

# CSE546 Homework 3

Chuanmudi Qin

May 28, 2020

## Problem A1

*Proof.* a. True

The new coordinate after PCA can be obtained by diagonalized the covariance matrix  $C = X^T X = P D P^{-1}$ , where  $P$  is the new coordinate who has  $\text{rank}(C) = d = \text{rank}(X)$  the full  $P$  matrix preserves all the information in  $X$ , therefore reconstruction error is 0.

□

*Proof.* b. False

It really depends on the data to see what are the most suitable linear classifiers.

□

*Proof.* c. True

Due to the fact that usually we randomly choosing  $B$  sets containing  $k$  samples, there are chances that  $x_i$  could occur multiple times

□

*Proof.* d. False, under the premise that each **column** of  $X$  is a observation. the matrix  $U$  should be the eigenvectors

□

*Proof.* e. True

noises in the new coordinates are associated with the directions of smaller eigenvalues. Denoise will be helpful for extracting useful information

□

*Proof.* f. False.

Supposed that we have  $n$  data points and  $k$  classe in the original data set. We could take  $n$  classes to get zero error but the result will not be meaningfull. □

*Proof.* g. I should decrease  $\sigma$

□

## Problem A2

*Proof.* known that:

$$\begin{aligned} K(x, x') &= e^{-\frac{(x-x')^2}{2}} = e^{\frac{-x^2+2xx'-x'^2}{2}} = e^{\frac{-x^2-x'^2}{2}} * e^{-2xx'} \\ &= e^{\frac{-x^2-x'^2}{2}} \left(1 + \frac{2xx'}{1!} + \frac{(2xx')^2}{2!} + \frac{(2xx')^3}{3!} + \dots\right) \\ &= e^{\frac{-x^2-x'^2}{2}} \left(1 * 1 + \sqrt{\frac{2}{1!}}x * \sqrt{\frac{2}{1!}}x' + \sqrt{\frac{2}{2!}}x^2 * \sqrt{\frac{2}{2!}}x'^2 + \sqrt{\frac{2}{3!}}x^3 * \sqrt{\frac{2}{3!}}x'^3 + \dots\right) \\ &= \langle e^{-x^2} [\sqrt{\frac{2}{1!}}x, \sqrt{\frac{2}{2!}}x^2, \sqrt{\frac{2}{3!}}x^3, \dots], e^{-x'^2} [\sqrt{\frac{2}{1!}}x', \sqrt{\frac{2}{2!}}x'^2, \sqrt{\frac{2}{3!}}x'^3, \dots] \rangle \\ &= \phi(x)^T \phi(x') \end{aligned}$$

□

### Problem A3

*Proof.* A3a)

1) For rbf kernel the hyperparameters are chose from  $d = [.01, .1, 1.10, 20, 50, 90, 100, 105, 200, 500]$

and  $\lambda = np.arange(0, 1, 0.05)$

the optimal lambda value is 0.19, the optimal d value is 100

2) For polynomial kernel the hyperparameters are chosen from  $d = np.arange(10, 90, 1)$

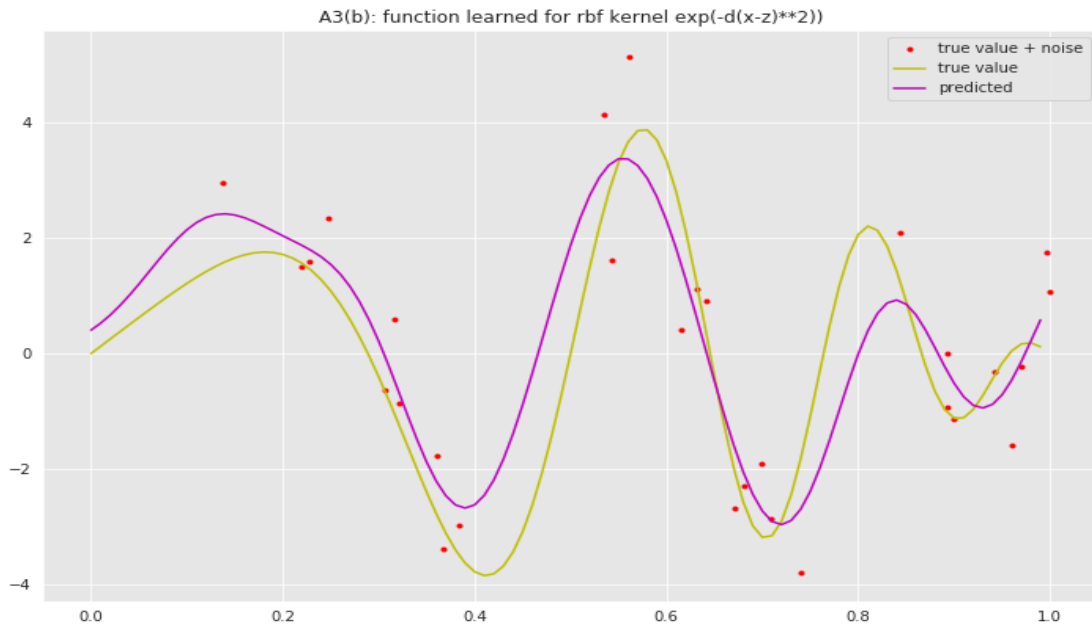
and  $\lambda = np.arange(0.01, .5, .05)$

the optimal lambda value is 0.46, the optimal d value is 45

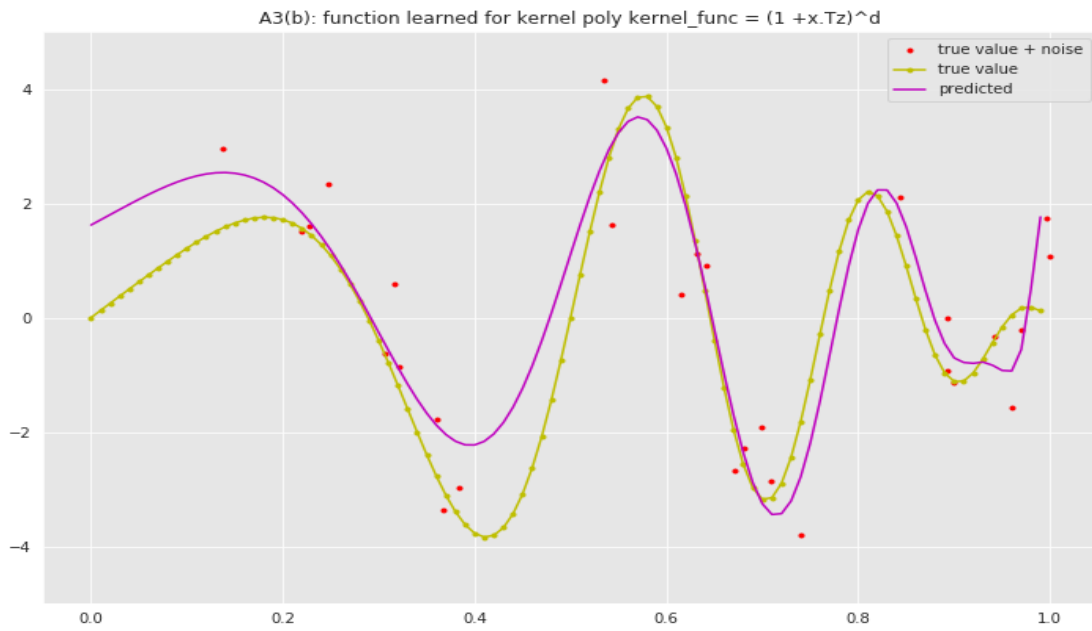
□

*Proof.* A3b)

1) for RBF kernel,  $f(x)$  and  $\hat{f}(x)$  are plotted as follow:



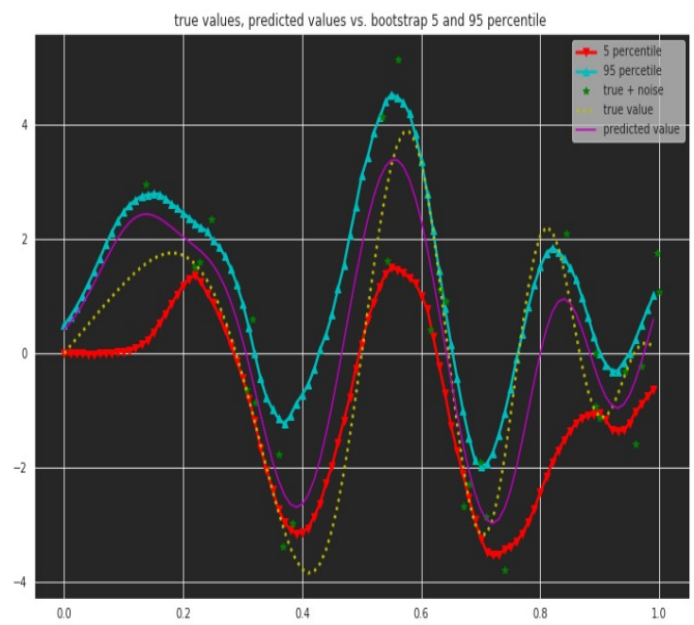
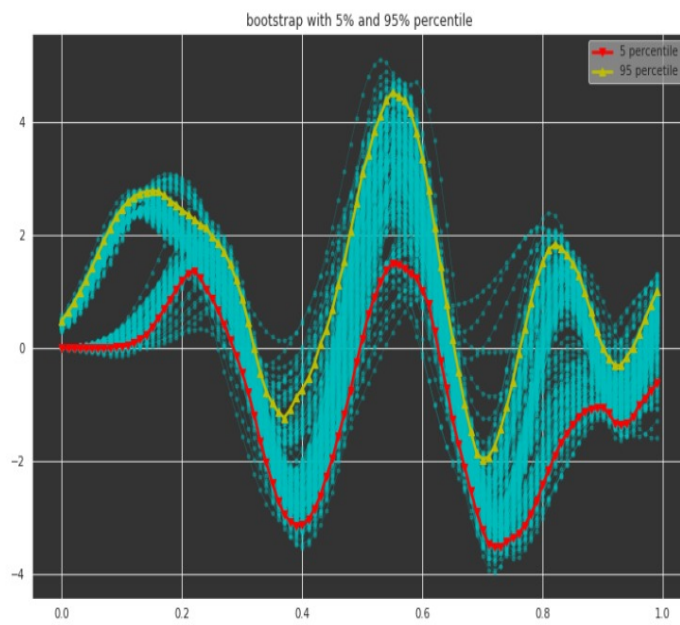
1) for polynomial kernel,  $f(x)$  and  $\hat{f}(x)$  are plotted as follow:



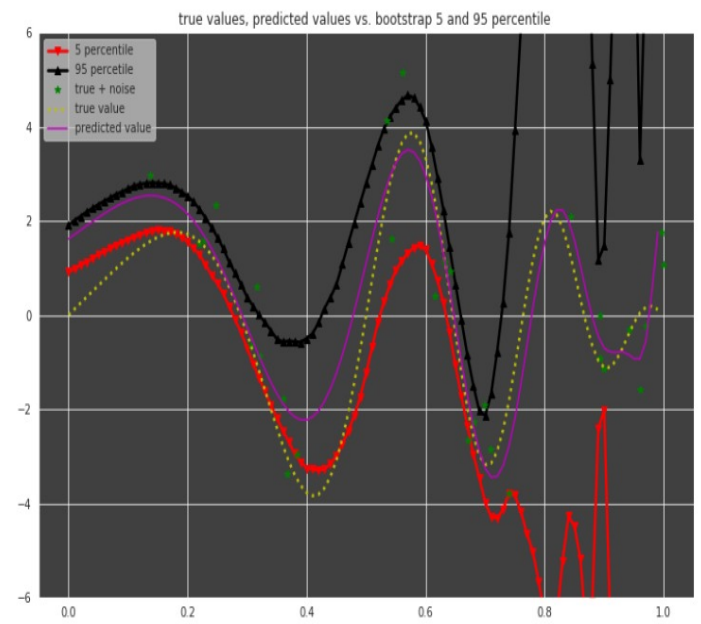
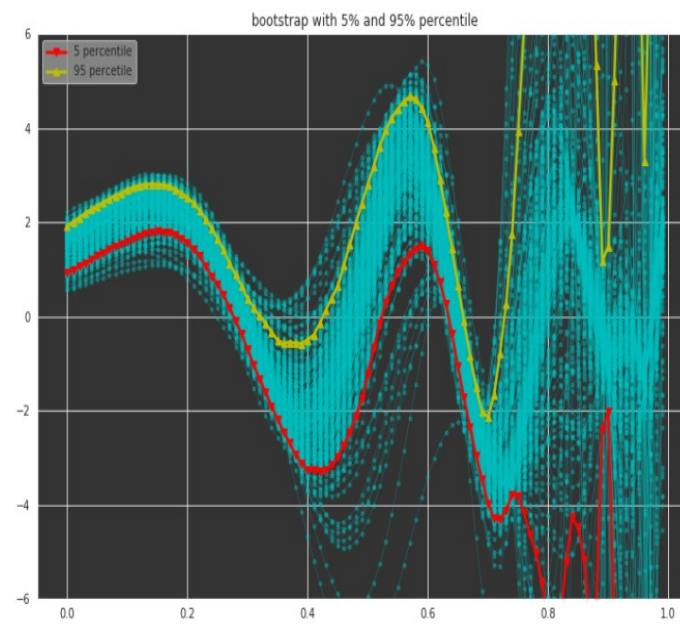
□

*Proof.* A3c)

1) for RBF kernel:



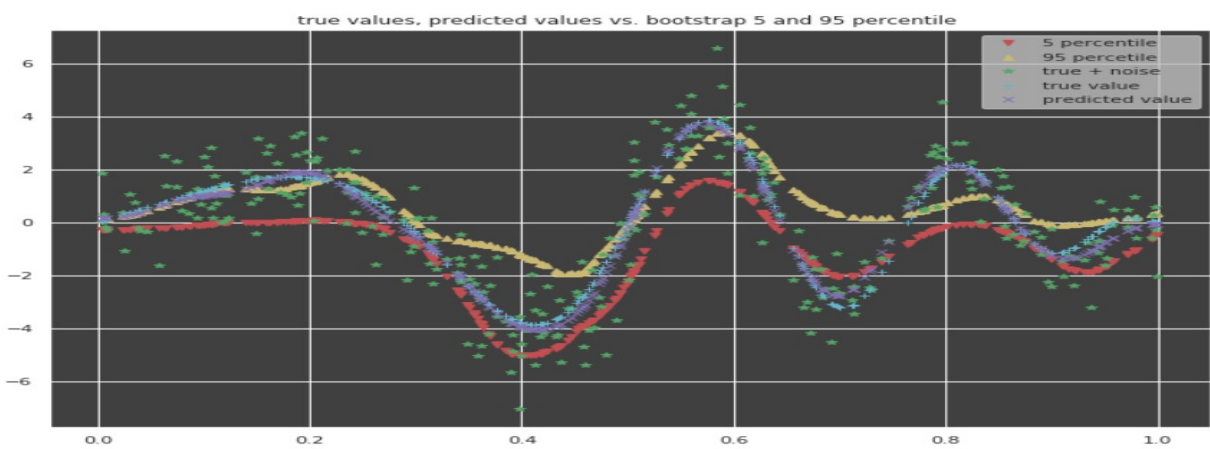
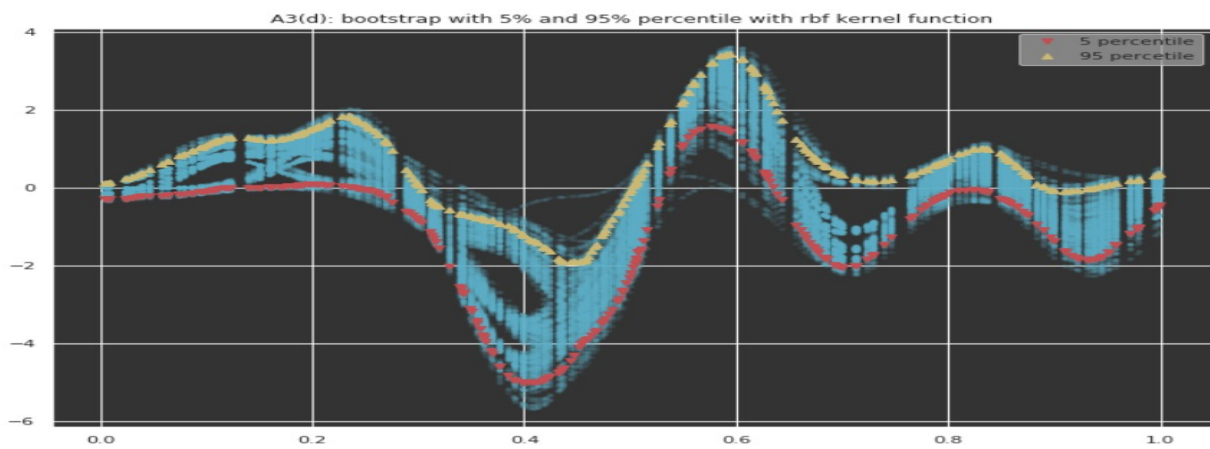
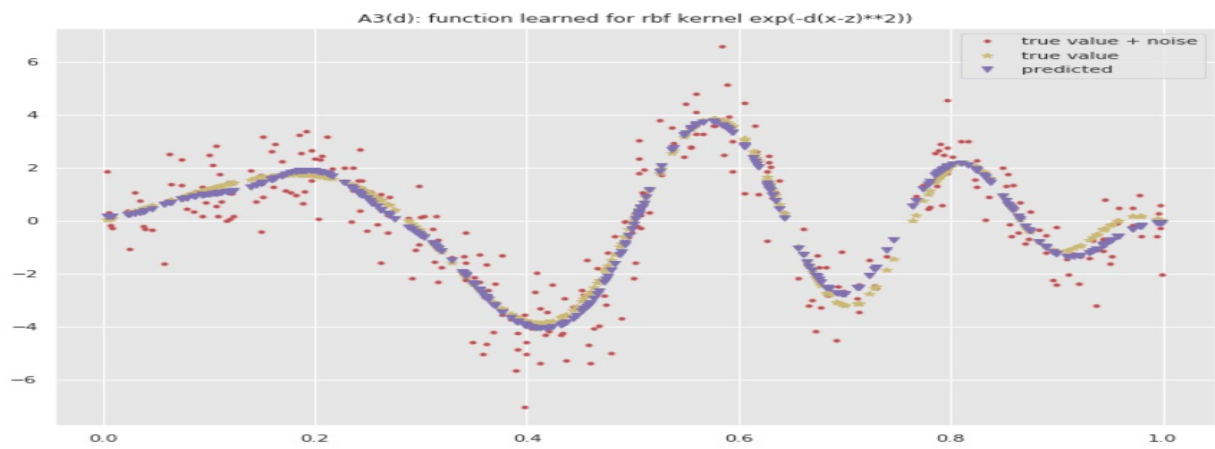
1) for polynomial kernel:



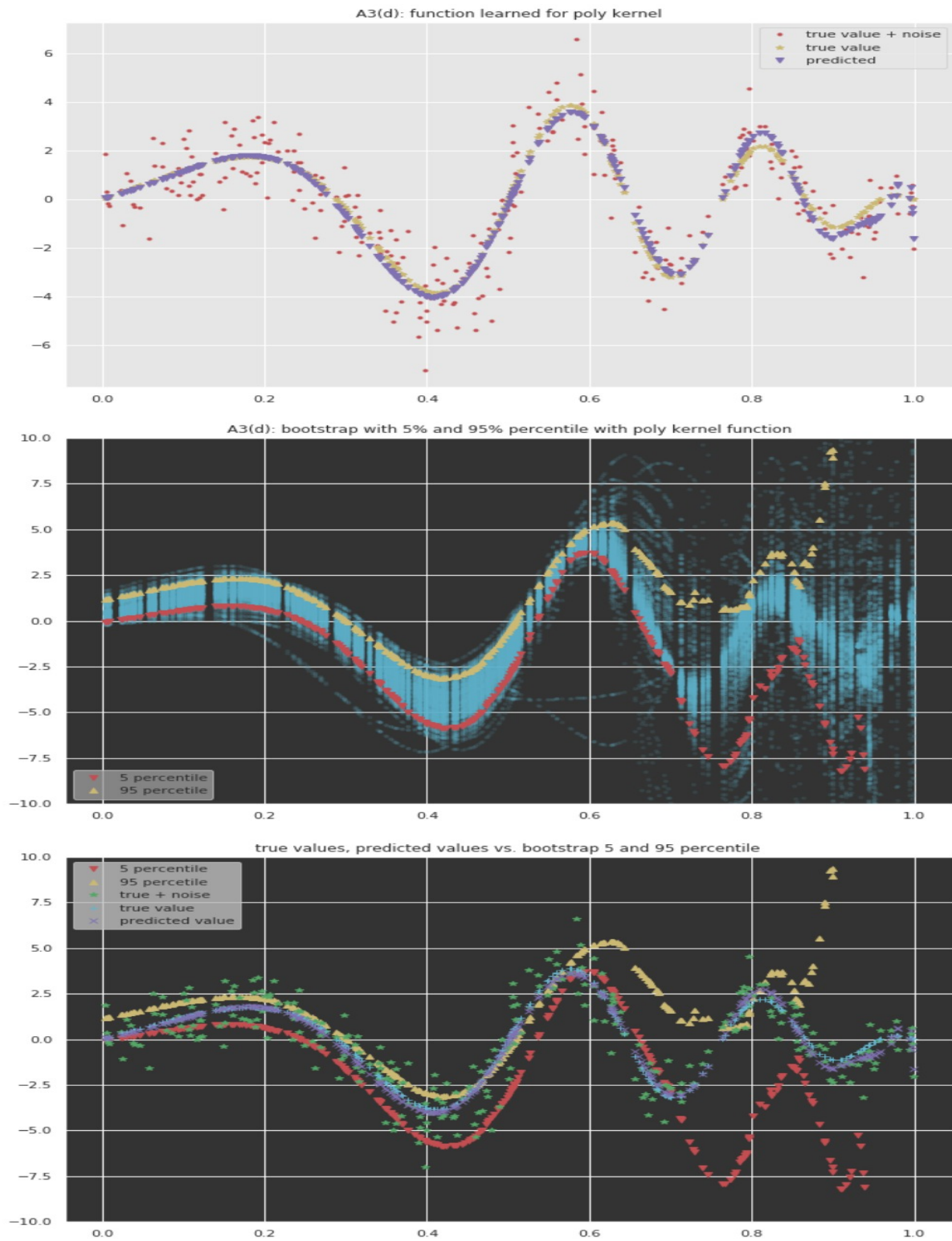
□

Proof. A3d)

1) RBF kernel:



1) Polynomial kernel:



□

*Proof.* A3e)

the 5% and 95% percentile is

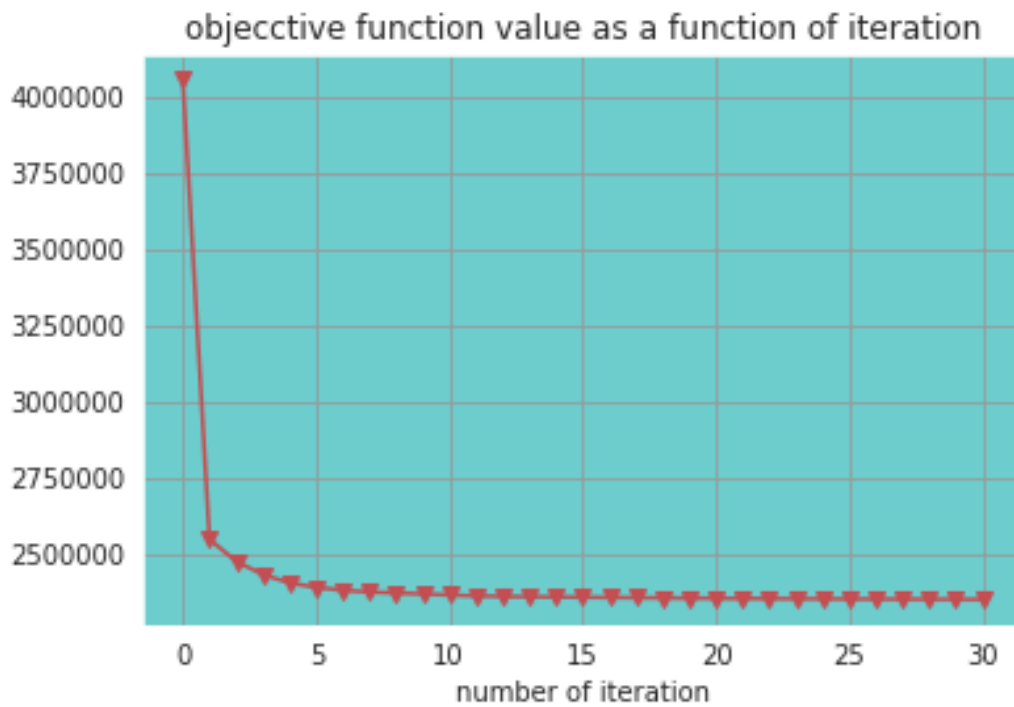
$$(0.003351860663601623, 0.07656613032802137)$$

the interval does not contain 0 which suggests that under the current optimal hyper parameters, polynomial has a better performance. However, we can not say for sure that which one of the  $f_{rbf}$  or  $f_{poly}$  is better, since hyperparameters of one might be tuned better than the other, which means it might not be a fair criteria to judge which one is better.  $\square$

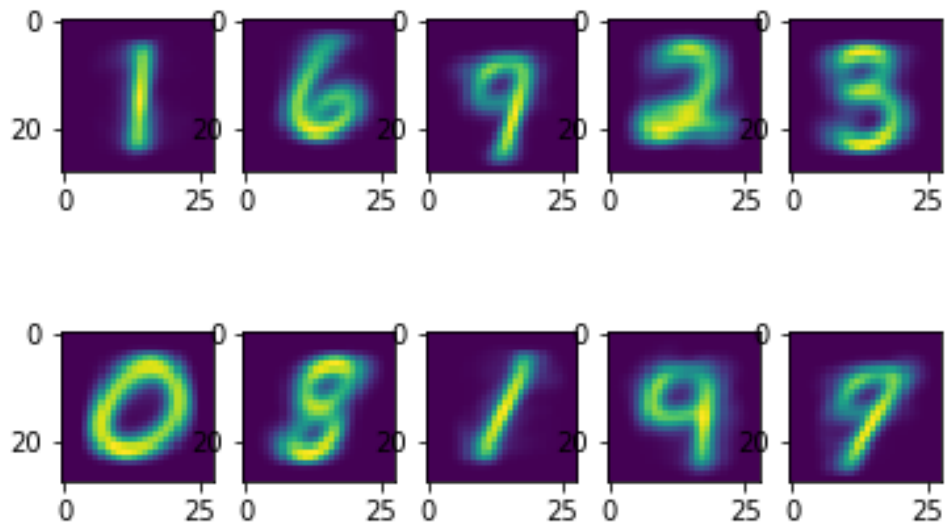


## Problem A4

*Proof.* A4b) I randomly chose initial centroids. and the result are as follow:



The final centroids are:



□

*Proof.* A4c)

The training error and testing error as a function of  $k$  will be:



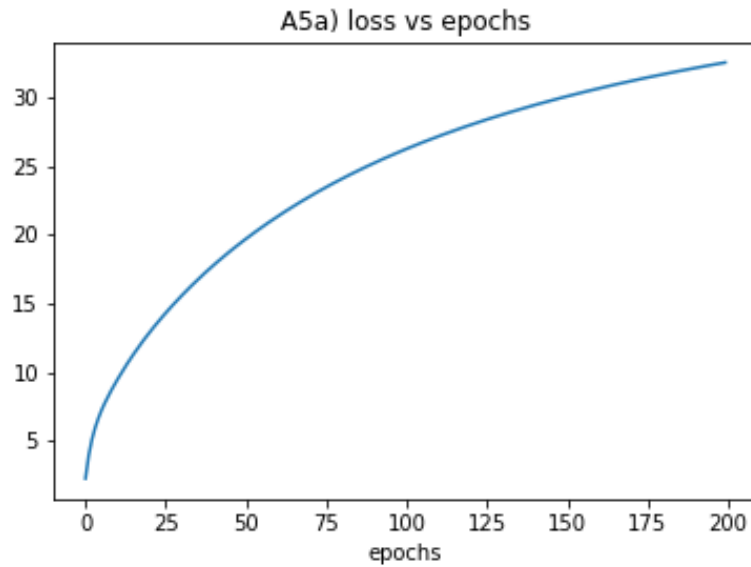
□

## Problem A5

*Proof.* A5a) learning rate is .01

The final training accuracy is 0.99005

The final testing final accuracy is :0.968

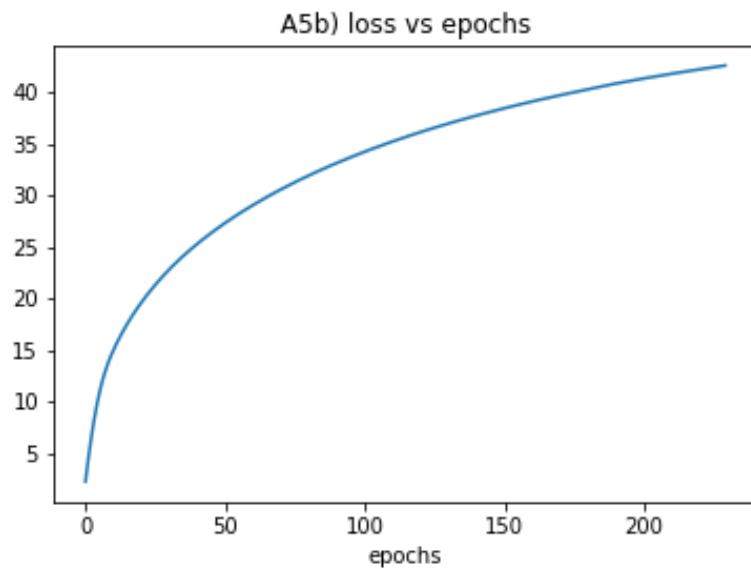


□

*Proof.* A5b) learning rate is .01

The final training accuracy is 0.9904833333333334

The final testing final accuracy is 0.9595



□

*Proof.* A5c)

The shallower one takes approximately 198 iterations to converge, with a learning rate of 0.01. parameters count for each epoch is  $784*64 + 64 + 64*10 + 10 = 50890$ .

The deeper one takes approximately 243 iterations to converge with a learning rate of .01. The parameters count for this is  $784*32 + 32*32 + 32*10 + 32+32+10 = 26505$ .

I would say the second one, which is the one with more layers is better. Due to the fact that, it reaches as good accuracy as the first one but is more memory efficient. For every epoach, it stores less coefficient and it reaches to the target accuracy faster. □

## Problem A6

*Proof.* A6a)

$\lambda_k$  represents how much variance of the data are explained in the direction of the  $k^{th}$  eigenvector.

For exaple,  $\lambda_1$  is telling how much variabce is explained by the most significant principle component.  $\lambda_2$  is telling how much variabce is explained by the second most significant principle component,....

$\sum_{i=1}^k \lambda_i$  means how much variance in the data are explained by the first  $k$  eigenvectors. □

*Proof.* A6b)

It depends on whether X's columns are the observations or rows are the observations.

For this problem, I am defining X to be:

$$X = [x_1, x_2, x_3, \dots] = \begin{pmatrix} | & | & | & & | \\ x_1 & x_2 & x_3 & \dots & x_{60000} \\ | & | & | & & | \end{pmatrix}$$

where each column of X is a observation

$$[U, S, V] = \text{svd}(X)$$

where  $U_{d-by-d}$  and each column of  $U$  is the eigenvector. Columns in  $U$ , in this case, is the new coordinates.

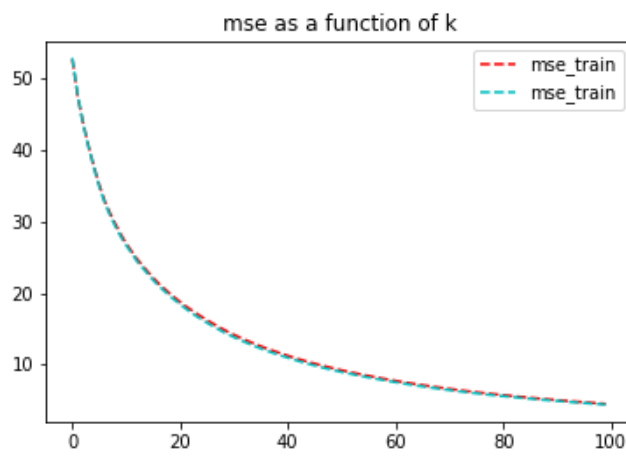
Using the first  $r$  vectors to reconstruct:

$$\hat{Y}_r = U[:, : r](U[:, : r])^T X$$

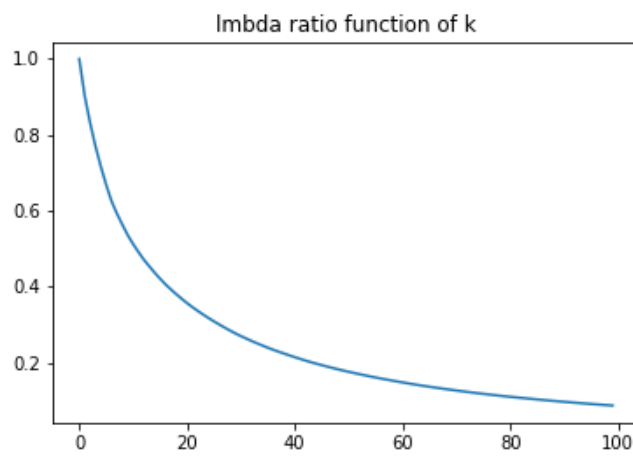
□

*Proof.* A6c)

MSE as a function of  $k$  is plotted as follow:



the ratio as a function of  $k$  is plotted as follow:

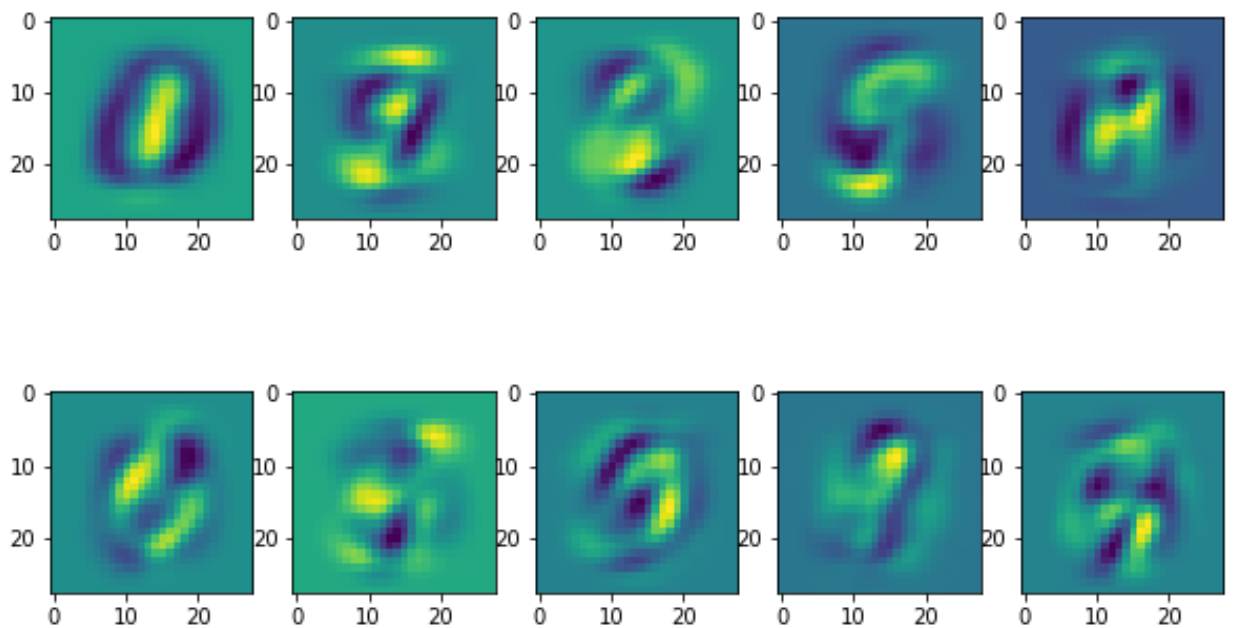


□

*Proof.* A6d)

the eigenvectors captures the most prominent/obvious differences of all images. For example, eigenvector that associated with the largest eigenvalue is the images that captures the pixels which are different the most among all the images. The brighter the pixels are, the more different all the images varies on them.

first 10 eigenVectors as images



□

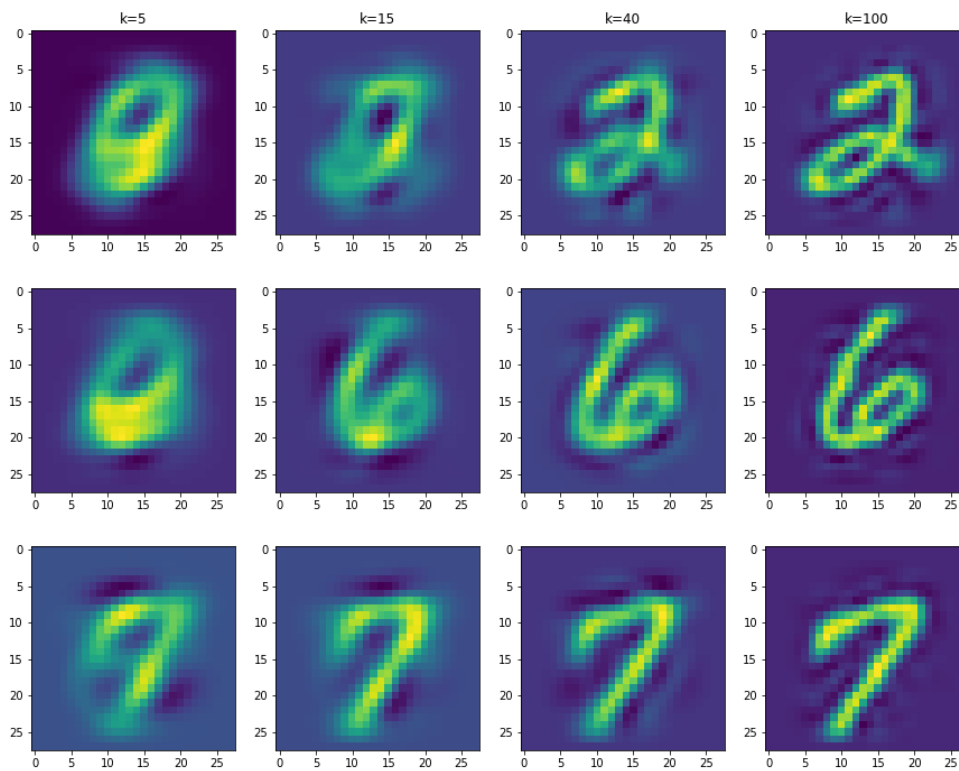


*Proof.* A6e)

As we already known, the eigenvectors in  $U$  are othogonal and span  $R^d$ , therefore, we have a new coordinates/basis, where we can represent everything in  $R^d$  with the linear combination of the set of eigenvectors

We can see that the images in the first column is the reconstructed with only the first eigenvector. the  $k=15$  is reconstructed with the linear combination of the first 15 eigenvectors. Same argument with  $k=40$  and  $k = 100$ , the image is the linear combination of the first 40 ir 100 eigenvectors

reconstruction of '2','6','7', using  $k=5,15,40,100$



□

## Problem B1

*Proof.* a)

To show that  $P(\widehat{R}_n(f) = 0) \leq e^{-n\sigma}$ , First I realize that :

$$\begin{aligned} E_{XY}[1\{f(X) \neq Y\}] &= E_X[E_{Y|X}[1\{f(X) \neq Y\}|X = x]] \\ &= E_{Y|X}[1\{f(X) \neq Y\}|X = x] \\ &= 1 - P(\{f(X) = Y\}|X = x) \geq \sigma \end{aligned}$$

which implies that:  $P(\{f(X) = Y\}|X = x) \leq 1 - \sigma$

$$\begin{aligned} P(\widehat{R}_n(f) = 0) &= P\left(\sum_{i=1}^n 1\{f(x_i) \neq y_i\}\right) \\ &= P(\{f(x_1) = y_1\}|X = x_1) * P(\{f(x_2) = y_2\}|X = x_2) * \dots * P(\{f(x_n) = y_n\}|X = x_n) \\ &\leq (1 - \sigma)^n \\ &= e^{-n\sigma} \end{aligned}$$

□

*Proof.* b)

Supposed that the set  $\mathcal{F} = \mathcal{A} \cup \mathcal{D}$ , where  $\mathcal{A} \cap \mathcal{D} = \emptyset$  and  $\mathcal{A} = \{f \in \mathcal{F} | R(f) > \sigma, \widehat{R}_n(f) = 0\}$

noticed  $\mathcal{A} \subset \mathcal{F} \implies |\mathcal{A}| \leq |\mathcal{F}|$  then our objective:

$$\begin{aligned} P(f \in \mathcal{F} \text{ s.t. } R(f) > \sigma, \widehat{R}_n(f) = 0) \\ &= P(\mathcal{A}) = P(\{f_1 | f_1 \in \mathcal{A}\} \cup \{f_2 | f_2 \in \mathcal{A}\} \cup \dots \cup \{f_k | f_k \in \mathcal{A}\}) \\ &= |\mathcal{A}|e^{-n\sigma} \leq |\mathcal{F}|e^{-n\sigma} \end{aligned}$$

□

*Proof.* c)

$$\begin{aligned} |\mathcal{F}|e^{-n\sigma} &\leq \delta \\ \implies \log(e^{-n\sigma}) &\leq \log\left(\frac{\delta}{|\mathcal{F}|}\right) \\ \implies \sigma &\geq -\frac{\log\left(\frac{\delta}{|\mathcal{F}|}\right)}{n} \\ \implies \sigma &\geq \frac{\log\left(\frac{|\mathcal{F}|}{\delta}\right)}{n} \end{aligned}$$

therefore,  $\sigma_{min} = \frac{\log(\frac{|\mathcal{T}|}{\delta})}{n}$

□

*Proof.* B1d)

Let  $M = |\mathcal{F}|$ , and  $c$  to be a number

suppose I have 2 sets  $A = \{R(\hat{f}) < c\}$  and  $B = \{R(\hat{f}) - R(f^*) < c\}$

I know that  $R(f^*) \geq 0$ , then

$$(*) = R(\hat{f}) - R(f^*) \leq R(\hat{f}) \leq c$$

To prove that  $c = \frac{\log(M/\delta)}{n}$ , I claim that  $P(R(\hat{f}) \leq c) > 1 - \delta$ .

This implies that:  $P(R(\hat{f}) > c) \leq 1 - \delta$

As has been proven in part (b),

$$\begin{aligned} P(R(\hat{f}) > c) &= P(\cup_{i=1}^M \{\hat{R}(f_i) = 0, R(f_i) > c\}) \\ &\leq \sum P(\{\hat{R}(f_i) = 0, R(f_i) > c\}) \\ &\leq M e^{-nc} = \delta \end{aligned}$$

As proven in part(c),

$$(**) = c_{min} = \frac{\log(\frac{|\mathcal{F}|}{\delta})}{n}$$

combining (\*) and (\*\*), I know

$$R(\hat{f}) - R(f^*) \leq R(\hat{f}) \leq \frac{\log(\frac{|\mathcal{F}|}{\delta})}{n}$$

□

Collaborator: Zidan Luo