

# STAT 154 Spring 2019: Sample Mid-term Exam

Instructor: Prof. Bin Yu  
Mar 21, 8:10 AM–9:30 AM

Your Name: \_\_\_\_\_ Student ID: \_\_\_\_\_

**Maximum Points: 50, Time: 80 minutes**

## Instructions

- Do not turn the page until you are told do so.
- **WRITE your student ID CLEARLY** on each page in the space provided at the top.
- This exam has **5 questions** with **24 parts** in total.
- Your answer **will be graded only** if it is written in the space provided after the question. You may use the blank pages at the end for rough work, your rough work will **NOT** be graded.
- Try not to get stuck in a particular problem.
- *When time is up*, please take your exam in your **left hand** and raise it. When asked pass it on to your left (facing the board). Please maintain the sequence of your seating. Instructor/GSI will collect the copies from the leftmost student in the row.
- As promised, there is a help-sheet (last page) for your comfort.

## Pre-Exam Questions

1. What is (are) your favorite movie(s)?
2. What do you like the most about this class?

## 1 True or False (with justification, 10 pts)

Please justify your reasoning in one line.

- (a) For a model to make meaningful predictions on the future data, we need some similarity between the representative data that was used to build the model and the future data.

True (Hw1) - The new data has to be similar  
else  $\rightarrow$  wrong conclusion

- (b) If the rank of a  $d \times d$  projection matrix is  $k$ , then its trace is equal to  $d - k$ .

False. ( $\text{tr} = k$ )

- (c) K-means is guaranteed to return the global minimizer of the following objective:

$$\min_{C_1, \dots, C_K} \min_{\mu_1, \dots, \mu_K} \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - \mu_k\|_2^2.$$

True  
False (local minima)  $\rightarrow$  Non-convex problem

- (d) For any  $d \times d$  matrix, we have  $\nabla_\theta(\theta^\top \mathbf{A} \theta) = 2\mathbf{A}\theta$ .

False

$(A + A^\top)\mathbf{e}$ .

$$\nabla_\theta \theta^\top \mathbf{A} \theta = (\mathbf{A} + \mathbf{A}^\top)$$

- (e) If some features ( $\{x_1, \dots, x_n\}$ ) are linearly dependent, at least one singular value of the feature matrix  $\mathbf{X}$  is zero.

*True ( $\mathbf{X}$  is not full rank of cols are linearly dependent).*

- (f) Lasso always produces features that are subset of original features.

*True Some coeff were  $\rightarrow$  zero.*

- (g) Number of features in a model or the number of iterations when using a method like gradient descent can act as regularization for training the model.

*True Complexity of model vs. appeared.*

- (h) Gradient descent on the ordinary least squares objective function is guaranteed to converge to the ridge estimate as long as the step size is positive.

*False Converge  $\rightarrow$  OLS estimate  
Stepsize can be too large*

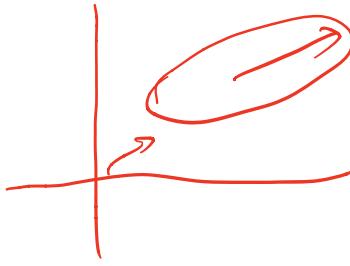
- (i) The output of PCA is a new representation of the data that is always in terms of the principal components that are best at predicting the output variable.

*False, PCA is unsupervised*

*$\hookrightarrow$  Not for prediction.  
No explicit output  $y$ .*

- (j) When the features in the data matrix have same units and are comparable to each other, simply centering them and not scaling them before doing PCA is a good idea.

*~~False~~. True - Only scale when features have different units  $\rightarrow$  so that data is comparable.*



## 2 EM with Poisson mixture (10 pts)

We work with the following simple two mixture model:

$$\begin{aligned} Z &\sim \text{Bernoulli}(1 - w) + 1 \\ X|Z = 1 &\sim \text{Poisson}(\lambda) \quad \text{and} \\ X|Z = 2 &\sim \text{Poisson}(\mu) \end{aligned} \tag{1}$$

where  $Z$  denotes the label of the Poisson distribution from which  $X$  is drawn. Given  $n$  observations  $\{x_1, \dots, x_n\}$  only for  $X$  (i.e., the labels are unobserved), our goal is to infer the maximum-likelihood parameters for  $\lambda_1, \lambda_2$ , and  $w$  using EM.

Note that for a random variable  $Y \sim \text{Poisson}(\lambda)$ , we have

$$\mathbb{P}(Y = k) = e^{-\lambda} \frac{\lambda^k}{k!}, \quad \text{for } k = 0, 1, \dots$$

- (a) (2 pts) Derive the complete log-likelihood for the observed data , i.e., derive the expression for  $\sum_{i=1}^n \log \mathbb{P}(X = x_i; \lambda, \mu, w)$  in terms of the parameters and the data  $\{x_1, \dots, x_n\}$ .

- (b) (2 pts) Write down the expression for the lower bound  $\mathcal{F}(\lambda, \mu, w, q)$  where  $q = (q_1, \dots, q_n)$  where  $q = (q_1, \dots, q_n)$  and  $q_i = p(Z_i = 1)$  denotes the parameter for an appropriate distribution on the hidden label  $Z_i$ .

- (c) (2 pts) Given the iterates  $w_t, \lambda_t$  and  $\mu_t$ , compute the E-step updates, i.e., **write the expressions for**

$$q^{t+1} = \arg \max_q \mathcal{F}(\lambda_t, \mu_t, w_t, q)$$

**No need to show any derivation. Simply writing the expressions suffices.**

*Hint:  $q_i^{t+1}$  is related to a posterior distribution.*

- (d) (4 pts) **Derive the M-step updates, i.e., compute the expressions**

$$\lambda_{t+1}, \mu_{t+1}, w_{t+1} = \arg \max_{\lambda, \mu, w} \mathcal{F}(\lambda, \mu, w, q^{t+1}),$$

*Hint: Some gradients have to be set to zero.*

### 3 Weighted Least Squares (8 pts)

The weighted least squares (WLS) is a slight modification on least squares:

$$\min_{\theta} \sum_{i=1}^n w_i (x_i^\top \theta - y_i)^2 \quad (2)$$

i.e., the  $i$ -th term in the least squares objective is weighted by a term  $w_i$ .

Define the matrices:

$$\mathbf{W} = \begin{bmatrix} w_1 & & \\ & w_2 & \\ & & \ddots & \\ & & & w_n \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix} \in \mathbb{R}^{n \times d} \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

- (a) (2 pts) **Write the objective in equation (2) using the matrix-vector notation.**  
 Your answer should be in terms of  $\mathbf{W}$ ,  $\mathbf{X}$ ,  $\mathbf{y}$  and  $\theta$ .

- (b) (3 pts) Let the matrix  $\sqrt{\mathbf{W}}\mathbf{X}$  be full column rank. **Derive the closed form solution for the WLS estimate obtained by minimizing the objective in equation (2).**  
*Hint: You may use the OLS formula for this part.*

- (c) (3 pts) Consider the following linear model:

$$y_i = x_i^\top \theta + \varepsilon_i, \quad (3)$$

where  $\varepsilon_i \sim \mathcal{N}(0, \sigma_i^2)$ . Assume that the noise variables  $\varepsilon_i$  are independent (not identically distributed) of each other and independent of  $x_i$ 's. Also assume that the variances  $\sigma_i^2$  are known. Given  $n$  independent samples  $\{(x_i, y_i), i = 1, \dots, n\}$  from the model (3), **show that computing the maximum likelihood estimate of  $\theta$  for this data is equivalent to solving a weighted least squares problem. What are the weights  $w_i$ ?**

## 4 Ridge Regression vs PCA (12 pts)

Assume we are given  $n$  training data points  $(\mathbf{x}_i, y_i)$ . We collect the responses into  $\mathbf{y} \in \mathbb{R}^n$ , and the features into the matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  where the rows are the  $d$ -dimensional feature vectors  $\mathbf{x}_i^\top$  corresponding to each training point. Furthermore, assume that  $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \mathbf{0}$ ,  $n > d$  and  $\mathbf{X}$  has rank  $d$ .

In this problem we want to compare two procedures: The first is ridge regression with hyper-parameter  $\lambda$ , while the second is applying ordinary least squares after using PCA to reduce the feature dimension from  $d$  to  $k$  (we give this latter approach the short-hand name  $k$ -PCA-OLS where  $k$  is the hyper-parameter).

Notation: The singular value decomposition of  $\mathbf{X}$  reads  $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$  where  $\mathbf{U} \in \mathbb{R}^{n \times n}$ ,  $\Sigma \in \mathbb{R}^{n \times d}$  and  $\mathbf{V} \in \mathbb{R}^{d \times d}$ . We denote by  $\mathbf{u}_j$  the  $n$ -dimensional column vectors of  $\mathbf{U}$  and by  $\mathbf{v}_j$  the  $d$ -dimensional column vectors of  $\mathbf{V}$ . Furthermore the diagonal entries  $\sigma_j = \Sigma_{j,j}$  of  $\Sigma$  satisfy  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d > 0$ . For notational convenience, assume that  $\sigma_j = 0$  for  $j > d$ .

- (a) (4 pts) It turns out that the ridge regression optimizer (with  $\lambda > 0$ ) in the  $\mathbf{V}$ -transformed coordinates

$$\hat{\theta}_{\text{ridge}} = \arg \min_{\theta} \|\mathbf{X}\mathbf{V}\theta - \mathbf{y}\|_2^2 + \lambda \|\theta\|_2^2$$

has the following expression:

$$\hat{\theta}_{\text{ridge}} = \text{diag}\left(\frac{\sigma_j}{\lambda + \sigma_j^2}\right) \mathbf{U}^\top \mathbf{y}. \quad (4)$$

Use  $\hat{y}_{test} = \mathbf{x}_{test}^\top \mathbf{V} \hat{\theta}_{\text{ridge}}$  to denote the resulting prediction for a hypothetical  $\mathbf{x}_{test}$ . Using (4) and the appropriate scalar  $\{\alpha_i\}$ , this can be written as:

$$\hat{y}_{test} = \mathbf{x}_{test}^\top \sum_{j=1}^d \mathbf{v}_j \alpha_j \mathbf{u}_j^\top \mathbf{y}. \quad (5)$$

**What are the  $\alpha_j \in \mathbb{R}$  for this to correspond to (4) from ridge regression? Show your work.**

- (b) (5 pts) Suppose that we do k-PCA-OLS — i.e. ordinary least squares on the reduced  $k$ -dimensional feature space obtained by projecting the raw feature vectors onto the  $k < d$  principal components of the covariance matrix  $\mathbf{X}^\top \mathbf{X}$ . Use  $\hat{y}_{test}$  to denote the resulting prediction for a hypothetical  $\mathbf{x}_{test}$ ,

It turns out that the learned k-PCA-OLS predictor can be written as:

$$\hat{y}_{test} = \mathbf{x}_{test}^\top \sum_{j=1}^d \mathbf{v}_j \alpha_j \mathbf{u}_j^\top \mathbf{y}. \quad (6)$$

**Give the  $\alpha_j \in \mathbb{R}$  coefficients for k-PCA-OLS. Show work. These  $\alpha_j$ 's would be different compared to the previous part.**

*Hint 1: some of these  $\alpha_j$  will be zero. Also, if you want to use the compact form of the SVD, feel free to do so if that speeds up your derivation.*

*Hint 2: some inspiration may be possible by looking at the next part for an implicit clue as to what the answer might be.*

- (c) (3 pts) For the following part,  $d = 5$ . The following  $\alpha := (\alpha_1, \dots, \alpha_5)$  (written out to two significant figures) are the results of OLS (i.e. what we would get from ridge regression in the limit  $\lambda \rightarrow 0$ ),  $\lambda$ -ridge-regression, and  $k$ -PCA-OLS for some  $\mathbf{X}, \mathbf{y}$  (identical for each method) and  $\lambda = 1, k = 3$ . **Write down which procedure was used for each of the three sub-parts below.**

We hope this helps you intuitively see the connection between these three methods.

*Hint: It is not necessary to find the singular values of  $\mathbf{X}$  explicitly, or to do any numerical computations at all.*

- (i)  $\alpha = (0.01, 0.1, 0.5, 0.1, 0.01)$
- (ii)  $\alpha = (0.01, 0.1, 1, 0, 0)$
- (iii)  $\alpha = (0.01, 0.1, 1, 10, 100)$

ii) PCA - Exacer ress

i) Ridge - Should have smaller coefficient.

iii) OLS -

## 5 EDA with Ames Data (10 pts)

In the two figures below, you are given two scatter-plots with the same data namely, two features “SalePrice” and “Gr.Liv.Area” of the Ames dataset.

- (a) (2 pts) Figure 1 was generated using ggplot2 library with the command

```
ggplot(x) + geom_point(aes(x=Ames[, ``Gr.Liv.Area``], y=Ames[, ``SalePrice``]))
```

but without changing any parameters. However to produce Figure 2, certain aesthetic changes were done (no data processing was done). Can you guess what kind of changes in the aesthetics were made? Which figure is more informative? And why?

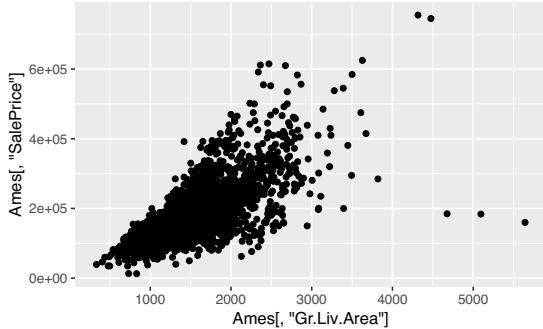


Fig 1.

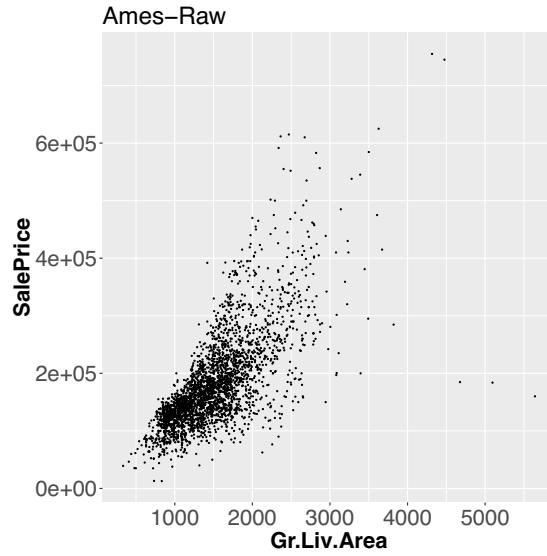


Fig. 2

- Point Size
- X Y- Axis Ratio
- Title
- Font Size.

- (b) (2 pts) The data from the previous part was *transformed* in two different ways to obtain the following two figures. **Can you guess what kind of transformations were made to obtain these plots?** To answer this question, you should compare these figures with Figure 2 in the previous part.

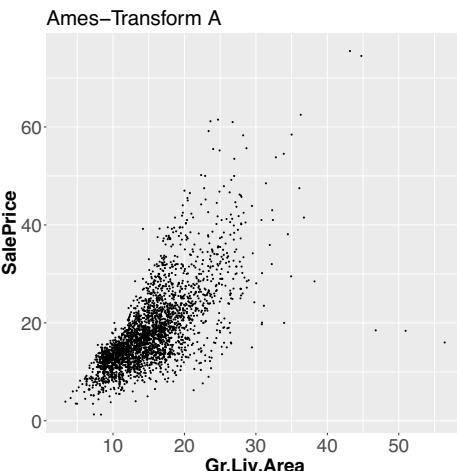


Fig. 3.

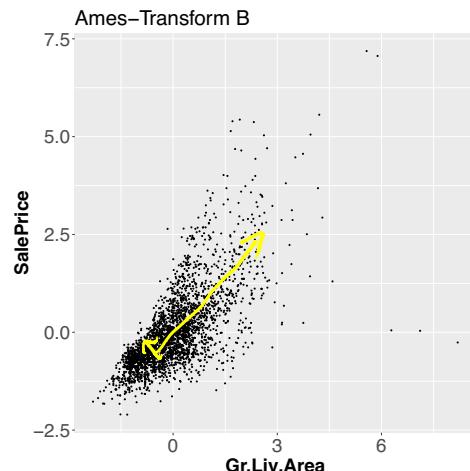


Fig. 4

- Scaling

- Cencoring

- Scaling

} Standardizing.

- (c) (3 pts) Between the datasets Ames-Raw (Figure 2), Ames-Transform A (Figure 3) and Ames-Transform B (Figure 3), which version of the dataset is most suitable for doing a PCA analysis? Why? Also plot the principal vectors on the figure that you think is most appropriate.

Figure 4

- XY Axes are of same size
- Cenred & Scaling

- (d) (3 pts) Do you think using K-means or EM would provide some insight with any of the datasets (Figure 2/3/4)? Justify.

No. No. visible clusters

## Help-sheet

- Given features  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and responses  $\mathbf{y} \in \mathbb{R}^n$ , if the feature matrix  $\mathbf{X}$  is full column rank, the OLS solution is given by

$$\hat{\beta}^{\text{OLS}} = \arg \min_{\beta} \|\mathbf{X}\beta - \mathbf{y}\|_2^2 = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

- The pdf/PMF of a few distributions is given below:

Distribution	Notation	pdf ( $p$ ) / PMF ( $\mathbb{P}$ )
Multi-variate Gaussian	$Z \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$	$p(Z = \mathbf{z}) = \frac{1}{\sqrt{\det(2\pi\boldsymbol{\Sigma})}} \exp\left(-\frac{(\mathbf{z} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu})}{2}\right)$
Exponential	$Z \sim \text{Exponential}(\lambda)$	$p(Z = z) = \lambda e^{-\lambda z}, \quad z \geq 0.$
Bernoulli	$Z \sim \text{Bernoulli}(\alpha)$	$\mathbb{P}(Z = z) = \alpha^z (1 - \alpha)^{(1-z)}, \quad z \in \{0, 1\}$
Poisson	$Z \sim \text{Poisson}(\lambda)$	$\mathbb{P}(Z = z) = e^{-\lambda} \frac{\lambda^z}{z!}, \quad z \in \{0, 1, 2, \dots\}.$

## Fun Space

Feel free to draw/write something if you want to or give us suggestions or complaints. You can also use this space to report anything suspicious that you might have noticed.

SID:

Rough space

---

Space for rough work (will not be graded)

SID:

Rough space

---

Space for rough work (will not be graded)

## 2 EM with Poisson mixture (10 pts)

We work with the following simple two mixture model:

$$\begin{aligned} Z &\sim \text{Bernoulli}(1-w) + 1 \\ X|Z=1 &\sim \text{Poisson}(\lambda) \quad \text{and} \\ X|Z=2 &\sim \text{Poisson}(\mu) \end{aligned}$$

$$\begin{aligned} X^\top A X > 0 & \quad \cancel{\underline{\underline{X^\top A X > 0}}} \\ \cancel{\underline{\underline{Z^\top A^\top B A Z > 0}}} & \downarrow \\ \cancel{\underline{\underline{(A^\top B A Z)^2 > 0}}} & \quad (1) \\ \cancel{\underline{\underline{Z = \tau}}} & \quad \cancel{\underline{\underline{X = \tau}}} \end{aligned}$$

where  $Z$  denotes the label of the Poisson distribution from which  $X$  is drawn. Given  $n$  observations  $\{x_1, \dots, x_n\}$  only for  $X$  (i.e., the labels are unobserved), our goal is to infer the maximum-likelihood parameters for  $\lambda_1, \lambda_2$ , and  $w$  using EM.

Note that for a random variable  $Y \sim \text{Poisson}(\lambda)$ , we have

$$\mathbb{P}(Y=k) = e^{-\lambda} \frac{\lambda^k}{k!}, \quad \text{for } k=0,1,\dots$$

$$\cancel{\underline{\underline{X^\top B X > 0}}}$$

- (a) (2 pts) Derive the complete log-likelihood for the observed data , i.e., derive the expression for  $\sum_{i=1}^n \log \mathbb{P}(X=x_i; \lambda, \mu, w)$  in terms of the parameters and the data  $\{x_1, \dots, x_n\}$ .

$$\theta = (\lambda, \mu, w)$$

$$P(X=x; \theta)$$

$$= \sum_{k=1}^2 P(X=x_i, Z=k; \theta)$$

$$= \sum_{k=1}^2 P(Z=k; \theta) P(X=x_i | Z=k; \theta)$$

$$= w \cdot e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} + (1-w) e^{-\mu} \frac{\mu^{x_i}}{x_i!}$$

$$\log P(x_1=x_1, \dots, x_n=x_n; \theta)$$

$$= \sum_{i=1}^n \log P(X=x_i; \theta)$$

$$= \sum_{i=1}^n \log \left( \sum_{k=1}^2 P(X=x_i, Z=k; \theta) \right)$$

$$A^{-1} = A$$

$$A$$

$$A^{-1} \rightarrow A^\top A$$

$$X^\top (A^\top A) X \quad A^{-1} \neq 0$$

$$A^\top = B^\top$$

$$\begin{aligned} & \arg \max \mathcal{L}(\theta) \\ \mathcal{L}(\theta) &= \sum_{i=1}^n \log \left( \sum_{k=1}^2 P(X=x_i, Z=k; \theta) \right) \\ &= \sum_{i=1}^n \log \left( \sum_{k=1}^2 g_i(k) \underbrace{\frac{P(X=x_i, Z=k; \theta)}{g_i(k)}}_{\sim \sim} \right) \\ & \left( \begin{array}{l} \log X \leq \log E(X) \\ \sim \sim \end{array} \right) \\ & \geq \sum_{i=1}^n \sum_{k=1}^2 g_i(k) \log \underbrace{\frac{P(X=x_i, Z=k; \theta)}{g_i(k)}}_{\mathcal{F}(\theta, g_i)} \\ & \left( \begin{array}{l} \mathcal{F}(\theta, g_i) \\ \sim \sim \end{array} \right) \\ &= \sum_{i=1}^n \sum_{k=1}^2 g_i(k) \log P(X=x_i; Z=k; \theta) \\ &+ \sum_{i=1}^n \sum_{k=1}^2 g_i(k) \log \frac{1}{g_i(k)} \end{aligned}$$

$$A^\top$$

$$A^\top \rightarrow A^\top A \text{ is PD}$$

$A \neq 0$

~~$A = 0$~~

$A \neq 0$

$$\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

$$\begin{aligned} & x^\top A^\top A x \cancel{=} \\ & = (Ax)^\top Ax > 0 \\ & \sum (A_{ii} x_i)^2 > 0 \\ & x^\top A^\top A x \end{aligned}$$

$$\begin{aligned} & Ax \quad n \times 1 = n \times 1 \quad = (Ax)^\top Ax \leq \|Ax\|_2^2 \geq 0 \\ & \left[ \begin{array}{c} A \\ \vdots \end{array} \right] \left[ \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array} \right] \quad \|Ax\|_2^2 > 0 \\ & \sum (a_i x_i)^2 \end{aligned}$$

$$Ax^\top = [a_1 \ a_2 \ \dots]$$

$$\begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}$$

$$Ax = Ax^\top, \quad \nabla_x x^\top a = a^\top$$

$$\arg \min_{x \in \mathbb{R}^n} \|x - y\|_2^2 \rightarrow \boxed{a^\top x} = a^\top$$

$$\mathcal{L}(\theta) \geq \mathcal{F}(\theta, q^*) \quad \rightarrow \quad \mathcal{L}(\theta) = \mathcal{F}(\theta, q)$$

E-Step:  $q^{t+1} \in \operatorname{argmax} \mathcal{F}(\theta^*, q)$

M-Step:  $\theta^{t+1} \in \operatorname{argmax} \mathcal{F}(\theta, q^{t+1})$

$q^{t+1}$  is a maximizer of  
 $\mathcal{F}(\theta, q) = \mathcal{L}(\theta)$

S-Step:  $q_i^{t+1}(k) = p(Z_i=k | X=x_i)$

$$\sum_{i=1}^n \log p(X=x_i; \theta) = \sum_{i=1}^n \sum_{k=1}^2 p(Z_i=k | X=x_i; \theta^*) \log \frac{p(X=x_i, Z_i=k; \theta^*)}{p(Z_i=k | X=x_i; \theta^*)}$$

$$\begin{aligned} \underbrace{q_i^{t+1}(k)}_{k=1} &= p(Z_i=k | X=x_i; \theta^*) = \frac{p(X=x_i, Z_i=k; \theta^*)}{p(X=x_i; \theta^*)} \\ &= w_k e^{-\lambda_k} \frac{\lambda_k^{x_i}}{(x_i)!} \\ &\quad \overline{w_k e^{-\lambda_k} \frac{\lambda_k^{x_i}}{(x_i)!} + ((1-w_k) e^{-\mu} \frac{\mu^{x_i}}{(x_i)!})} \end{aligned}$$

$k=2$ :

$$q_i^{t+1}(2) = 1 - q_i^{t+1}(1)$$

$$F(\theta, g_i^{(t+1)}) := \sum_{i=1}^n g_i^{(t+1)}(\cdot) \log [w - \lambda + x_i \log \lambda + \\ g_i^{(t+1)}(\cdot) \underbrace{\log ((1-w)e^{-\mu} \frac{\mu^{x_i}}{(x_i)!})}_{\log (1-w) - \mu + x_i \log \mu} + C$$

$$\nabla_w F = 0 \Rightarrow \sum_{i=1}^n \left( \frac{g_i^{(t+1)}(1)}{w} - \frac{g_i^{(t+1)}(2)}{1-w} = 0 \right)$$

$$\Rightarrow \sum_{i=1}^n g_i^{(t+1)}(1) - w \left( \sum_{i=1}^n g_i^{(t+1)}(1) + \sum_{i=1}^n g_i^{(t+1)}(2) \right) = 0$$

$$\Rightarrow w_{t+1} = \frac{\sum_{i=1}^n g_i^{(t+1)}(1)}{n}$$

$$\nabla_\lambda F = 0 \Rightarrow \sum_{i=1}^n g_i^{(t+1)}(1) \left( -1 + \frac{x_i}{\lambda} \right) = 0$$

$$\Rightarrow \sum_{i=1}^n \lambda g_i^{(t+1)}(1) = \sum_{i=1}^n x_i g_i^{(t+1)}(1)$$

$$\Rightarrow \lambda_{t+1} = \frac{\sum_{i=1}^n x_i g_i^{(t+1)}(1)}{\sum_{i=1}^n g_i^{(t+1)}(1)}$$

$$\mu_{t+1} = \frac{\sum_{i=1}^n x_i g_i^{(t+1)}(2)}{\sum_{i=1}^n g_i^{(t+1)}(2)}$$

- (b) (2 pts) Write down the expression for the lower bound  $\mathcal{F}(\lambda, \mu, w, q)$  where  $q = (q_1, \dots, q_n)$  where  $q = (q_1, \dots, q_n)$  and  $q_i = p(Z_i = 1)$  denotes the parameter for an appropriate distribution on the hidden label  $Z_i$ .

### 3 Weighted Least Squares (8 pts)

The weighted least squares (WLS) is a slight modification on least squares:

$$\min_{\theta} \sum_{i=1}^n w_i (x_i^\top \theta - y_i)^2 \quad (2)$$

i.e., the  $i$ -th term in the least squares objective is weighted by a term  $w_i$ .

Define the matrices:

$$\mathbf{W} = \begin{bmatrix} w_1 & & \\ & w_2 & \\ & & \ddots \\ & & & w_n \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix} \in \mathbb{R}^{n \times d} \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

- (a) (2 pts) Write the objective in equation (2) using the matrix-vector notation.

Your answer should be in terms of  $\mathbf{W}$ ,  $\mathbf{X}$ ,  $\mathbf{y}$  and  $\theta$ .

$$\begin{aligned}
 & \sum_{i=1}^n w_i (x_i^\top \theta - y_i)^2 \\
 &= \sum_{i=1}^n \underbrace{(\sqrt{w_i} x_i^\top \theta)}_{\tilde{x}_i^\top \theta} - \underbrace{\sqrt{w_i} y_i}_{\tilde{y}_i} \quad \text{OLS Objective} \\
 &= \| \tilde{\mathbf{X}} \theta - \tilde{\mathbf{y}} \|^2 = \mathbf{e}^\top \mathbf{e} \\
 &= \sum_{i=1}^n (\tilde{x}_i^\top \theta - \tilde{y}_i)^2 \quad \text{WLS Objective;} \\
 &= \|\tilde{\mathbf{X}}\theta - \tilde{\mathbf{y}}\|_2^2 \\
 &= (\tilde{\mathbf{X}}\theta - \tilde{\mathbf{y}})^\top \begin{pmatrix} \tilde{\mathbf{X}}\theta \\ \tilde{\mathbf{y}} \end{pmatrix} \\
 &= (\mathbf{X}\theta - \mathbf{y})^\top \sqrt{\mathbf{W}} \sqrt{\mathbf{W}} (\mathbf{X}\theta - \mathbf{y}) \\
 &= (\mathbf{X}\theta - \mathbf{y})^\top \mathbf{W} (\mathbf{X}\theta - \mathbf{y})
 \end{aligned}$$

$$b) \hat{\theta}_{OLS} = (X^T X)^{-1} X^T Y$$

$$\hat{\theta}_{WLS} = (X^T W X)^{-1} X^T W Y$$

$$c) y_i = x_i^T \theta + \epsilon_i$$

$$y_i - x_i^T \theta = \epsilon_i \sim N(0, \sigma_i^2)$$

$$\Rightarrow \hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} P(y_1, \dots, y_n | x_1, \dots, x_n; \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^n P(Y=y_i | X=x_i; \theta) \quad ?$$

$$= \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^n P(\epsilon_i = y_i - x_i^T \theta)$$

$$= \sum_{i=1}^n \underbrace{(y_i - x_i^T \theta)^2}_{= \sigma_i^2}$$

$$w_i = \frac{1}{\sigma_i^2} \text{ or } \frac{1}{2\sigma^2}.$$

## 4 Ridge Regression vs PCA (12 pts)

Assume we are given  $n$  training data points  $(\mathbf{x}_i, y_i)$ . We collect the responses into  $\mathbf{y} \in \mathbb{R}^n$ , and the features into the matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  where the rows are the  $d$ -dimensional feature vectors  $\mathbf{x}_i^\top$  corresponding to each training point. Furthermore, assume that  $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \mathbf{0}$ ,  $n > d$  and  $\mathbf{X}$  has rank  $d$ .

In this problem we want to compare two procedures: The first is ridge regression with hyper-parameter  $\lambda$ , while the second is applying ordinary least squares after using PCA to reduce the feature dimension from  $d$  to  $k$  (we give this latter approach the short-hand name  $k$ -PCA-OLS where  $k$  is the hyper-parameter).

Notation: The singular value decomposition of  $\mathbf{X}$  reads  $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$  where  $\mathbf{U} \in \mathbb{R}^{n \times n}$ ,  $\Sigma \in \mathbb{R}^{n \times d}$  and  $\mathbf{V} \in \mathbb{R}^{d \times d}$ . We denote by  $\mathbf{u}_j$  the  $n$ -dimensional column vectors of  $\mathbf{U}$  and by  $\mathbf{v}_j$  the  $d$ -dimensional column vectors of  $\mathbf{V}$ . Furthermore the diagonal entries  $\sigma_j = \Sigma_{j,j}$  of  $\Sigma$  satisfy  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d > 0$ . For notational convenience, assume that  $\sigma_j = 0$  for  $j > d$ .

- (a) (4 pts) It turns out that the ridge regression optimizer (with  $\lambda > 0$ ) in the  $\mathbf{V}$ -transformed coordinates

$$\hat{\theta}_{\text{ridge}} = \arg \min_{\theta} \|\mathbf{X}\mathbf{V}\theta - \mathbf{y}\|_2^2 + \lambda \|\theta\|_2^2$$

has the following expression:

$$\hat{\theta}_{\text{ridge}} = \text{diag}\left(\frac{\sigma_j}{\lambda + \sigma_j^2}\right) \mathbf{U}^\top \mathbf{y}. \quad (4)$$

Use  $\hat{y}_{\text{test}} = \mathbf{x}_{\text{test}}^\top \mathbf{V} \hat{\theta}_{\text{ridge}}$  to denote the resulting prediction for a hypothetical  $\mathbf{x}_{\text{test}}$ . Using (4) and the appropriate scalar  $\{\alpha_i\}$ , this can be written as:

$$\hat{y}_{\text{test}} = \mathbf{x}_{\text{test}}^\top \sum_{j=1}^d \mathbf{v}_j \alpha_j \mathbf{u}_j^\top \mathbf{y}. \quad (5)$$

What are the  $\alpha_j \in \mathbb{R}$  for this to correspond to (4) from ridge regression?  
Show your work.

$$\begin{aligned} \hat{y}_{\text{test}} &= \mathbf{x}_{\text{test}}^\top \mathbf{V} \hat{\theta}_{\text{ridge}} \\ &= \mathbf{x}_{\text{test}}^\top \sum_{j=1}^d \mathbf{v}_j \alpha_j \mathbf{u}_j^\top \mathbf{y} \\ \hat{\theta}_{\text{ridge}} &= \begin{pmatrix} \sigma_1 \\ \sigma_1^2 + \lambda \\ \vdots \\ \sigma_d \\ \sigma_d^2 + \lambda \end{pmatrix} \end{aligned}$$

$$\alpha_j = ?$$

$$\mathbf{V} \hat{\theta}_{\text{ridge}} = \sum_{j=1}^d \mathbf{v}_j \alpha_j \mathbf{u}_j^\top \mathbf{y}$$

$$\begin{aligned} &= \mathbf{V} \tilde{\mathbf{D}} \mathbf{V}^\top \mathbf{y} \\ &= [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_d] \begin{bmatrix} \frac{\sigma_1}{\sigma_1^2 + \lambda} & & & \\ & \ddots & & \\ & & \frac{\sigma_d}{\sigma_d^2 + \lambda} & \\ & & & \end{bmatrix} \begin{bmatrix} \mathbf{u}_1^\top \mathbf{y} \\ \vdots \\ \mathbf{u}_d^\top \mathbf{y} \end{bmatrix} \end{aligned}$$

$$= \begin{bmatrix} v_1 & v_2 & \dots & v_d \end{bmatrix} \begin{bmatrix} \frac{\sigma_1}{\sigma_1^2 + \lambda} v_1^T y \\ \vdots \\ \frac{\sigma_d}{\sigma_d^2 + \lambda} v_d^T y \end{bmatrix} = \sum_{j=1}^d v_j \frac{\sigma_j}{\sigma_j^2 + \lambda}$$

$$\alpha_j = \frac{\sigma_j}{\sigma_j^2 + \lambda}$$

$$\hat{\theta}_{\text{Ridge}} : \alpha_j = \frac{\sigma_j}{\sigma_j^2 + \lambda}$$

$$\hat{\theta}_{\text{OLS}} : \alpha_j|_{\lambda=0} = \frac{\sigma_j}{\sigma_j^2} = \frac{1}{\sigma_j}$$

- (b) (5 pts) Suppose that we do k-PCA-OLS — i.e. ordinary least squares on the reduced  $k$ -dimensional feature space obtained by projecting the raw feature vectors onto the  $k < d$  principal components of the covariance matrix  $\mathbf{X}^\top \mathbf{X}$ . Use  $\hat{y}_{test}$  to denote the resulting prediction for a hypothetical  $\mathbf{x}_{test}$ ,

It turns out that the learned k-PCA-OLS predictor can be written as:

$$\hat{y}_{test} = \mathbf{x}_{test}^\top \sum_{j=1}^d \mathbf{v}_j \alpha_j \mathbf{u}_j^\top \mathbf{y}. \quad (6)$$

Give the  $\alpha_j \in \mathbb{R}$  coefficients for k-PCA-OLS. Show work. These  $\alpha_j$ 's would be different compared to the previous part.

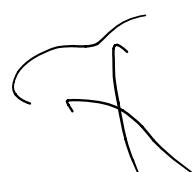
*Hint 1: some of these  $\alpha_j$  will be zero. Also, if you want to use the compact form of the SVD, feel free to do so if that speeds up your derivation.*

*Hint 2: some inspiration may be possible by looking at the next part for an implicit clue as to what the answer might be.*

$$\hat{y}_{test} = \mathbf{x}_{test}^\top \sum_{j=1}^d \mathbf{v}_j \alpha_j \mathbf{u}_j^\top \mathbf{y}$$

$$\alpha_j = ? \quad \begin{cases} \frac{1}{\sigma_j} & \text{for } j \leq k \\ 0 & \text{for } j > k. \end{cases}$$

Top  $K$  PCs are corresponding to  $\mathbf{e}_1, \dots, \mathbf{e}_k$   
 and in V-Space,  $\hat{\theta}_j$  is independent of  $\hat{\theta}_i$   
 for  $i \neq j$ .



$$V = \{V_1, \dots, V_k, V_{k+1}, \dots, V_d\}$$

$$\min \| \underbrace{XV}_{\text{circled}} \theta - y \|_2^2$$

$$V = [v_1, \dots, v_k, \dots, v_d]$$

$$X = U \Sigma V^T$$

$$X^T X = V \Sigma^2 V^T \quad \checkmark$$

$$\hat{\theta} = (X^T X)_{kk}^{-1} V_k^T X^T y$$

$$\hat{\theta} = (V_k^T V \Sigma^2 V^T V_k)^{-1} V_k^T V \Sigma U^T y \quad = [V_k^T \theta]$$

$$= (\underbrace{[I_n, 0]}_{= \begin{pmatrix} I_n & 0 \\ 0 & 0 \end{pmatrix}} \begin{pmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_d^2 \end{pmatrix} \begin{pmatrix} I_n \\ 0 \end{pmatrix})^{-1} (\underbrace{[I_n, 0]}_{= \begin{pmatrix} I_n & 0 \\ 0 & 0 \end{pmatrix}} \begin{pmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_d^2 \end{pmatrix} V^T y)$$

$$y_{\text{ren}} = X_{\text{ren}} V_k \hat{\theta}$$

$$V_k^T V = V_k^T [V_k \quad V_{-k}]$$

$$= \begin{pmatrix} V_k^T V_k & V_{-k}^T V_k \\ V_k^T V_{-k} & V_{-k}^T V_{-k} \end{pmatrix} = [I_n, 0]$$

$$\begin{pmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_d^2 \\ & & & 0 \end{pmatrix} \begin{pmatrix} I_n \\ 0 \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_d^2 \\ & & & 0 \end{pmatrix}$$

$$(A \ B) \begin{pmatrix} C \\ D \end{pmatrix} = AC + BD$$

$$= \begin{pmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_d^2 \end{pmatrix}^{-1} \begin{pmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \ddots & 0 \end{pmatrix} \begin{bmatrix} u_1^T y \\ \vdots \\ u_d^T y \end{bmatrix}$$

$$* = \begin{pmatrix} \frac{1}{\sigma_1} & & & \\ & \ddots & & \\ & & \frac{1}{\sigma_k} & \\ & & & 0 \end{pmatrix} \begin{pmatrix} u_1^T y \\ \vdots \\ u_k^T y \\ y^T y \end{pmatrix} = \begin{pmatrix} u_1^T y / \sigma_1 \\ u_2^T y / \sigma_2 \\ \vdots \\ u_k^T y / \sigma_k \end{pmatrix}$$

$$V_k \hat{\theta} = [v_1, \dots, v_k] \begin{bmatrix} u_1^T y / \sigma_1 \\ \vdots \\ u_k^T y / \sigma_k \end{bmatrix}$$

$$= \sum_{j=1}^k v_j \frac{1}{\sigma_j} u_j^T y = \text{Xren} \sum_{j=1}^d v_j d_j u_j^T y.$$

Redo Sample Midterm:

1. a) True - Predictions are only meaningful if our data is representative of current & future population.

1. b) False. Its value = k

$$H\hat{Y} = \hat{Y} = X\hat{\beta} = X(X^T X)^{-1} X^T Y$$

$$H = X(X^T X)^{-1} X^T$$

1. c) False - Local minimizer

1. d) False -  $\nabla_{\theta} (\theta^T A \theta) = \theta^T A^T + \theta^T A \neq 2A\theta$

1. e) True - Linearly dependent  $\rightarrow$  Not full rank

1. f) True - See some books to see

1. g) True -

1. h) False - No relationship  $\rightarrow$  Ridge

1. i) False - PCA is unsupervised learning  $\rightarrow$  not used for prediction

1.j) True - Only scale when units are different

2. a)

$$\begin{aligned}
 P(X) &= \frac{P(X, Z)}{P(Z)} = \sum_{i=1}^2 P(Z=z_i) = 1 \\
 &\sum_{i=1}^n \sum_{j=1}^2 P(X=x_i, Z=z_j) \\
 &= \sum_{i=1}^n \left[ P(X=x_i, Z=1) + P(X=x_i, Z=2) \right] \\
 &= \sum_{i=1}^n \left[ P(X=x_i | Z=1) P(Z=1) + P(X=x_i | Z=2) P(Z=2) \right] \\
 &= \sum_{i=1}^n \left[ e^{-x_i} \frac{\lambda^{x_i}}{(x_i)!} w + e^{-x_i} \frac{\mu^{x_i}}{(x_i)!} (1-w) \right] \\
 &= \sum_{i=1}^n \frac{e^{-x_i}}{(x_i)!} (\lambda^{x_i} w + \mu^{x_i} (1-w)) \\
 &\Rightarrow \sum_{i=1}^n \frac{e^{-x_i}}{(x_i)!} (\lambda^{x_i} w + \mu^{x_i} (1-w)) \\
 &\Rightarrow \sum_{i=1}^n -x_i + \log \left( \frac{\lambda^{x_i}}{(x_i)!} w \right) - x_i - \log \left( \frac{\mu^{x_i}}{(x_i)!} (1-w) \right)
 \end{aligned}$$

$$\begin{aligned}
 3.a) \quad & \min_{\theta} \sum_{i=1}^n w_i (x_i^\top \theta - y_i)^2 \\
 & = \min_{\theta} \left( w_1 (x_1^\top \theta - y_1)^2 + w_2 (x_2^\top \theta - y_2)^2 + \dots + w_n (x_n^\top \theta - y_n)^2 \right) \\
 & = \min_{\theta} [ \|w \parallel x \theta - y \parallel_2^2 ] \quad (x^\top w x)^\top x^\top w y \\
 & = \min_{\theta} [(x \theta - y)^\top w (x \theta - y)]
 \end{aligned}$$

$$\begin{aligned}
 3.b) \quad & \min_{\theta} [(x \theta - y)^\top w (x \theta - y)] \\
 & \nabla_{\theta} (\theta^\top x^\top w x \theta - \theta^\top x^\top w y - y^\top w x \theta - y^\top w y) \\
 & = \theta^\top x^\top w x + \theta^\top x^\top w x - y^\top w^\top x - y^\top w x \\
 & = 2 \theta^\top w x - 2 y^\top w^\top x = 0 \\
 & \theta^\top w x = y^\top w^\top x
 \end{aligned}$$

$$3.a) \min_{\theta} [\|w\| \|x\theta - y\|_2^2]$$

$$= \min_{\theta} [(x\theta - y)^T w (x\theta - y)]$$

$$\Rightarrow \nabla_{\theta} [\theta^T \underline{x^T w x} \theta - \theta^T x^T w y - y^T w x \theta + y^T w y]$$

$$= 2\theta^T x^T w x - y^T w^T x - y^T w x = 0$$

$$2\theta^T x^T w x = y^T (w^T + w) x$$

~~$$2\theta^T x^T w x = 2y^T w x$$~~

$$\theta^T = (x^T w x)^{-1} y^T w x$$

$$\theta = (x^T w x)^{-1} x^T w y$$

$$3.c) \quad y_i = x_i^T \theta + \varepsilon_i$$

Maximum Likelihood

$$\underset{\theta}{\operatorname{argmax}} \quad P(y_1=y_1, y_2=y_2, \dots | x_1=x_1, x_2=x_2, \dots)$$

$$= \underset{\theta}{\operatorname{argmax}} \quad P(\varepsilon_1=y_1-x_1^T \theta; \varepsilon_2=y_2-x_2^T \theta \dots)$$

$$= \underset{\theta}{\operatorname{argmax}} \left[ \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_1-x_1^T \theta)^2}{2\sigma^2}} \dots \right]$$

$$= \underset{\theta}{\operatorname{argmax}} \left[ \sum_{i=1}^n -\frac{(y_i-x_i^T \theta)^2}{2\sigma^2} \right]$$

$$= \underset{\theta}{\operatorname{argmin}} \left[ \sum_{i=1}^n w_i (x_i^T \theta - y_i)^2 \right]$$

$$w_i = \frac{1}{2\sigma_i^2}$$

4.a)

$$\hat{y}_{\text{true}} = \cancel{x^T} \cancel{V} \hat{\theta}_{\text{ridge}} = \cancel{x^T} \cancel{V} \sum_{j=1}^d v_j \alpha_j u_j^T y$$

$$V \hat{\theta}_{\text{ridge}} = V \text{diag} \left( \frac{\sigma_i}{\lambda + \sigma_i^2} \right) U^T y$$

$$\begin{aligned} \sum_{j=1}^d v_j \alpha_j u_j^T y &= v_1 \alpha_1 u_1^T y + v_2 \alpha_2 u_2^T y + \dots \\ &\quad + v_d \alpha_d u_d^T y \\ &= \left( \sum_{j=1}^d v_j \alpha_j \right) U^T y \\ &= V \left( \sum_{j=1}^d \alpha_j \right) U^T y \end{aligned}$$

$$\alpha_1 + \alpha_2 + \dots + \alpha_d = \text{diag} \left( \frac{\sigma_i}{\lambda + \sigma_i^2} \right)$$

4.a)

$$\hat{y}_{\text{err}} = \hat{x}^\top \text{err} V \hat{\theta}_{\text{ridge}}$$

$$V \hat{\theta}_{\text{ridge}} = V \text{diag}\left(\frac{\sigma^2}{\sigma^2 + \lambda}\right) U^\top Y$$

$$V \hat{\theta}_{\text{ridge}} = V \begin{bmatrix} \frac{\sigma^2}{\sigma^2 + \lambda} & & \\ & \ddots & \\ & & \frac{\sigma^2}{\sigma^2 + \lambda} \end{bmatrix}$$

#### 4 Ridge Regression vs PCA (12 pts)

Assume we are given  $n$  training data points  $(\mathbf{x}_i, y_i)$ . We collect the responses into  $\mathbf{y} \in \mathbb{R}^n$ , and the features into the matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  where the rows are the  $d$ -dimensional feature vectors  $\mathbf{x}_i^\top$  corresponding to each training point. Furthermore, assume that  $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \mathbf{0}$ ,  $n > d$  and  $\mathbf{X}$  has rank  $d$ .

In this problem we want to compare two procedures: The first is ridge regression with hyper-parameter  $\lambda$ , while the second is applying ordinary least squares after using PCA to reduce the feature dimension from  $d$  to  $k$  (we give this latter approach the short-hand name  $k$ -PCA-OLS where  $k$  is the hyper-parameter).

Notation: The singular value decomposition of  $\mathbf{X}$  reads  $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$  where  $\mathbf{U} \in \mathbb{R}^{n \times n}$ ,  $\Sigma \in \mathbb{R}^{n \times d}$  and  $\mathbf{V} \in \mathbb{R}^{d \times d}$ . We denote by  $\mathbf{u}_j$  the  $n$ -dimensional column vectors of  $\mathbf{U}$  and by  $\mathbf{v}_j$  the  $d$ -dimensional column vectors of  $\mathbf{V}$ . Furthermore the diagonal entries  $\sigma_j = \Sigma_{j,j}$  of  $\Sigma$  satisfy  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d > 0$ . For notational convenience, assume that  $\sigma_j = 0$  for  $j > d$ .

- (a) (4 pts) It turns out that the ridge regression optimizer (with  $\lambda > 0$ ) in the  $\mathbf{V}$ -transformed coordinates

$$\hat{\theta}_{\text{ridge}} = \arg \min_{\theta} \|\mathbf{X}\mathbf{V}\theta - \mathbf{y}\|_2^2 + \lambda \|\theta\|_2^2$$

has the following expression:

$$\hat{\theta}_{\text{ridge}} = \text{diag}\left(\frac{\sigma_j}{\lambda + \sigma_j^2}\right) \mathbf{U}^\top \mathbf{y}. \quad (4)$$

$d \times d \quad n \times n \quad n \times 1$

Use  $\hat{y}_{\text{test}} = \mathbf{x}_{\text{test}}^\top \hat{\theta}_{\text{ridge}}$  to denote the resulting prediction for a hypothetical  $\mathbf{x}_{\text{test}}$ . Using (4) and the appropriate scalar  $\{\alpha_i\}$ , this can be written as:

$$\hat{y}_{\text{test}} = \mathbf{x}_{\text{test}}^\top \sum_{j=1}^d \mathbf{v}_j \alpha_j \mathbf{u}_j^\top \mathbf{y}. \quad (5)$$

What are the  $\alpha_j \in \mathbb{R}$  for this to correspond to (4) from ridge regression?  
Show your work.

$$\begin{aligned} \cancel{\mathbf{x}^\top \hat{\theta}_{\text{ridge}}} &= \cancel{\mathbf{x}^\top \sum_{j=1}^d \mathbf{v}_j \alpha_j \mathbf{u}_j^\top \mathbf{y}} \\ \hat{\theta}_{\text{ridge}} &= \sum_{j=1}^d \mathbf{v}_j \alpha_j \mathbf{u}_j \\ \begin{bmatrix} 1 & | & | & | \\ \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_d \\ 1 & | & | & | \end{bmatrix} &\quad \begin{bmatrix} \frac{\sigma_j}{\lambda + \sigma_j^2} \\ \vdots \\ \frac{\sigma_d}{\lambda + \sigma_d^2} \end{bmatrix} \quad \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \\ \vdots \\ \mathbf{u}_n \end{bmatrix} \quad \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_n \end{bmatrix} \\ d \times d &\quad d \times n \quad n \times n \quad n \times 1 \end{aligned}$$

$$\begin{array}{c} \text{diag} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} \begin{pmatrix} a_1 & 0 \\ a_2 & 0 \\ a_3 & 0 \end{pmatrix} \\ \left[ \begin{array}{ccc} a_1 & & 0 \\ a_2 & & 0 \\ a_3 & & 0 \end{array} \right] \end{array}$$

$$V\hat{\theta}_{\text{ridge}} = \sum_{j=1}^d v_j \alpha_j u_j y$$

$$\text{LHS} = V\hat{\theta}_{\text{ridge}}$$

$$= V \text{diag}\left(\frac{\sigma_j}{\sigma_j^2 + \lambda}\right) U^T y$$

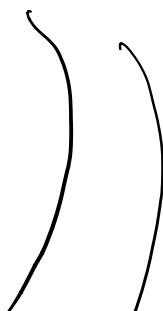
$$= \begin{bmatrix} | & | & | \\ v_1 & v_2 & \dots & v_d \\ | & | & | \end{bmatrix} \begin{bmatrix} \frac{\sigma_1}{\sigma_1^2 + \lambda} & & & \\ & \ddots & & \\ & & \frac{\sigma_d}{\sigma_d^2 + \lambda} & \\ & & & 0 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$\text{RHS} = \sum_{j=1}^d v_j \alpha_j u_j$$

$$= v_1 \alpha_1 u_1 y_1 + v_2 \alpha_2 u_2 y_2 + \dots$$

=

$$= v_1 \cdot \frac{\sigma_1}{\sigma_1^2 + \lambda} \cdot u_1 y_1$$



$$\begin{array}{c}
 \text{diag}(\mathbf{\Sigma}) \quad \mathbf{U}^T \mathbf{y} \\
 \underbrace{\hspace{10em}}_{\text{nx1}}
 \end{array}
 \quad
 \begin{array}{l}
 \mathbf{u} \in \mathbb{R}^n \\
 \mathbf{d} \in \mathbb{R}^m
 \end{array}$$

$$\mathbf{V} \hat{\boldsymbol{\beta}}_{\text{ridge}} = \mathbf{d} \in \mathbb{R}^m \quad \text{nx1}$$

$$\mathbf{X} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T$$

$\mathbf{U}$   $\boldsymbol{\Sigma}$   $\mathbf{V}$   
 $n \times n$        $n \times d$        $d \times m$

- (b) (5 pts) Suppose that we do k-PCA-OLS — i.e. ordinary least squares on the reduced  $k$ -dimensional feature space obtained by projecting the raw feature vectors onto the  $k < d$  principal components of the covariance matrix  $\mathbf{X}^\top \mathbf{X}$ . Use  $\hat{y}_{test}$  to denote the resulting prediction for a hypothetical  $\mathbf{x}_{test}$ ,

It turns out that the learned k-PCA-OLS predictor can be written as:

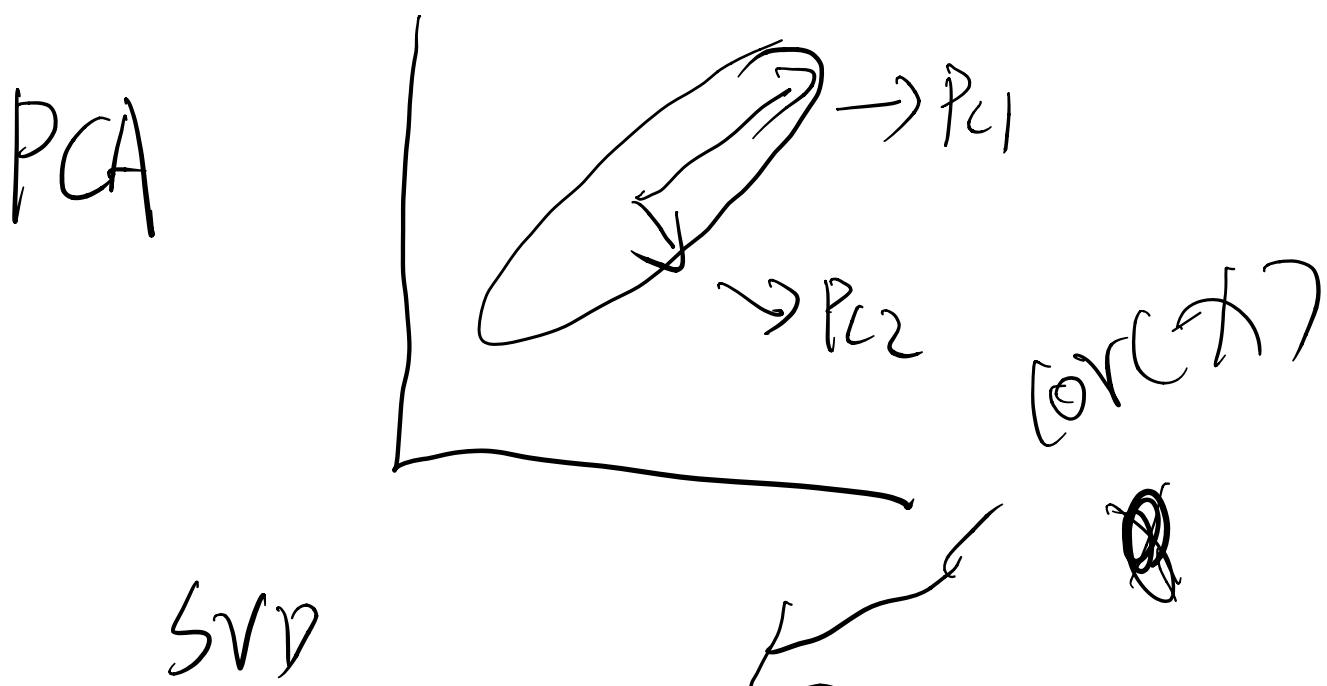
$$\hat{y}_{test} = \mathbf{x}_{test}^\top \sum_{j=1}^d \mathbf{v}_j \alpha_j \mathbf{u}_j^\top \mathbf{y}. \quad (6)$$

Give the  $\alpha_j \in \mathbb{R}$  coefficients for k-PCA-OLS. Show work. These  $\alpha_j$ 's would be different compared to the previous part.

*Hint 1: some of these  $\alpha_j$  will be zero. Also, if you want to use the compact form of the SVD, feel free to do so if that speeds up your derivation.*

*Hint 2: some inspiration may be possible by looking at the next part for an implicit clue as to what the answer might be.*

4.b)  $\hat{y}_{test} = \mathbf{x}_{test}^\top \sum_{j=1}^d \mathbf{v}_j \alpha_j \mathbf{u}_j^\top \mathbf{y}$



$$X = U\Sigma V^T \quad X^T X = V \Sigma^2 V^T$$

$\Sigma$

eigenvalue

Scores aka principal components

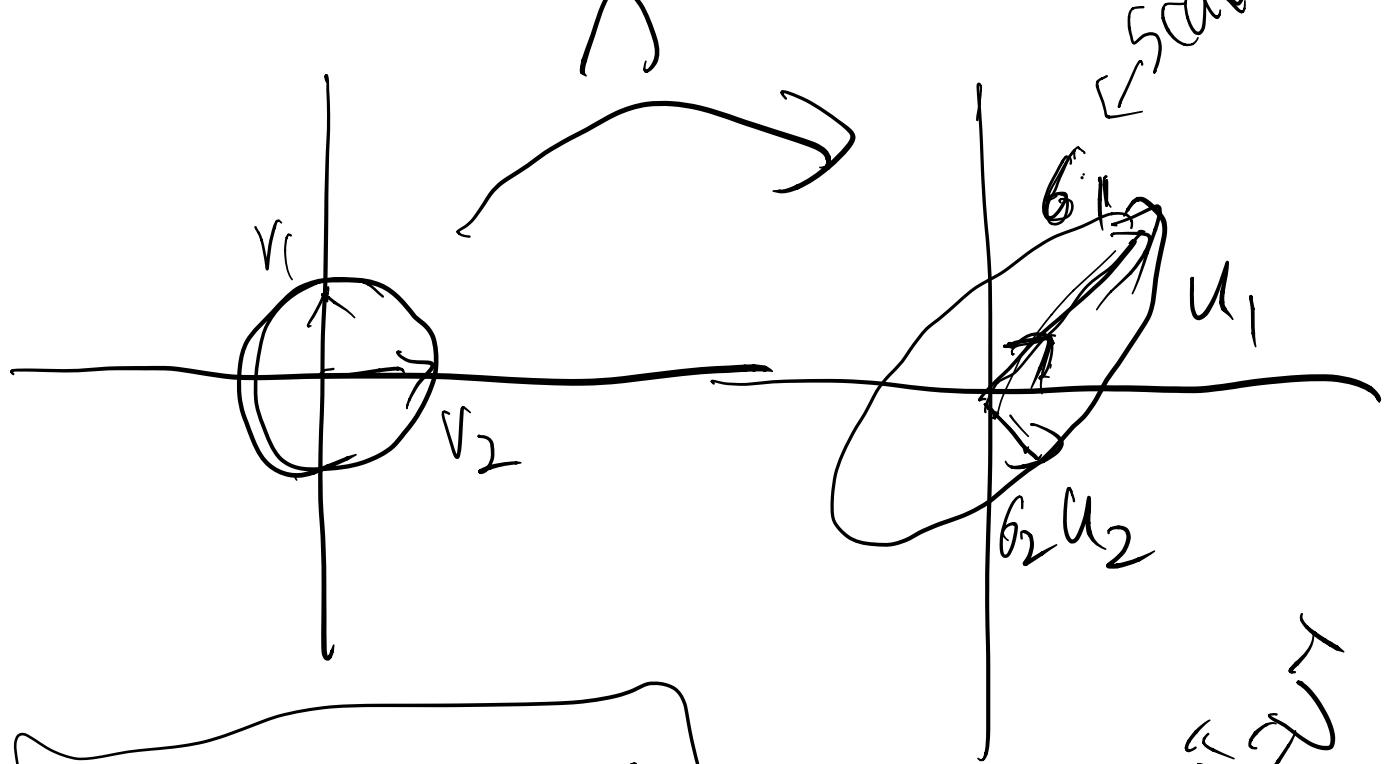
$V$  is called rotation matrix  $\rightarrow$  loadings

its eigenvectors of  $X^T X \rightarrow$  principal vectors

$\lambda_1 = 6$

$v_1$  unit vector

$X$



$$XV = US$$

$[V, U]$

$U^T$   $U$   $U^T$   $U$

$X = USV^T$



$(x_1, \dots, x_n)$





$$y = X\beta \quad \left\| Y - XU\beta \right\|_2^2 = L(\beta)$$

$\hat{\beta}_{OLS}$

$$\hat{\beta}_{OLS} = (X^T X)^{-1} X^T Y$$

$$(V^T X^T X V)^{-1} V^T X^T Y$$

$$= (V^T V \Sigma^2 V^T V)^{-1} V^T X^T Y$$

$$X^T X = \underbrace{\dots}_{n \times n \text{ and } d \times d}$$

$$X = U \Sigma V^T = \underbrace{(\Sigma^2)^{-1} V^T X^T Y}_{d \times d}$$

$$X^T X = \overbrace{V \Sigma^2 V^T}^V U^T U \Sigma V^T$$

$$X^T X = V \Sigma^2 V^T$$

$$= (\Sigma)^{-1} \cancel{\frac{1}{n} V^T U^T X} \\ = (\Sigma)^{-1} \Sigma^T U^T Y$$

$$\xrightarrow{n \times d} \begin{bmatrix} v_{11} & \dots & v_{1d} \\ \vdots & \ddots & \vdots \\ v_{n1} & \dots & v_{nd} \end{bmatrix} \xrightarrow{d \times d} \begin{bmatrix} \Sigma_{11} & \dots & \Sigma_{1d} \\ \vdots & \ddots & \vdots \\ \Sigma_{d1} & \dots & \Sigma_{dd} \end{bmatrix} \xrightarrow{d \times 1} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$3.6) \min_{\theta} \| \underline{J} \bar{w} X \theta - y \|_2^2$$

$$\hat{\theta}_{OLS} = (X^T X)^{-1} X^T Y$$

$$= (X^T \bar{w}^T \bar{w} X)^{-1} \bar{w} X^T Y$$

$$= (X^T w X)^{-1} \bar{w} X^T Y$$

$\Sigma = \Sigma^{-1}$

$$(\Sigma)^{-1} \Sigma^T U^T Y = \sum U_j \alpha_j \underline{U_j^T Y}$$

$$\begin{bmatrix} \frac{1}{\sigma_1^2} & & & \\ & \frac{1}{\sigma_2^2} & & \\ & & \ddots & \\ & & & \frac{1}{\sigma_d^2} \end{bmatrix} \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_d \end{bmatrix} \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$$

$$\begin{bmatrix} \frac{1}{\sigma_1^2} & & & \\ & \frac{1}{\sigma_2^2} & & \\ & & \ddots & \\ & & & \frac{1}{\sigma_d^2} \end{bmatrix} = \begin{bmatrix} 1 & y_1 & \dots & y_n \\ u_1 & u_2 & \dots & u_n \\ 1 & 1 & \dots & 1 \end{bmatrix} \alpha,$$

$$\alpha_j = \frac{1}{\sigma_j} u_j \quad \alpha_0 = \frac{1}{\sigma_0}$$

$\hat{Y}$

$$F = \hat{Y}^T Y$$

$$X V \hat{\beta} = X V \alpha U^T Y$$

~~$$\hat{X} V (\Sigma^{-2}) \Sigma^T U^T Y = \cancel{X} \cancel{V} \cancel{\alpha} \underline{U^T Y}$$~~

$$\alpha = (\Sigma^{-2}) \Sigma^T$$

normal

~~$$\alpha = \begin{bmatrix} \frac{1}{\sigma_1^2} & & & \\ & \frac{1}{\sigma_2^2} & & \\ & & \ddots & \\ & & & \frac{1}{\sigma_n^2} \end{bmatrix} \begin{bmatrix} \sigma_1 \\ \sigma_2 \\ \vdots \\ \sigma_n \end{bmatrix}$$~~

~~$$\alpha \rightarrow I$$~~

$$\alpha_i = \frac{1}{\sigma_i}$$

KLD

$$V = [V_1 \dots V_K, 0, 0]$$

$$\hat{Y} = XV\hat{\beta}_{PCA}$$

$$= X V \cdot \sum_{n=1}^2 \sum_{d=1}^r U^T Y_{dn}^{nxn}$$

||

$$\left( \frac{1}{\sigma_j^2}, \dots \right) \left( \begin{matrix} \theta_1 \\ \vdots \\ \theta_r \end{matrix} \right)$$

$$= X V \cdot \left( \begin{matrix} \theta_1 \\ \vdots \\ \theta_r \end{matrix} \right) U^T Y$$

$$= \text{Até} \sum_{i=1}^d V_i \left( \begin{matrix} \theta_1 \\ \vdots \\ \theta_r \end{matrix} \right) u_i v_i^T$$

$$\left. \begin{array}{c} i=0 \\ \vdots \\ i=k \\ \theta \end{array} \right\}$$