

# HOMWORK 5

Lucas Poon  
llpoon

**Instructions:** Use this latex file as a template to develop your homework. Submit your homework on time as a single pdf file. Please wrap your code and upload to a public GitHub repo, then attach the link below the instructions so that we can access it. Answers to the questions that are not within the pdf are not accepted. This includes external links or answers attached to the code implementation. Late submissions may not be accepted. You can choose any programming language (i.e. python, R, or MATLAB). Please check Piazza for updates about the homework. It is ok to share the experiments results and compare them with each other.

Github: <https://github.com/poonlucas/CS760>

## 1 Clustering

### 1.1 K-means Clustering (14 points)

1. **(6 Points)** Given  $n$  observations  $X_1^n = \{X_1, \dots, X_n\}$ ,  $X_i \in \mathcal{X}$ , the K-means objective is to find  $k (< n)$  centres  $\mu_1^k = \{\mu_1, \dots, \mu_k\}$ , and a rule  $f: \mathcal{X} \rightarrow \{1, \dots, K\}$  so as to minimize the objective

$$J(\mu_1^K, f; X_1^n) = \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}(f(X_i) = k) \|X_i - \mu_k\|^2 \quad (1)$$

Let  $\mathcal{J}_K(X_1^n) = \min_{\mu_1^K, f} J(\mu_1^K, f; X_1^n)$ . Prove that  $\mathcal{J}_K(X_1^n)$  is a non-increasing function of  $K$ .

$$\mathcal{J}_K(X_1^n) = \min_{\mu_1^K, f} \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}(f(X_i) = k) \|X_i - \mu_k\|^2$$

Using induction:

Assume that  $\mathcal{J}$  has been non-increasing up to some  $k$ . Now we need to show that  $\mathcal{J}$  is non-increasing up to some  $k + 1$ .

For  $k + 1$ , we add a new cluster centered at some arbitrary point. Since K-means has not converged yet,  $\mathcal{J}$  is not at its minimum for  $k + 1$ . With the new cluster center, we know that there is some point where its  $\mathcal{J}$  term is now reduced to 0. This means  $\mathcal{J}$  has decreased when we added the  $k + 1$ th cluster. Hence, by induction, QED.

2. **(8 Points)** Consider the K-means (Lloyd's) clustering algorithm we studied in class. We terminate the algorithm when there are no changes to the objective. Show that the algorithm terminates in a finite number of steps.

Since we know that K-means is non-increasing, that means its steps cannot revisit the same cluster state twice unless the state is a converged state. Since the number of states is finite, this implies that K-means must terminate in a finite number steps.

## 1.2 Experiment (20 Points)

In this question, we will evaluate K-means clustering and GMM on a simple 2 dimensional problem. First, create a two-dimensional synthetic dataset of 300 points by sampling 100 points each from the three Gaussian distributions shown below:

$$P_a = \mathcal{N}\left(\begin{bmatrix} -1 \\ -1 \end{bmatrix}, \sigma \begin{bmatrix} 2 & 0.5 \\ 0.5 & 1 \end{bmatrix}\right), \quad P_b = \mathcal{N}\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix}, \sigma \begin{bmatrix} 1 & -0.5 \\ -0.5 & 2 \end{bmatrix}\right), \quad P_c = \mathcal{N}\left(\begin{bmatrix} 0 \\ 1 \end{bmatrix}, \sigma \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}\right)$$

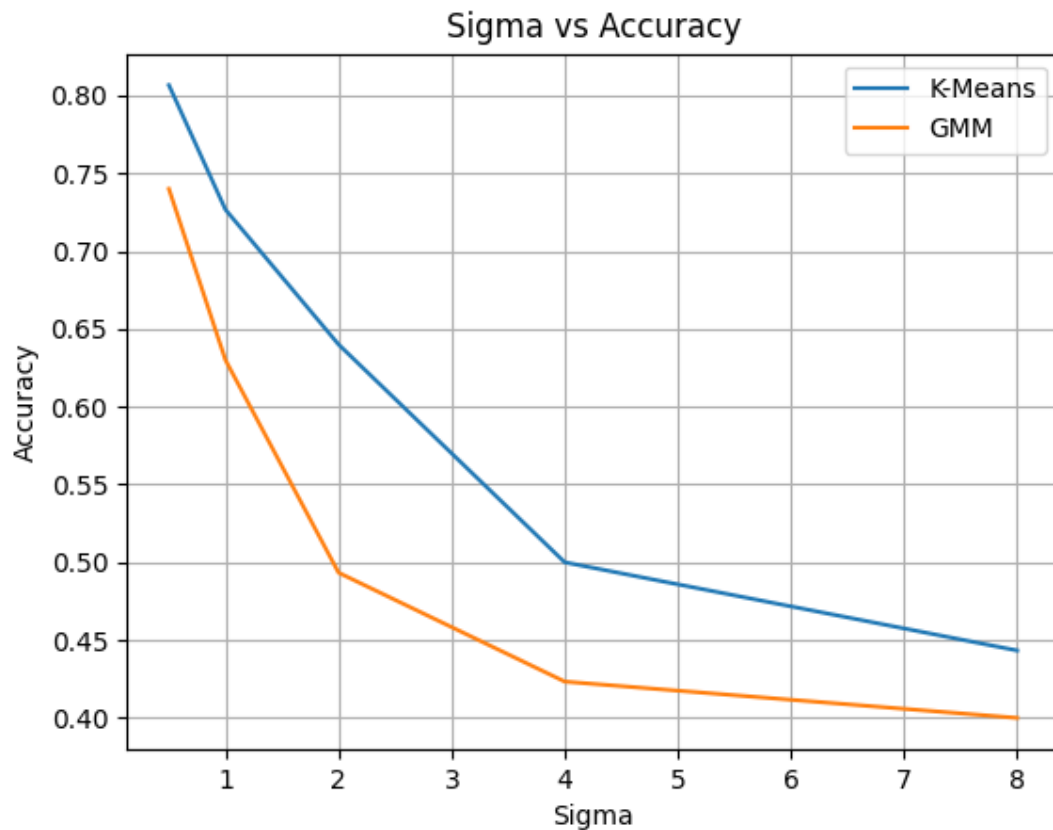
Here,  $\sigma$  is a parameter we will change to produce different datasets.

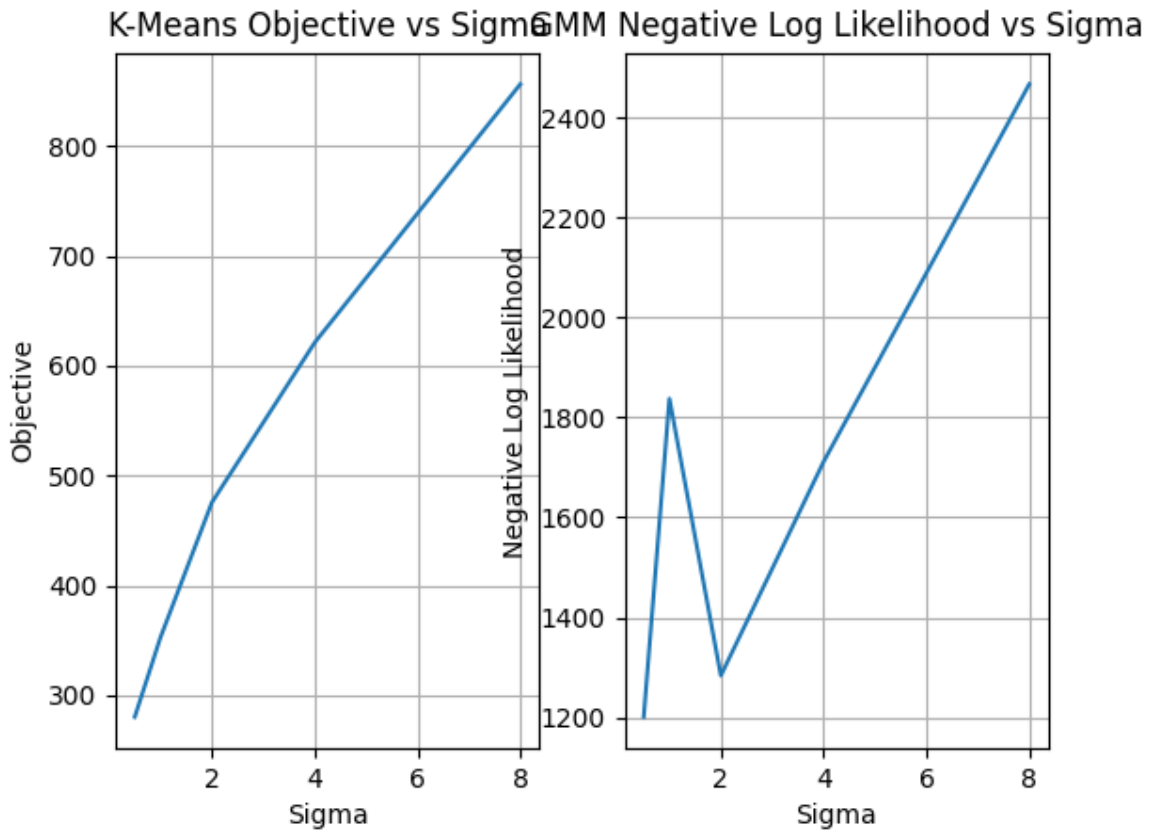
First implement K-means clustering and the expectation maximization algorithm for GMMs. Execute both methods on five synthetic datasets, generated as shown above with  $\sigma \in \{0.5, 1, 2, 4, 8\}$ . Finally, evaluate both methods on (i) the clustering objective (1) and (ii) the clustering accuracy. For each of the two criteria, plot the value achieved by each method against  $\sigma$ .

Guidelines:

- Both algorithms are only guaranteed to find only a local optimum so we recommend trying multiple restarts and picking the one with the lowest objective value (This is (1) for K-means and the negative log likelihood for GMMs). You may also experiment with a smart initialization strategy (such as kmeans++).
- To plot the clustering accuracy, you may treat the ‘label’ of points generated from distribution  $P_u$  as  $u$ , where  $u \in \{a, b, c\}$ . Assume that the cluster id  $i$  returned by a method is  $i \in \{1, 2, 3\}$ . Since clustering is an unsupervised learning problem, you should obtain the best possible mapping from  $\{1, 2, 3\}$  to  $\{a, b, c\}$  to compute the clustering objective. One way to do this is to compare the clustering centers returned by the method (centroids for K-means, means for GMMs) and map them to the distribution with the closest mean.

Points break down: 7 points each for implementation of each method, 6 points for reporting of evaluation metrics.





## 2 Linear Dimensionality Reduction

### 2.1 Principal Components Analysis (10 points)

Principal Components Analysis (PCA) is a popular method for linear dimensionality reduction. PCA attempts to find a lower dimensional subspace such that when you project the data onto the subspace as much of the information is preserved. Say we have data  $X = [x_1^\top; \dots; x_n^\top] \in \mathbb{R}^{n \times D}$  where  $x_i \in \mathbb{R}^D$ . We wish to find a  $d$  ( $< D$ ) dimensional subspace  $A = [a_1, \dots, a_d] \in \mathbb{R}^{D \times d}$ , such that  $a_i \in \mathbb{R}^D$  and  $A^\top A = I_d$ , so as to maximize  $\frac{1}{n} \sum_{i=1}^n \|A^\top x_i\|^2$ .

1. **(4 Points)** Suppose we wish to find the first direction  $a_1$  (such that  $a_1^\top a_1 = 1$ ) to maximize  $\frac{1}{n} \sum_i (a_1^\top x_i)^2$ . Show that  $a_1$  is the first right singular vector of  $X$ .

We wish to maximize the following:

$$\begin{aligned}
 \frac{1}{n} \sum_{i=1}^n (a_1^\top x_i)^2 &= \frac{1}{n} \sum_{i=1}^n a_1^\top x_i x_i^\top a_1 \\
 &= a_1^\top \left( \frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right) a_1 \\
 &= a_1^\top \left( \frac{1}{n} X^\top X \right) a_1
 \end{aligned}$$

Using SVD,  $X = U \Sigma V^\top$ , then the eigen-decomposition of the covariance matrix  $S = \frac{1}{n} X^\top X$  is:

$$S = \frac{1}{n} V \Sigma^\top U^\top U \Sigma V^\top$$

Since  $U$  is orthogonal, and  $U^\top U = I$

$$S = \frac{1}{n} V \Sigma^\top \Sigma V^\top$$

$$= V \left( \frac{1}{n} \Sigma^2 \right) V^\top$$

Therefore the eigenvectors of  $S$  are the right singular vectors of  $X$ . Then,

$$\max_{a_1} \frac{1}{n} \sum_{i=1}^n (a_1^\top x_i)^2 = \max_{a_1} a_1^\top \left( V \left( \frac{1}{n} \Sigma^2 \right) V^\top \right) a_1$$

subject to the constraint  $\|a_1\| = 1$ , which is the first eigenvector of  $V \left( \frac{1}{n} \Sigma^2 \right) V^\top$ , which constraints  $a_1$  to be the first right singular vector of  $X$ .

2. **(6 Points)** Given  $a_1, \dots, a_k$ , let  $A_k = [a_1, \dots, a_k]$  and  $\tilde{x}_i = x_i - A_k A_k^\top x_i$ . We wish to find  $a_{k+1}$ , to maximize  $\frac{1}{n} \sum_i (a_{k+1}^\top \tilde{x}_i)^2$ . Show that  $a_{k+1}$  is the  $(k+1)^{th}$  right singular vector of  $X$ .

Let  $\tilde{X} = X - X A A^\top = X(I - A A^\top)$ . Find the orthonormal  $a_{k+1}$  to maximize  $J(a) = a^\top (I - A A^\top)^\top X^\top X (I - A A^\top) a$ . For all given  $a$ ,  $J(a) = 0$ . Therefore to maximize  $a_{k+1}^\top \Sigma a_{k+1}$  is subject to the constraints  $\|a_{k+1}\| = 1$  and  $a_{k+1} \neq a_i$  for  $1 \leq i \leq k$ . Then this is the  $(k+1)^{st}$  eigenvector of  $V \left( \frac{1}{n} \Sigma^2 \right) V^\top$ .

## 2.2 Dimensionality reduction via optimization (22 points)

We will now motivate the dimensionality reduction problem from a slightly different perspective. The resulting algorithm has many similarities to PCA. We will refer to method as DRO.

As before, you are given data  $\{x_i\}_{i=1}^n$ , where  $x_i \in \mathbb{R}^D$ . Let  $X = [x_1^\top; \dots; x_n^\top] \in \mathbb{R}^{n \times D}$ . We suspect that the data actually lies approximately in a  $d$  dimensional affine subspace. Here  $d < D$  and  $d < n$ . Our goal, as in PCA, is to use this dataset to find a  $d$  dimensional representation  $z$  for each  $x \in \mathbb{R}^D$ . (We will assume that the span of the data has dimension larger than  $d$ , but our method should work whether  $n > D$  or  $n < D$ .)

Let  $z_i \in \mathbb{R}^d$  be the lower dimensional representation for  $x_i$  and let  $Z = [z_1^\top; \dots; z_n^\top] \in \mathbb{R}^{n \times d}$ . We wish to find parameters  $A \in \mathbb{R}^{D \times d}$ ,  $b \in \mathbb{R}^D$  and the lower dimensional representation  $Z \in \mathbb{R}^{n \times d}$  so as to minimize

$$J(A, b, Z) = \frac{1}{n} \sum_{i=1}^n \|x_i - A z_i - b\|^2 = \|X - Z A^\top - \mathbf{1} b^\top\|_F^2. \quad (2)$$

Here,  $\|A\|_F^2 = \sum_{i,j} A_{ij}^2$  is the Frobenius norm of a matrix.

1. **(3 Points)** Let  $M \in \mathbb{R}^{d \times d}$  be an arbitrary invertible matrix and  $p \in \mathbb{R}^d$  be an arbitrary vector. Denote,  $A_2 = A_1 M^{-1}$ ,  $b_2 = b_1 - A_1 M^{-1} p$  and  $Z_2 = Z_1 M^\top + \mathbf{1} p^\top$ . Show that both  $(A_1, b_1, Z_1)$  and  $(A_2, b_2, Z_2)$  achieve the same objective value  $J(2)$ .

Since both achieve the same objective value  $J$ , we can form the equation  $Z_1 A_1^\top + \mathbf{1} b_1^\top = Z_2 A_2^\top + \mathbf{1} b_2^\top$ . Then,

$$\begin{aligned} Z_2 A_2^\top + \mathbf{1} b_2^\top &= (Z_1 M^\top + \mathbf{1} p^\top) M^{-\top} A_1^\top + \mathbf{1} (b_1 - A_1 M^{-1} p)^\top \\ &= Z_1 M^\top M^{-\top} A_1^\top + \mathbf{1} p^\top M^{-\top} A_1^\top + \mathbf{1} b_1^\top - \mathbf{1} p^\top M^{-\top} A_1^\top \\ &= Z_1 A_1^\top + \mathbf{1} b_1^\top \end{aligned}$$

Then,

$$J(A_2, b_2, Z_2) = \|X - Z_2 A_2^\top - \mathbf{1} b_2^\top\|_F^2 = \|X - Z_1 A_1^\top - \mathbf{1} b_1^\top\|_F^2 = J(A_1, b_1, Z_1)$$

Therefore, in order to make the problem determined, we need to impose some constraint on  $Z$ . We will assume that the  $z_i$ 's have zero mean and identity covariance. That is,

$$\bar{Z} = \frac{1}{n} \sum_{i=1}^n z_i = \frac{1}{n} Z^\top \mathbf{1}_n = 0, \quad S = \frac{1}{n} \sum_{i=1}^n z_i z_i^\top = \frac{1}{n} Z^\top Z = I_d$$

Here,  $\mathbf{1}_d = [1, 1, \dots, 1]^\top \in \mathbb{R}^d$  and  $I_d$  is the  $d \times d$  identity matrix.

2. **(16 Points)** Outline a procedure to solve the above problem. Specify how you would obtain  $A, Z, b$  which minimize the objective and satisfy the constraints.

**Hint:** The rank  $k$  approximation of a matrix in Frobenius norm is obtained by taking its SVD and then zeroing out all but the first  $k$  singular values.

First we find the value for  $b$ , we take the derivative w.r.t  $b$  and set to zero:

$$\begin{aligned}\frac{\partial}{\partial b} J &= 2(X - ZA^\top - \mathbf{1}b^\top)^\top \mathbf{1} = 0 \\ X^\top \mathbf{1}^\top - AZ^\top \mathbf{1} &= b^\top \mathbf{1}^\top \mathbf{1} \\ X^\top \mathbf{1}^\top - AZ^\top \mathbf{1} &= bn\end{aligned}$$

And assuming that we can find some  $Z$  that satisfies the constraint  $Z^\top \mathbf{1} = 0$ ,

$$\begin{aligned}bn &= X^\top \mathbf{1}^\top \\ b &= \frac{1}{n} X^\top \mathbf{1} \\ b &= \bar{X}\end{aligned}$$

where  $\bar{X}$  is the mean of  $X$ .

To minimize our function  $\|X - ZA^\top - \mathbf{1}b^\top\|_F^2$  w.r.t  $A, Z$ , we want  $ZA^\top$  to be the best rank  $k$  approximation of  $X - \mathbf{1}b^\top$ .

We can then take the SVD of  $\tilde{X} = X - \mathbf{1}b^\top$  and then zeroing out all but the first  $k$  singular values. Deconstructing using SVD gives us  $\tilde{X} = U\Sigma V^\top$  where  $U \in \mathbb{R}^{n \times n}$ ,  $\Sigma$  is diagonal,  $\Sigma \in \mathbb{R}^{n \times D}$  and  $V \in \mathbb{R}^{D \times D}$ . Let the first  $d$  columns of  $U$  be denoted by  $U_d$ , the top  $d \times d$  block of  $\Sigma$  be denoted by  $\Sigma_d$  and the first  $d$  columns of  $V$  be denoted by  $V_d$ , which gives us the rank  $d$  approximation of  $\tilde{X}$  to be  $U_d \Sigma_d V_d^\top$ .

Now we have to find  $Z$  that satisfies the constraints. Let  $Z = \sqrt{n}U_d I_d^{1/2}$  and  $A^\top = \frac{1}{\sqrt{n}} I_d^{1/2} \Sigma_d V_d^\top$ , then we can satisfy the identity covariance since:

$$\frac{1}{n} Z^\top Z = \frac{1}{n} n I_d^{1/2} U_d^\top U_d I_d^{1/2} = I_d$$

We can also show that  $Z$  satisfies  $Z^\top \mathbf{1} = \mathbf{0}$ . Since  $\tilde{X}^\top \mathbf{1} = V\Sigma U^\top \mathbf{1} = \mathbf{0}$  (from SVD). Since  $V\Sigma$  is full rank, it implies that the following must be true:  $U^\top \mathbf{1} = \mathbf{1}$ . Hence  $U_d^\top \mathbf{1} = \mathbf{0}$ , which also means  $Z^\top \mathbf{1} = \mathbf{0}$ .

3. **(3 Points)** You are given a point  $x_*$  in the original  $D$  dimensional space. State the rule to obtain the  $d$  dimensional representation  $z_*$  for this new point. (If  $x_*$  is some original point  $x_i$  from the  $D$ -dimensional space, it should be the  $d$ -dimensional representation  $z_i$ .)

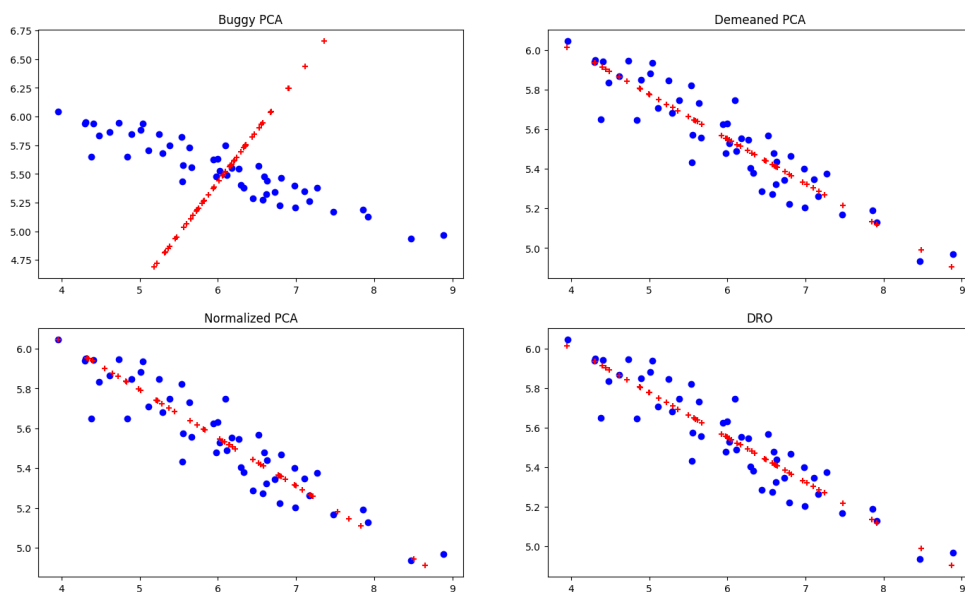
$$z_* = (A^\top A)^{-1} A^\top (x_* - b)$$

## 2.3 Experiment (34 points)

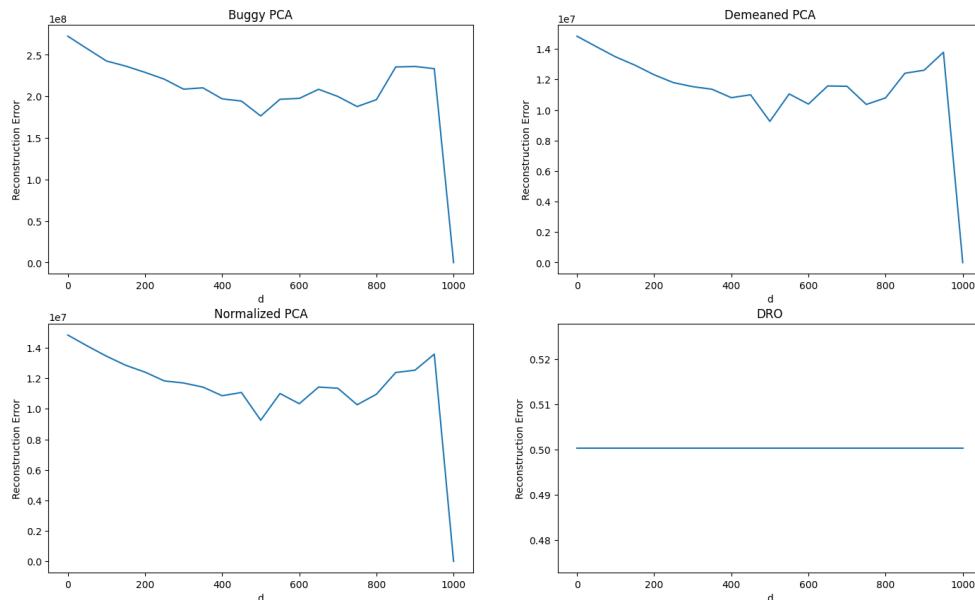
Here we will compare the above three methods on two data sets.

- We will implement three variants of PCA:
  1. "buggy PCA": PCA applied directly on the matrix  $X$ .
  2. "demeaned PCA": We subtract the mean along each dimension before applying PCA.
  3. "normalized PCA": Before applying PCA, we subtract the mean and scale each dimension so that the sample mean and standard deviation along each dimension is 0 and 1 respectively.
- One way to study how well the low dimensional representation  $Z$  captures the linear structure in our data is to project  $Z$  back to  $D$  dimensions and look at the reconstruction error. For PCA, if we mapped it to  $d$  dimensions via  $z = Vx$  then the reconstruction is  $V^\top z$ . For the preprocessed versions, we first do this and then reverse the preprocessing steps as well. For DRO we just compute  $Az + b$ . We will compare all methods by the reconstruction error on the datasets.
- Please implement code for the methods: Buggy PCA (just take the SVD of  $X$ ), Demeaned PCA, Normalized PCA, DRO. In all cases your function should take in an  $n \times d$  data matrix and  $d$  as an argument. It should return the  $d$  dimensional representations, the estimated parameters, and the reconstructions of these representations in  $D$  dimensions.
- You are given two datasets: A two Dimensional dataset with 50 points `data2D.csv` and a thousand dimensional dataset with 500 points `data1000D.csv`.
- For the 2D dataset use  $d = 1$ . For the 1000D dataset, you need to choose  $d$ . For this, observe the singular values in DRO and see if there is a clear "knee point" in the spectrum. Attach any figures/ Statistics you computed to justify your choice.
- For the 2D dataset you need to attach the a plot comparing the original points with the reconstructed points for all 4 methods. For both datasets you should also report the reconstruction errors, that is the squared sum of differences  $\sum_{i=1}^n \|x_i - r(z_i)\|^2$ , where  $x_i$ 's are the original points and  $r(z_i)$  are the  $D$  dimensional points reconstructed from the  $d$  dimensional representation  $z_i$ .

2D Results: (Demeaned PCA and DRO performed best)



Algorithm	Reconstruction Error
Buggy PCA	44.3451541867397
Demeaned PCA	0.500304281425646
Normalized PCA	2.4736041727385345
DRO	0.500304281425646



For buggy PCA, demeaned PCA and normalized PCA, there is a clear "knee point" in the spectrum at  $d=500$  as there is a big dip for those 3 algorithms. However, there isn't a clear "knee point" for DRO as the reconstruction error seems to be not changing for DRO.

- **Questions:** After you have completed the experiments, please answer the following questions.
  1. Look at the results for Buggy PCA. The reconstruction error is bad and the reconstructed points don't seem to well represent the original points. Why is this?  
**Hint:** Which subspace is Buggy PCA trying to project the points onto?  
**Buggy PCA is projecting the points onto a subspace that passes through the origin. The other PCA algorithms allow for variance maximization through normalization allowing the points to be projected onto directions that maximizes the variance.**
  2. The error criterion we are using is the average squared error between the original points and the reconstructed points. In both examples DRO and demeaned PCA achieves the lowest error among all methods. Is this surprising? Why?  
**Not surprising as DRO and demeaned PCA are equivalent as derived earlier with zero mean and identity covariance.**
- **Point allocation:**
  - Implementation of the three PCA methods: **(6 Points)**
  - Implementation of DRO: **(6 points)**
  - Plots showing original points and reconstructed points for 2D dataset for each one of the 4 methods: **(10 points)**
  - Implementing reconstructions and reporting results for each one of the 4 methods for the 2 datasets: **(5 points)**
  - Choice of  $d$  for 1000D dataset and appropriate justification: **(3 Points)**
  - Questions **(4 Points)**