstrap_predicted_poly_matrix = np.array(bootstrap_predicted_poly_matrix)
104
105 percent_5_list_poly = []
106 percent_95_list_poly = []
107 for i in range(n_fine):
108     sorted_xi_from_300_B_sample = np.sort(bootstrap_predicted_poly_matrix[:, i])
109     x_percentile_5 = sorted_xi_from_300_B_sample[int(B * 0.05)]
110     x_percentile_95 = sorted_xi_from_300_B_sample[int(B * 0.95)]
111     percent_5_list_poly.append(x_percentile_5)
112     percent_95_list_poly.append(x_percentile_95)
113
114 # x_fine = np.array(list(range(0, 100, 1))) / 100
115 y_fine_true = 4*np.sin(np.pi*x_fine)*np.cos(6*np.pi*(x_fine**2))
116 plt.plot(x_fine, y_fine_true, label = 'True Model')
117 plt.plot(np.array(list(range(0, 100, 1))) / 100, f_poly_predicted, label = 'Poly Kernel Prediction')
118 plt.fill_between(x_fine, percent_5_list_poly, percent_95_list_poly, alpha=0.4, label="90% CI")
119 plt.plot(x, y_observed,'bo', alpha=0.2, label = 'Observed data')
120 plt.ylim(-6, 6)
121 plt.legend()
122 plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A3e_1_test.png")
123 plt.show()
124
125
126
127 #####
128
129 # A3.d - 2
130
131 def k_rbf(x, z, gamma):

```

```

132     return np.exp(-gamma*(x-z)*(x-z))
133
134 n = 300
135 np.random.seed(1)
136 cv_size = int(n/10)
137 x = np.random.uniform(0,1,n)
138 y_true = 4*np.sin(np.pi*x)*np.cos(6*np.pi*(x**2))
139 y = y_true + np.random.randn(n)
140 error_validation_list = []
141 lamb_list = []
142 gamma_list = []
143 d_list = []
144
145 lamb = 1
146 for lamb in list(500 * (1/2)**(np.arange(0,30))):
147     for gamma in list(50 * (1/1.1)**(np.arange(0,30))):
148         print("Lam: ", lamb, ", gamma: ", gamma)
149         error_validation = 0
150         for i in range(0, n, cv_size):
151             x_train = np.append(x[0:i], x[i+cv_size:n])
152             y_train = np.append(y[0:i], y[i+cv_size:n])
153             x_validation = x[i:i+cv_size]
154             y_validation = y[i:i+cv_size]
155             K = k_rbf(x_train[:,np.newaxis],x_train[np.newaxis,:], gamma)
156             alpha = np.linalg.pinv(K + lamb) @ y_train
157             k_xi_x = np.exp(-gamma*(x_validation - np.repeat(x_train[np.newaxis,:],
158                 cv_size).reshape((n-cv_size, cv_size))))**2)
159             y_predicted = alpha[np.newaxis, :] @ k_xi_x
160             error_validation += np.sum((y_predicted - y_validation).T @ (y_predicted- y_validation))
161         error_validation /= n
162         error_validation_list.append(error_validation)
163         print("error_validation: ", error_validation)
164         lamb_list.append(lamb)
165         gamma_list.append(gamma)
166
167 min_error = min(error_validation_list)
168 index_bootstrap_sample_min_error = error_validation_list.index(min_error)
169 lamb_best_rbf = lamb_list[index_bootstrap_sample_min_error]
170 gamma_best = gamma_list[index_bootstrap_sample_min_error]
171 print('Best gamma for RBF kernel is : ', gamma_best)
172 print('Best Lambda for RBF kernel is :', lamb_best_rbf)
173
174 # Best gamma for RBF kernel is : 8.992939495460687
175 # Best Lambda for RBF kernel is : 1.862645149230957e-06
176 # plots the comparison
177 n = 100
178 np.random.seed(10)
179
180 x_fine = np.array(list(range(0, 100, 1))) / 100
181 y_fine_true = 4*np.sin(np.pi*x_fine)*np.cos(6*np.pi*(x_fine**2))
182 y_fine_grid = y_fine_true + np.random.standard_normal(n)
183
184 f_rbf_predicted = []
185 K_rbf = k_rbf(x_fine[:,np.newaxis],x_fine[np.newaxis,:], gamma_best)
186 alpha = np.linalg.solve((K_rbf + lamb_best_rbf*np.eye(n, n)), y_fine_grid)
187 for xi in x_fine:
188     f_rbf_predicted.append(np.sum(alpha * np.exp(-gamma_best*(xi-x_fine)**2)))
189

```

```

190 plt.plot(x_fine, y_fine_true, label = 'True Model')
191 plt.plot(x_fine, f_rbf_predicted, label = 'RBF Kernel Prediction')
192 plt.plot(x, y, 'bo', label = 'Observed data')
193 plt.legend()
194 plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A3d_2_test.png")
195 plt.show()
196
197
198 # A3.e2
199 B = 300
200 n=300
201 # m=1000
202 # n_fine = 100
203 np.random.seed(0)
204 bootstrap_predicted_rbf_matrix = []
205 # x = np.array(list(range(0, 100, 1))) / 100
206 x = np.random.uniform(0,1,n)
207 y_true_sample = 4*np.sin(np.pi*x)*np.cos(6*np.pi*(x**2))
208 y_observed = y_true_sample + np.random.randn(n)
209
210 for j in range(B):
211     index_bootstrap_sample = np.random.choice(n,n)
212     x_training = x[index_bootstrap_sample]
213     y_training = y_observed[index_bootstrap_sample]
214     K_rbf = k_rbf(x_training[:, np.newaxis], x_training[np.newaxis, :], gamma_best)
215     alpha = np.linalg.solve((K_rbf + lamb_best_rbf * np.eye(n, n)), y_training)
216
217     y_predicted_bootstrap_rbf = []
218     for xi in x_fine:
219         # k_xi_x = np.exp(-gamma * (xi - x_training[np.newaxis, :]) ** 2)
220         # predicted = k_xi_x@alpha
221         y_predicted_bootstrap_rbf.append(np.sum(alpha * np.exp(-gamma_best*(xi-x_training)**2)))
222     bootstrap_predicted_rbf_matrix.append(y_predicted_bootstrap_rbf)
223 bootstrap_predicted_rbf_matrix = np.array(bootstrap_predicted_rbf_matrix)
224
225 percent_5_list_rbf = []
226 percent_95_list_rbf = []
227 for i in range(len(x_fine)):
228     sorted_xi_from_300_B_sample = np.sort(bootstrap_predicted_rbf_matrix[:, i])
229     x_percentile_5 = sorted_xi_from_300_B_sample[int(B * 0.05)]
230     x_percentile_95 = sorted_xi_from_300_B_sample[int(B * 0.95)]
231     percent_5_list_rbf.append(x_percentile_5)
232     percent_95_list_rbf.append(x_percentile_95)
233
234 x_fine = np.array(list(range(0, 100, 1))) / 100
235 y_fine_true = 4*np.sin(np.pi*x_fine)*np.cos(6*np.pi*(x_fine**2))
236 plt.plot(x_fine, y_fine_true, label = 'True Model')
237 plt.plot(x_fine, f_rbf_predicted, label = 'rbf Kernel Prediction')
238 plt.fill_between(x_fine, percent_5_list_rbf, percent_95_list_rbf, alpha=0.4, label="90% CI")
239 plt.plot(x, y_observed, 'bo', alpha=0.2, label = 'Observed data')
240 plt.ylim(-6, 6)
241 plt.legend()
242 plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A3e_2_test.png")
243 plt.show()
244
245
246
247

```

248

249

250

251 ##### A3 e #####

252 #

253 #

A3 e

254 #

255 ##### A3 e #####

256 d_best=33

257

258 B = 300

259 n=300

260 m=1000

261 np.random.seed(4)

262 x = np.random.uniform(0,1,m+n)

263 y_true_sample = 4*np.sin(np.pi*x)*np.cos(6*np.pi*(x**2))

264 y_observed = y_true_sample + np.random.randn(m+n)

265 error_difference_list = np.zeros(B)

266

267 for j in range(B):

268 print("B: ", j)

269 index_bootstrap_sample = np.random.choice(m+n,m)

270 x_training = x[index_bootstrap_sample]

271 y_training = y_observed[index_bootstrap_sample]

272 # compute poly kernel

273 K = k_poly(x_training[:, np.newaxis], x_training[:, np.newaxis], d_best)

274 alpha = np.linalg.pinv(K + lamb_best_poly) @ y_training

275 # compute rbf kernel

276 K_rbf = k_rbf(x_training[:, np.newaxis], x_training[np.newaxis, :], gamma_best)

277 alpha_rbf = np.linalg.solve((K_rbf + lamb_best_rbf * np.eye(m, m)), y_training)

278

279 error_difference = 0

280 for i in range(len(x_training)):

281 # poly predicted

282 k_xi_x = (1 + x_training[i] * x_training[np.newaxis, :]) ** d_best # use this when polynomial k

283 y_predicted = alpha @ k_xi_x.T

284 # rbf predicted

285 y_predicted_rbf = np.sum(alpha_rbf * np.exp(-gamma_best*(x_training[i]-x_training)**2))

286 # error difference

287 error_difference += (y_training[i] - y_predicted)**2 - (y_training[i] - y_predicted_rbf)**2

288 error_difference /= len(x_training)

289 # error_difference_list.append(error_difference)

290 print("error_difference: ", error_difference)

291 error_difference_list[j] = error_difference

292 error_difference_list = np.sort(error_difference_list)

293

294 print("5%: ", error_difference_list[int(B * 0.05)])

295 print("95%: ", error_difference_list[int(B * 0.95)])

A.4

a.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from mnist import MNIST
4 import random
```

5

```

6  def load_data(size=1):
7      print("Loading Data!")
8      mndata = MNIST('python-mnist/data/')
9      X_train, labels_train = map(np.array, mndata.load_training())
10     X_test, labels_test = map(np.array, mndata.load_testing())
11     X_train = X_train / 255.0
12     X_test = X_test / 255.0
13
14     if (size != 1):
15         return X_train[:int(len(X_train)*size)], \
16             labels_train[:int(len(labels_train)*size)], \
17             X_test[:int(len(X_test)*size)], \
18             labels_test[:int(len(labels_test)*size)]
19     else:
20         return X_train, labels_train, X_test, labels_test
21 X_train, y_train, X_test, y_test = load_data(size=0.1)
22 print("Load Data Complete!")
23
24 def compute_cluster(X_train, centers, k):
25     objective_value = 0
26     clusters = [[] for i in range(k)]
27     for i in np.arange(len(X_train)):
28         distance_list = []
29         for center in centers:
30             norm = np.linalg.norm(X_train[i] - center)
31             distance_list.append(norm)
32         closed_center_index = distance_list.index(min(distance_list))
33         objective_value += min(distance_list)**2
34         clusters[closed_center_index].append(X_train[i])
35     return clusters, objective_value
36
37 def compute_centers(classes):
38     centers = []
39     for i in range(len(classes)):
40         centers.append(np.mean(classes[i], axis = 0))
41     return centers
42
43 objective_value_list = []
44
45 k = 10
46 objs = []
47 iteration = 0
48 old_centers = random.sample(list(X_train), k)
49 new_centers = random.sample(list(X_train), k)
50 while not np.array_equal(old_centers, new_centers):
51     iteration += 1
52     print("Iteration: ", iteration)
53     print("----- Compute Clusters -----")
54     old_centers = new_centers
55     clusters, objective_value = compute_cluster(X_train, new_centers, k)
56     print("objective_value: ", objective_value)
57     new_centers = compute_centers(clusters)
58     objective_value_list.append(objective_value)

```

b.

```

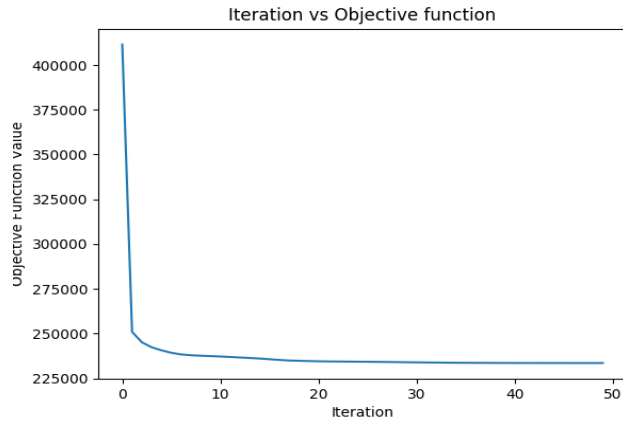
1  # plot center imgs
2  for i in range(k):

```

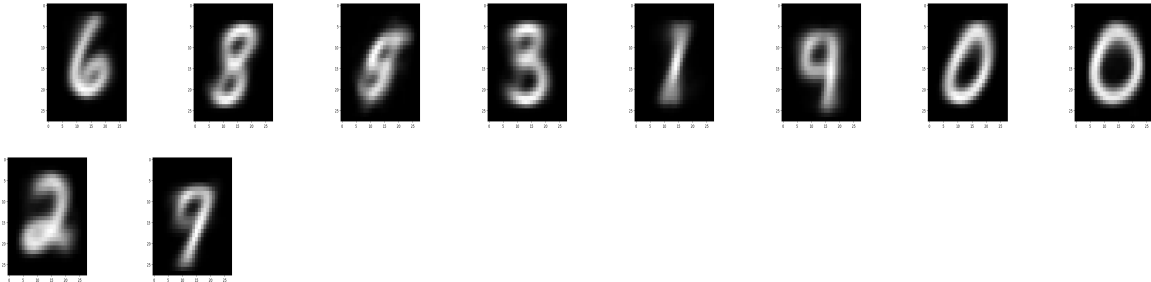
```

3 plt.imshow(new_centers[i].reshape((28, 28)))
4 plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A4/A4b_"+
5 str(i)+".png")
6 plt.show()

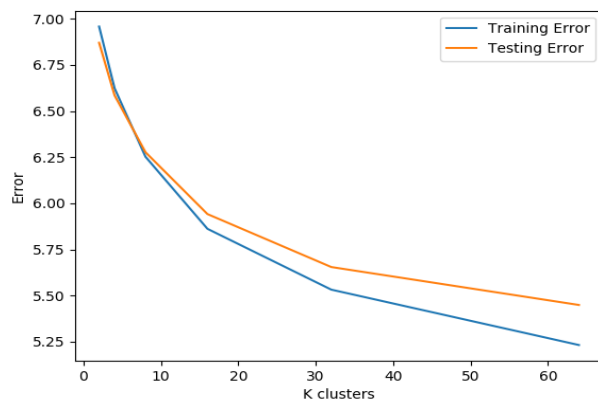
```



When $K = 10$, the centroids are:



c.



```

1 objective_value_list = []
2 training_error_list = []
3 testing_error_list = []
4 k_list = [2,4,8,16,32,64]
5
6 for k in k_list:
7     objs = []
8     iteration = 0
9     old_centers = random.sample(list(X_train), k)

```



```

10 new_centers = random.sample(list(X_train), k)
11 while not np.array_equal(old_centers, new_centers):
12     iteration += 1
13     print("Iteration: ", iteration)
14     print("----- Compute Clusters -----")
15     old_centers = new_centers
16     clusters, objective_value = compute_cluster(X_train, new_centers, k)
17     print("objective_value: ", objective_value)
18     new_centers = compute_centers(clusters)
19     objective_value_list.append(objective_value)
20 training_error = 0
21 for x_i in X_train:
22     distance_to_center = [np.linalg.norm(x_i - center) for center in new_centers]
23     training_error += min(distance_to_center)
24 training_error /= len(X_train)
25 training_error_list.append(training_error)
26
27 testing_error = 0
28 for x_i in X_test:
29     distance_to_center = [np.linalg.norm(x_i - center) for center in new_centers]
30     testing_error += min(distance_to_center)
31 testing_error /= len(X_test)
32 testing_error_list.append(testing_error)
33
34 plt.plot(k_list, training_error_list, label="Training Error")
35 plt.plot(k_list, testing_error_list, label="Testing Error")
36 plt.ylabel("Error")
37 plt.xlabel("K clusters")
38 plt.legend()
39 plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A4/A4c.png")
40 plt.show()

```

B.1

a.

From the hint, we get to know that:

$$1 - \epsilon \leq e^{-\epsilon}$$

Then based on for some $f_i \in \mathcal{F}$, we have $R(f_i) > \epsilon$,

$$1 - \epsilon \leq e^{-\epsilon}$$

$$1 - e^{-\epsilon} \leq \epsilon$$

$$1 - e^{-\epsilon} \leq \epsilon < R(f_i)$$

$$1 - e^{-\epsilon} \leq R(f_i)$$

$$e^{-\epsilon} - 1 \geq -R(f_i)$$

$$e^{-\epsilon} \geq 1 - R(f_i)$$

$$e^{-\epsilon} \geq 1 - E \left[\mathbf{1}_{\{f(x_1) \neq y_i\}} \right]$$

$$e^{-\epsilon} \geq 1 - P \left[\hat{R}_n(f_i) \neq 0 \right]$$

$$e^{-\epsilon} \geq P \left[\hat{R}_n(f_i) = 0 \right]$$

$$(e^{-\epsilon})^n \geq (P[\hat{R}_n(f_i) = 0])^n$$

$$e^{-n\epsilon} \geq P[\hat{R}_n(f) = 0]$$

b.

Suppose there are k f_i satisfy the requirement and Let event:

$$A_i = \{R(f_i) > \epsilon \text{ and } \hat{R}_n(f) = 0\}$$

And from previous question we know that:

$$Pr(\hat{R}_n(f) = 0) \leq e^{-n\epsilon}$$

So this works as well,

$$Pr(\hat{R}_n(f) = 0) * Pr(R(f_i) > \epsilon) \leq e^{-n\epsilon}$$

$$Pr(A_i) \leq e^{-n\epsilon}$$

$$Pr(A_1 \cup A_2 \cup A_3 \cdots \cup A_n) \leq \sum_{i=1}^{i=k} A_i$$

$$Pr(A_1 \cup A_2 \cup A_3 \cdots \cup A_n) \leq \sum e^{-n\epsilon}$$

$$Pr(A_1 \cup A_2 \cup A_3 \cdots \cup A_n) \leq ke^{-n\epsilon}$$

So,

$$\prod Pr(R(f_i) > \epsilon \text{ and } \hat{R}_n(f) = 0) \leq |\mathcal{F}|e^{-n\epsilon}$$

$$Pr(\exists f \in \mathcal{F} s.t R(f) > \epsilon \text{ and } \hat{R}_n(f) = 0) \leq |\mathcal{F}|e^{-n\epsilon}$$

c.

$$|\mathcal{F}|e^{-\epsilon n} \leq \delta$$

$$e^{-\epsilon n} \leq \frac{\delta}{|\mathcal{F}|}$$

$$-\epsilon n \leq \log\left(\frac{\delta}{|\mathcal{F}|}\right)$$

$$\epsilon n \geq \log\left(\frac{|\mathcal{F}|}{\delta}\right)$$

$$\epsilon \geq \frac{1}{n} \log\left(\frac{|\mathcal{F}|}{\delta}\right)$$

So,

$$\epsilon_{min} = \frac{1}{n} \log\left(\frac{|\mathcal{F}|}{\delta}\right)$$

d.

According to Chebyshev Inequality, for any random variable X and $\epsilon > 0$:

$$P(|X - E[X]| > \epsilon) \leq \frac{V[X]}{\epsilon^2}$$

$$P(|X - E[X]| \leq \epsilon) > \frac{V[X]}{\epsilon^2}$$

And here: X is $R(\hat{f})$,

From previous question we got:

$$\epsilon \geq \frac{1}{n} \log\left(\frac{|\mathcal{F}|}{\delta}\right)$$

So,

$$P(R(\hat{f}) - R(f^*) \leq \epsilon) \geq \frac{V[R(\hat{f})]}{\epsilon^2}$$

So now, we want to prove that:

$$\frac{V[R(\hat{f})]}{\epsilon^2} \geq 1 - \delta$$

So, here is what I will do, I will use the inequality from c:

$$\epsilon \geq \frac{1}{n} \log\left(\frac{|\mathcal{F}|}{\delta}\right)$$

$$(\epsilon)^2 \geq \left(\frac{1}{n} \log\left(\frac{|\mathcal{F}|}{\delta}\right)\right)^2$$

$$\frac{V[R(\hat{f})]}{\epsilon^2} \geq V[R(\hat{f})] / \left(\frac{1}{n} \log\left(\frac{|\mathcal{F}|}{\delta}\right)\right)^2$$

$$\frac{V[R(\hat{f})]}{\epsilon^2} \geq \frac{n^2 V[R(\hat{f})]}{(\log(|\mathcal{F}|) - \log(\delta))^2}$$

And the variance need to be noted, variance of distribution of the best function should be smaller than the variance of the original distribution, according to the variance of distribution of sample mean $\sigma_x^2 = \frac{\sigma_{x_i}^2}{n}$, and same idea applies here. However, it is not necessarily $1/n$ in this case, but I will just regard it this way and the inequality will still hold because we know that $V[R(f_i)]$ is bigger than $V[R(\hat{f})]$. And the reason I do this is because $V[R(\hat{f})]$ might be close to 0 and I don't want the numerator to risk to be 0, if the numerator is 0, then the proof failed. So, I will use the following:

$$V[R(\hat{f})] \approx \frac{V[R(f_i)]}{n}$$

So,

$$\frac{V}{\epsilon^2} \geq \frac{n^2 \frac{V[R(f_i)]}{n}}{(\log(|\mathcal{F}|) - \log(\delta))^2}$$

$$\frac{V}{\epsilon^2} \geq \frac{n V[R(f_i)]}{(\log(|\mathcal{F}|) - \log(\delta))^2}$$

Here, we know that $n > |\mathcal{F}|$ since $|\mathcal{F}|$ is finite, and n is infinity technically. And here $V[R(f_i)] \neq 0$. And since n is approximately infinity, and also the left hand side of the equation is positive, so:

$$n \frac{V[R(f_i)]}{(\log(|\mathcal{F}|) - \log(\delta))^2} \geq 1$$

$$n \frac{V[R(f_i)]}{(\log(|\mathcal{F}|) - \log(\delta))^2} \geq 1 - \delta$$

So, now we have proved that:

$$\frac{V[R(\hat{f})]}{\epsilon^2} \geq 1 - \delta$$

Then,

$$P(R(\hat{f}) - R(f^*) \leq \epsilon) \geq \frac{V[R(\hat{f})]}{\epsilon^2} \geq 1 - \delta$$

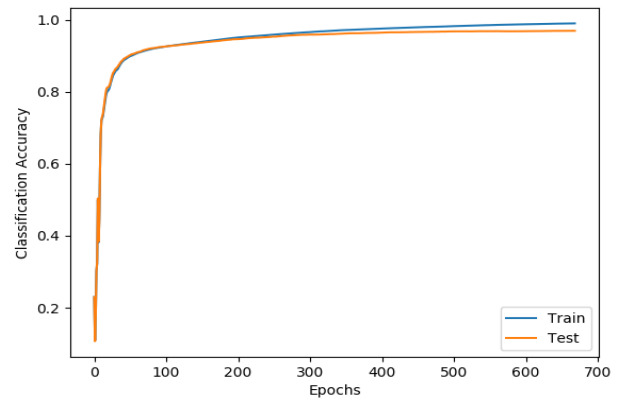
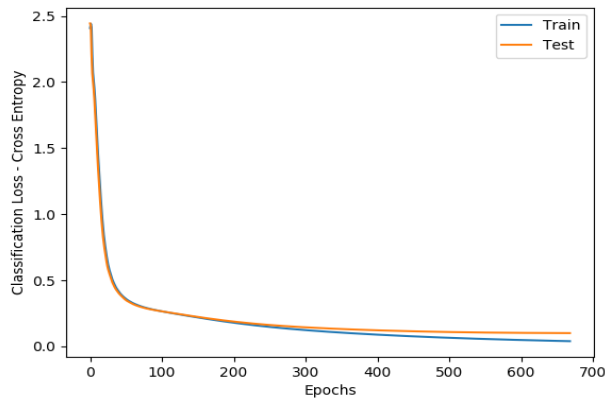
Then,

$$P(R(\hat{f}) - R(f^*) \leq \frac{\log(|\mathcal{F}|/\delta)}{n}) \geq \frac{V[R(\hat{f})]}{\epsilon^2} \geq 1 - \delta$$

A.5

a.

This part used a simple nn:



Train Accuracy: 0.99
Test Accuracy: 0.9701
Train Loss: (0.0408)
Test Loss: (0.1029)

```
1  import numpy as np
2  import torch
3  import matplotlib.pyplot as plt
4  from mnist import MNIST
5
6  def load_data(size=1):
7      print("Loading Data!")
8      mndata = MNIST('python-mnist/data/')
9      X_train, labels_train = map(np.array, mndata.load_training())
10     X_test, labels_test = map(np.array, mndata.load_testing())
11     X_train = X_train / 255.0
12     X_test = X_test / 255.0
13
14     if (size != 1):
15         return X_train[:int(len(X_train)*size)], \
16                labels_train[:int(len(labels_train)*size)], \
17                X_test[:int(len(X_test)*size)], \
18                labels_test[:int(len(labels_test)*size)]
19     else:
20         return X_train, labels_train, X_test, labels_test
21
22     # One hot encoding
23     def one_hot(y_train_, m):
24         n = len(y_train_)
25         reformed_tensor = torch.zeros(n, m)
26         for i in range(n):
27             index = y_train_[i]
28             reformed_tensor[i][index] = 1
29         return reformed_tensor
30
31     if __name__ == "__main__":
32         X_train, y_train, X_test, y_test = load_data(size=1)
33         print("Load Data Complete!")
34         # convert to tensor
```

```

35 dtype = torch.FloatTensor
36
37 X_train_ = torch.tensor(X_train, dtype=torch.double).type(dtype)
38 y_train_ = torch.tensor(y_train, dtype=torch.int64)
39
40 X_test_ = torch.tensor(X_test, dtype=torch.double).type(dtype)
41 y_test_ = torch.tensor(y_test, dtype=torch.int64)
42
43 def ReLU(x):
44     return abs(x) * (x > 0)
45
46 def model(x, w0, w1, b0, b1):
47     return (w1 @ ReLU(w0 @ x.T + b0) + b1).T
48
49 def alpha(d):
50     return 1 / np.sqrt(d)
51 h = 64
52 k = 10
53 n_train, d_train = X_train.shape
54 w0_data = (alpha(d_train) + alpha(d_train))*(torch.rand(h, d_train) -alpha(d_train)).type(dtype)
55 w0 = torch.autograd.Variable(w0_data, requires_grad=True)
56 b0_data = (alpha(d_train) + alpha(d_train))*(torch.rand(h, 1) -alpha(d_train)).type(dtype)
57 b0 = torch.autograd.Variable(b0_data,requires_grad=True)
58 w1_data = (alpha(d_train) + alpha(d_train))*(torch.rand(k, h) -alpha(d_train)).type(dtype)
59 n_test, d_test = X_test.shape
60 w1 = torch.autograd.Variable(w1_data,requires_grad=True)
61 b1_data = (alpha(d_train) + alpha(d_train))*(torch.rand(k, 1) -alpha(d_train)).type(dtype)
62 b1 = torch.autograd.Variable(b1_data, requires_grad=True)
63 step_size = 0.001
64 epochs = 150
65 train_accuracy_list = []
66 test_accuracy_list = []
67 loss_train_list = []
68 loss_test_list = []
69
70 optim = torch.optim.Adam([w0, w1, b0, b1], lr=step_size)
71
72 train_accu = 0
73 test_accu = 0
74 epochs_list = list(range(epochs))
75 # for epoch in epochs_list:
76 iter = 0
77 while train_accu < 0.99:
78     iter += 1
79     # print("Epoch: ", epoch)
80     optim.zero_grad()
81     y_hat = model(X_train_, w0, w1, b0, b1)
82     y_hat_index = torch.max(y_hat, dim=0).indices
83     loss = torch.nn.functional.cross_entropy(y_hat, y_train_)
84     loss.backward()
85     optim.step()
86     # Cross Entropy
87     max_index_train = torch.max(model(X_train_, w0, w1, b0, b1), dim=1).indices.numpy()
88     num_corrected_prediction_train = sum(max_index_train == y_train)
89     train_accu = num_corrected_prediction_train / len(y_train)
90     train_accuracy_list.append(train_accu)
91     loss_train_list.append(loss)
92

```

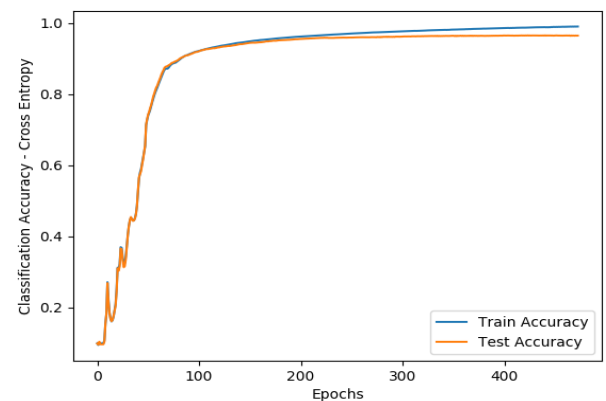
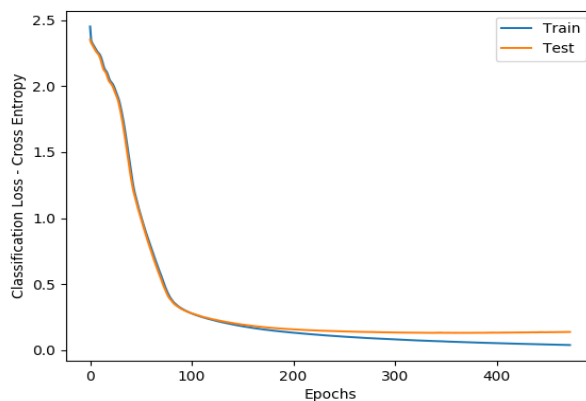
```

93     max_index_test = torch.max(model(X_test_, w0, w1, b0, b1), dim=1).indices.numpy()
94     num_corrected_prediction_test = sum(max_index_test == y_test)
95     test_accu = num_corrected_prediction_test / len(y_test)
96     test_accuracy_list.append(test_accu)
97     loss_test_list.append(torch.nn.functional.cross_entropy(model(X_test_, w0, w1, b0, b1), y_test_))
98     print("#####", iter, "#####")
99     print("CROSS ENTROPY: Train Accuracy: ", train_accu)
100    print("CROSS ENTROPY: Test Accuracy: ", test_accu)
101
102
103    print("CROSS ENTROPY: Train Accuracy: ", train_accuracy_list[-1])
104    print("CROSS ENTROPY: Test Accuracy: ", test_accuracy_list[-1])
105    print("Train Loss: ", loss_train_list[-1])
106    print("Test Loss: ", loss_test_list[-1])
107
108    plt.plot(range(iter), train_accuracy_list, label="Train")
109    plt.plot(range(iter), test_accuracy_list, label="Test")
110    plt.xlabel("Epochs")
111    plt.ylabel("Classification Accuracy - Cross Entropy")
112    plt.legend()
113    plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A5a.png")
114    plt.show()
115
116    w0.shape[0] * w0.shape[1] + w1.shape[0] * w1.shape[1] + len(b0) + len(b1)

```

b.

This part used a little more complex nn than last part:



Train Accuracy: 0.99005
 Test Accuracy: 0.9624
 Train Loss: (0.0409)
 Test Loss: (0.1531)

```

1  import numpy as np
2  import torch
3  import matplotlib.pyplot as plt
4  from mnist import MNIST
5
6
7  if __name__ == "__main__":
8      X_train, y_train, X_test, y_test = load_data(size=1)
9      print("Load Data Complete!")
10     # convert to tensor

```

```

11 dtype = torch.FloatTensor
12
13 X_train_ = torch.tensor(X_train, dtype=torch.double).type(dtype)
14 y_train_ = torch.tensor(y_train, dtype=torch.int64)
15
16 X_test_ = torch.tensor(X_test, dtype=torch.double).type(dtype)
17 y_test_ = torch.tensor(y_test, dtype=torch.int64)
18
19 def ReLU(x):
20     return abs(x) * (x > 0)
21
22 def model(x, w0, w1, b0, b1):
23     return (w1 @ ReLU(w0 @ x.T + b0) + b1).T
24
25 # In this question, use this 2 layers model.
26 def model_2_layers(x, w0, w1, w2, b0, b1, b2):
27     return (w2 @ ReLU(w1 @ ReLU(w0 @ x.T + b0) + b1) + b2).T
28
29 def alpha(d):
30     return 1 / np.sqrt(d)
31
32
33 def get_n_params(model):
34     pp = 0
35     for p in list(model.parameters()):
36         nn = 1
37         for s in list(p.size()):
38             nn = nn * s
39         pp += nn
40     return pp
41     h0 = 32
42     h1 = 32
43     h2 = 32
44     k = 10
45     n_train, d_train = X_train.shape
46     w0_data = (alpha(d_train) + alpha(d_train))*(torch.rand(h0, d_train) -alpha(d_train)).type(dtype)
47     w0 = torch.autograd.Variable(w0_data, requires_grad=True)
48     b0_data = (alpha(d_train) + alpha(d_train))*(torch.rand(h0, 1) -alpha(d_train)).type(dtype)
49     b0 = torch.autograd.Variable(b0_data,requires_grad=True)
50     w1_data = (alpha(d_train) + alpha(d_train))*(torch.rand(h1, h0) -alpha(d_train)).type(dtype)
51     n_test, d_test = X_test.shape
52     w1 = torch.autograd.Variable(w1_data,requires_grad=True)
53     b1_data = (alpha(d_train) + alpha(d_train))*(torch.rand(h1, 1) -alpha(d_train)).type(dtype)
54     b1 = torch.autograd.Variable(b1_data, requires_grad=True)
55
56     w2_data = (alpha(d_train) + alpha(d_train))*(torch.rand(k, h2) -alpha(d_train)).type(dtype)
57     n_test, d_test = X_test.shape
58     w2 = torch.autograd.Variable(w2_data,requires_grad=True)
59     b2_data = (alpha(d_train) + alpha(d_train))*(torch.rand(k, 1) -alpha(d_train)).type(dtype)
60     b2 = torch.autograd.Variable(b2_data, requires_grad=True)
61
62     step_size = 0.001
63     epochs = 150
64     train_accuracy_list = []
65     test_accuracy_list = []
66
67     optim = torch.optim.Adam([w0, w1, w2, b1, b0, b2], lr=step_size)
68

```

```

69 train_accu = 0
70 test_accu = 0
71 epochs_list = list(range(epochs))
72 loss_train_list = []
73 loss_test_list = []
74 # for epoch in epochs_list:
75 iter = 0
76 while train_accu < 0.99:
77     iter += 1
78     # print("Epoch: ", epoch)
79     optim.zero_grad()
80     y_hat = model_2_layers(X_train_, w0, w1, w2, b0, b1, b2)
81     y_hat_index = torch.max(y_hat, dim=0).indices
82     loss = torch.nn.functional.cross_entropy(y_hat, y_train_)
83     loss.backward()
84     optim.step()
85     # Cross Entropy
86     max_index_train = torch.max(model_2_layers(X_train_,
87     w0, w1, w2, b0, b1, b2), dim=1).indices.numpy()
88     num_corrected_prediction_train = sum(max_index_train == y_train)
89     train_accu = num_corrected_prediction_train / len(y_train)
90     train_accuracy_list.append(train_accu)
91     loss_train_list.append(loss)
92
93     max_index_test = torch.max(model_2_layers(X_test_,
94     w0, w1, w2, b0, b1, b2), dim=1).indices.numpy()
95     num_corrected_prediction_test = sum(max_index_test == y_test)
96     test_accu = num_corrected_prediction_test / len(y_test)
97     test_accuracy_list.append(test_accu)
98     loss_test_list.append(torch.nn.functional.cross_entropy(model_2_layers(X_test_,
99     w0, w1, w2, b0, b1, b2), y_test_))
100     print("#####", iter, "#####")
101     print("CROSS ENTROPY:
102     Train Accuracy: ", train_accu)
103     print("CROSS ENTROPY: Test Accuracy: ", test_accu)
104
105     print("CROSS ENTROPY: Train Accuracy: ", train_accuracy_list[-1])
106     print("CROSS ENTROPY: Test Accuracy: ", test_accuracy_list[-1])
107     print("Train Loss: ", loss_train_list[-1])
108     print("Test Loss: ", loss_test_list[-1])
109
110
111
112 plt.plot(range(iter), train_accuracy_list, label="Train Accuracy")
113 plt.plot(range(iter), test_accuracy_list, label="Test Accuracy")
114 # plt.plot(range(iter), loss_train_list, label="Train Loss")
115 # plt.plot(range(iter), loss_test_list, label="Test Loss")
116 plt.xlabel("Epochs")
117 plt.ylabel("Classification Accuracy - Cross Entropy")
118 plt.legend()
119 plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A5b.png")
120 plt.show()
121
122 # w0.shape[0] * w0.shape[1] + w1.shape[0] * w1.shape[1] + w2.shape[0] * w2.shape[1] + len(b0) + len(b1) + len(b2)

```


c.

In this first model, we have 50890 parameters.
In the section model, I have 26506 parameters.

```
1  # one layer
2  w0.shape[0] * w0.shape[1] + w1.shape[0] * w1.shape[1] + len(b0) + len(b1)
3  # two layer
4  w0.shape[0] * w0.shape[1] + w1.shape[0] * w1.shape[1] +
5  w2.shape[0] * w2.shape[1] + len(b0) + len(b1) + len(b2)
```

We can see that the one layer one has almost double the number of parameter. Interestingly they have similar accuracy. The two approach has similar results. One way to interpret this is that either more parameters or deeper network will give better fit to the data set. And in this case, we had similar accuracy because they are compensating each other. In this case, one has more parameter but shallower, on the other side, one has less parameters but deeper network.

A.6

a.

Solution:

Lambda 1: 5.116787728342091
Lambda 2: 3.7413284788648014
Lambda 10: 1.24272937641733
Lambda 30: 0.36425572027888947
Lambda 50: 0.16970842700672756
Sum of Lambdas: 52.72503549512679

$$\sum_{i=1}^d \lambda_i = 52.72$$

```
1  import torch
2  import numpy as np
3  import matplotlib.pyplot as plt
4  from mnist import MNIST
5  import random
6
7  def load_data(size=1):
8      print("Loading Data!")
9      mndata = MNIST('python-mnist/data/')
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
15     if (size != 1):
16         return X_train[:int(len(X_train)*size)], \
17             labels_train[:int(len(labels_train)*size)], \
18             X_test[:int(len(X_test)*size)], \
19             labels_test[:int(len(labels_test)*size)]
20     else:
21         return X_train, labels_train, X_test, labels_test
22 X_train, y_train, X_test, y_test = load_data(size=1)
23 print("Load Data Complete!")
24
25 # A6.a
26 n, d = X_train.shape
```

```

27 mu = np.zeros(d)
28 for i in range(d):
29     mu[i] = np.mean(X_train[:, i])
30
31 demean_X_train = X_train[:]
32 for row in range(X_train.shape[0]):
33     demean_X_train[row] = X_train[row] - mu
34 demean_X_test = X_test[:]
35 for row in range(X_test.shape[0]):
36     demean_X_test[row] = X_test[row] - mu
37
38 sigma = (demean_X_train.T @ demean_X_train) / 60000
39
40
41 U, S, V = np.linalg.svd(demean_X_train / np.sqrt(60000), False)
42 eigenvalues = S**2
43 print("Lambda 1: ", eigenvalues[0])
44 print("Lambda 2: ", eigenvalues[1])
45 print("Lambda 10: ", eigenvalues[9])
46 print("Lambda 30: ", eigenvalues[29])
47 print("Lambda 50: ", eigenvalues[49])
48 print("Sum of Lambdas: ", sum(eigenvalues))

```

b.

Solution:

Reconstruct:

V_k is top k eigenvectors.

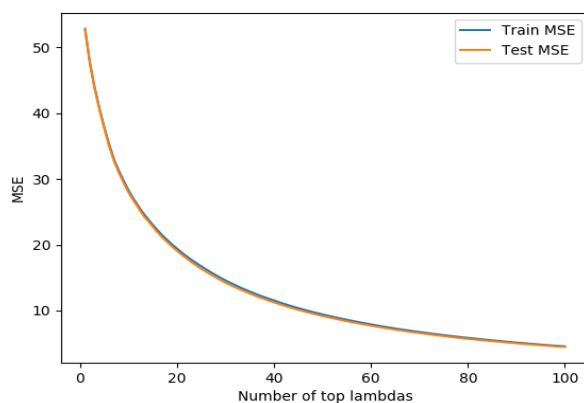
$$V_k V_k^T (x - \mu) + \mu$$

c.

Solution:

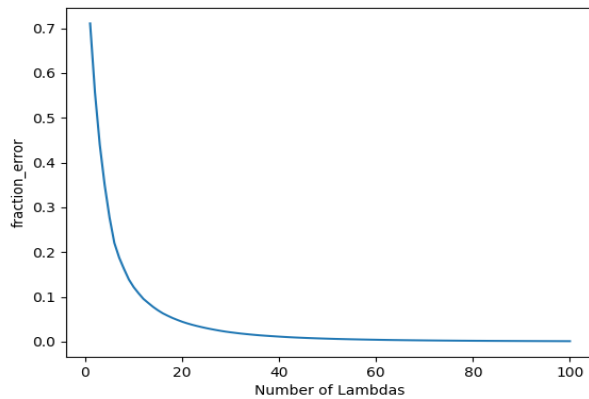
For reconstruction error (Mean Squared Error):

For $k = 1$ to 100, $d = 784$, plot $\frac{1}{n} \sum_{i=1}^n (\hat{X}_i - X_i)^2$



For fractional error:

For $k = 1$ to 100, $d = 784$, plot $1 - \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i}$



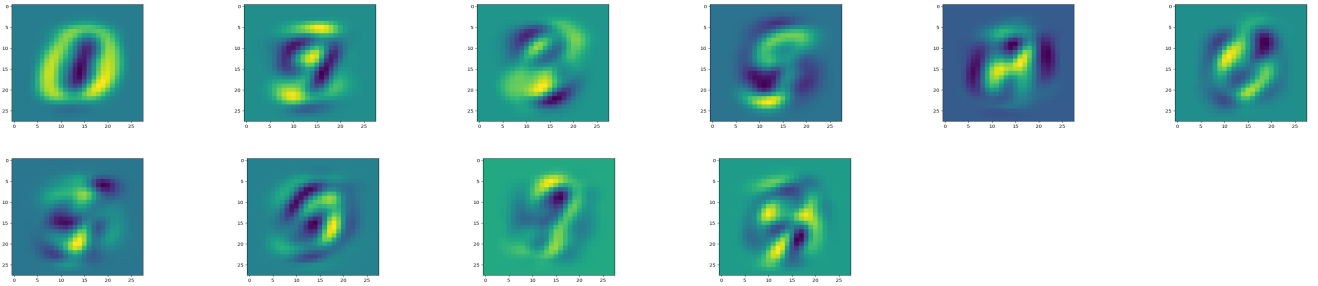
```

1  # plot mse
2  mean_squared_error_list_train = []
3  mean_squared_error_list_test = []
4  for k in range(0, 100):
5      print("K: ", k)
6      reconstructed = np.dot(V[:k, :].T.dot(V[:k, :]), demean_X_train[:, :].T).T
7      MSE_train = np.sum((reconstructed - demean_X_train) ** 2) / demean_X_train.shape[0]
8      mean_squared_error_list_train.append(MSE_train)
9
10     reconstructed_test = np.dot(V[:k, :].T.dot(V[:k, :]), demean_X_test[:, :].T).T
11     MSE_test = np.sum((reconstructed_test - demean_X_test) ** 2) / demean_X_test.shape[0]
12     mean_squared_error_list_test.append(MSE_test)
13
14     plt.plot(range(1, 101), mean_squared_error_list_train, label="Train MSE")
15     plt.plot(range(1, 101), mean_squared_error_list_test, label="Test MSE")
16     plt.xlabel("Number of top lambdas")
17     plt.ylabel("MSE")
18     plt.legend()
19     plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A6c_1.png")
20     plt.show()
21
22
23     # plot construction error
24     eigenvalues_sum = sum(eigenvalues)
25     fraction_error_list = []
26     for k in range(0, 100):
27         fraction_error = 1 - np.sum(eigenvalues[(k+1):]) / eigenvalues_sum
28         fraction_error_list.append(fraction_error)
29
30     plt.plot(range(1, 101), fraction_error_list)
31     plt.xlabel("Number of Lambdas")
32     plt.ylabel("fraction_error")
33     plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A6c_2.png")
34     plt.show()

```

d.

Solution:



Interpretation:

Each of the image represent a PCA direction. It is a feature space with the highest variation. We can capture most of the variation and information through the top directions and thus we can reconstruct fairly close data.

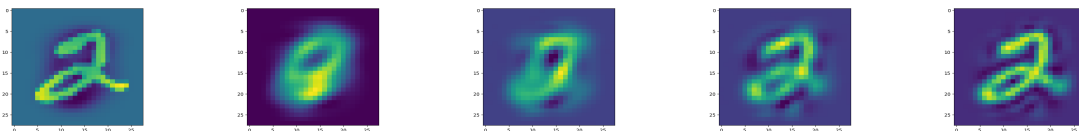
```
1 # ===== A6d =====
2 for k in range(10):
3     plt.imshow(V[k,:].reshape((28,28)))
4     plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A6d/A6d_"+
5                 str(k)+".png")
6     plt.show()
```

e.

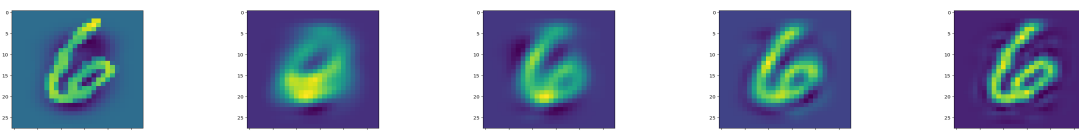
Solution:

Plot the comparison between different number of top eigenvalues for [5, 15, 40, 100]:

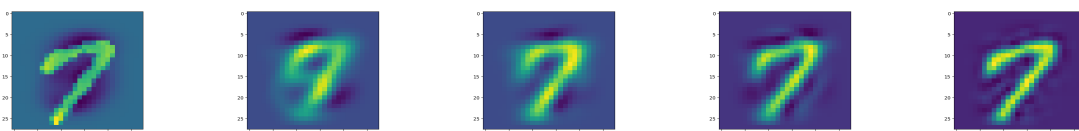
For $y = 2$:



For $y = 6$:



For $y = 7$:



Interpretation:

For this reconstruction process, we used 5, 15, 40, 100 top lambdas to reconstruct the image, and we can see that with more dimensionality we got better representation of the original image. So basically, more dimensionality we used, we could capture more variation in the data.

```
1 # ===== A6e =====
2 y_train[5] # =2
3 y_train[13] # =6
4 y_train[15] # =7
5
```

```

6  for y in [5, 13, 15]: # this is the index of y in the X_train
7      plt.imshow(X_train[y,:].reshape((28,28)))
8      plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A6e/A6e_"+
9          str(y)+".png")
10     plt.show()
11
12     for k in [5, 15, 40, 100]:
13         plt.imshow(np.dot(V[:,k,:].T.dot(V[:,k,:]), X_train[:, :].T.T[y].reshape((28, 28)))
14         plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw3/latex/plots/A6e/A6e_" +
15             str(y) + "_" + str(k) + ".png")
16         plt.show()
17
18     # test case
19     plt.imshow(np.dot(V[:,k,:].T.dot(V[:,k,:]), X_train[:, :].T.T[y].reshape((28, 28)))
20     plt.show()

```