

## Sample Midterm Solutions

### 1. *Whitening.*

- a. Let  $\Sigma$  and  $\Sigma'$  be the sample covariance matrices of the original and the expanded data respectively. We have

$$\Sigma' := \frac{1}{2n} \sum_{i=1}^n (x_i x_i^T + x_i x_i^T) \quad (1)$$

$$= \frac{1}{n} \sum_{i=1}^n x_i x_i^T \quad (2)$$

$$= \Sigma. \quad (3)$$

The covariance matrices are the same, so nothing changes.

- b. Let  $\Sigma$  and  $\Sigma'$  be the sample covariance matrices of the original and the transformed respectively. The sample covariance equals

$$\Sigma' = \frac{1}{n-1} \sum_{i=1}^n A x_i (A x_i)^T \quad (4)$$

$$= A \left( \frac{1}{n-1} \sum_{i=1}^n x_i x_i^T \right) A^T \quad (5)$$

$$= A \Sigma A^T \quad (6)$$

$$= A U \Lambda U^T A^T, \quad (7)$$

where  $U \Lambda U^T$  is the eigendecomposition of  $\Sigma$ . Setting  $A = U^T$  results in a diagonal sample covariance matrix.

- c. By the argument above, setting  $B := (\sqrt{\Lambda})^{-1} U^T$ , where  $\sqrt{\Lambda}$  is a diagonal matrix that contains the square root of the eigenvalues, results in a sample covariance matrix equal to

$$\Sigma' = B U \Lambda U^T B^T \quad (8)$$

$$= (\sqrt{\Lambda})^{-1} U^T U \Lambda U^T U (\sqrt{\Lambda})^{-1} \quad (9)$$

$$= I. \quad (10)$$

- d. It doesn't make a difference. Let  $X = U S V^T$  be the SVD of the original feature matrix. Then  $B := S^{-1} U^T$ . The OLS estimate with the modified data equals

$$y_{\text{OLS}} = (B X)^T (B X X^T B^T)^{-1} B X y \quad (11)$$

$$= X^T B^T B X y \quad (12)$$

$$= V V^T y, \quad (13)$$

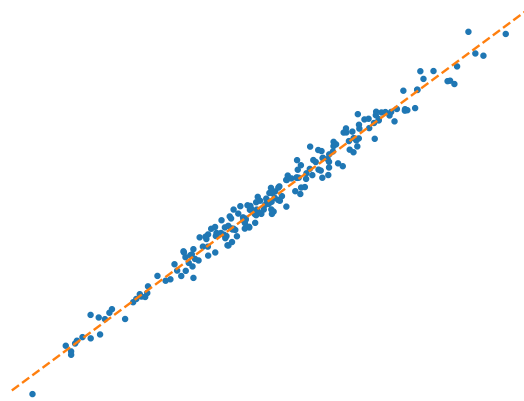
which is the exact same prediction as if we use  $X$  as the feature matrix.

### 2. *Quadratic form.*

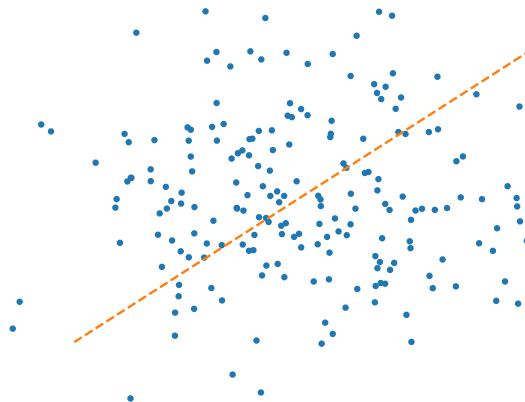
- a. By the spectral theorem the eigenvalues are the maximum and minimum of the quadratic form on the unit circle. These occur at points where the contour lines are tangent to the unit circle, so  $\lambda_1 = 1$  and  $\lambda_2 = -5$ .
- b. No. The quadratic form takes negative values, which would imply that the variance in certain directions is negative, which is impossible.
- c. No. The gradient equals  $2Av$ , so if the gradient is zero for a nonzero  $v$ , then this implies that  $A$  has a null space. However, this is impossible because it has two nonzero eigenvalues and is consequently full rank.

### 3. PCA.

a.



b.



- c. No. By linearity of expectation, the mean of the random variable that represents the data

equals

$$\mathbb{E}(\tilde{y}) = \mathbb{E}(\tilde{x}v + \tilde{z}) \quad (14)$$

$$= \mathbb{E}(\tilde{x})v + \mathbb{E}(\tilde{z}) \quad (15)$$

$$= 0. \quad (16)$$

d. By linearity of expectation and the fact that the mean of  $\tilde{x}$  and  $\tilde{z}$  is zero,

$$\mathbb{E}(\tilde{y}\tilde{y}^T) = \mathbb{E}((\tilde{x}v + \tilde{z})(\tilde{x}v + \tilde{z})^T) \quad (17)$$

$$= \mathbb{E}(\tilde{x}^2)vv^T + \mathbb{E}(\tilde{z}\tilde{z}^T) \quad (18)$$

$$= \sigma_{\text{signal}}^2 vv^T + \sigma_{\text{noise}}^2 I. \quad (19)$$

e. Setting  $U := [v/\|v\|_2 \ u_2 \ \cdots \ u_d]$ , we have  $UU^T = I$  by the orthonormality assumption. We can therefore rewrite the covariance matrix as

$$\mathbb{E}(\tilde{y}\tilde{y}^T) = \sigma_{\text{signal}}^2 vv^T + \sigma_{\text{noise}}^2 UU^T \quad (20)$$

$$= U \begin{bmatrix} \sigma_{\text{signal}}^2 \|v\|^2 + \sigma_{\text{noise}}^2 & 0_{d-1} \\ 0_{d-1} & \sigma_{\text{noise}}^2 I_{d-1} \end{bmatrix} U^T, \quad (21)$$

where  $0_{d-1}$  is a vector of  $d-1$  zeros and  $I_{d-1}$  is the  $d-1 \times d-1$  identity matrix.

f. Apply PCA and use the first principal direction as an estimate of the direction of  $v$ .

#### 4. Interference.

a. Note that all random variables are zero mean. Let  $\tilde{a}$  be equal to one if  $\tilde{x}[i] = y$  and to zero otherwise. By iterated expectation

$$\text{Var}(\tilde{x}[i]) = \mathbb{E}(\tilde{x}[i]^2) \quad (22)$$

$$= \mathbb{E}(\tilde{x}[i]^2 \mid \tilde{a} = 1)P(\tilde{x}[i] = y) + \mathbb{E}(\tilde{x}[i]^2 \mid \tilde{a} = 0)P(\tilde{x}[i] \neq y) \quad (23)$$

$$= \frac{\sigma^2 + \sigma^2}{2} \quad (24)$$

$$= \sigma^2. \quad (25)$$

Similarly, let  $\tilde{b}$  be equal to one if both  $\tilde{x}[1]$  and  $\tilde{x}[2]$  equal  $y$  and to zero otherwise. By iterated expectation

$$\text{Cov}(\tilde{x}[1]\tilde{x}[2]) = \mathbb{E}(\tilde{x}[1]\tilde{x}[2]) \quad (26)$$

$$= \mathbb{E}(\tilde{x}[1]\tilde{x}[2] \mid \tilde{b} = 1)P(\tilde{x}[1] = y)P(\tilde{x}[2] = y) \quad (27)$$

$$+ \mathbb{E}(\tilde{x}[1]\tilde{x}[2] \mid \tilde{b} = 0)(1 - P(\tilde{x}[1] = y)P(\tilde{x}[2] = y)) \quad (28)$$

$$= \frac{\mathbb{E}(\tilde{y}^2)}{4} + 0 \quad (29)$$

$$= \frac{\sigma^2}{4}. \quad (30)$$

The covariance matrix equals

$$\Sigma_{\tilde{x}} = \sigma^2 \begin{bmatrix} 1 & \frac{1}{4} \\ \frac{1}{4} & 1 \end{bmatrix}. \quad (31)$$

and its inverse equals

$$\Sigma_{\tilde{x}}^{-1} = \frac{16}{15\sigma^2} \begin{bmatrix} 1 & -\frac{1}{4} \\ -\frac{1}{4} & 1 \end{bmatrix}. \quad (32)$$

The crosscovariance equals

$$\text{Cov}(\tilde{x}[i]\tilde{y}) = \text{E}(\tilde{x}[i]\tilde{y}) \quad (33)$$

$$= \text{E}(\tilde{x}[i]\tilde{y} \mid \tilde{a} = 1)\text{P}(\tilde{x}[i] = y) + \text{E}(\tilde{x}[i]\tilde{y} \mid \tilde{a} = 0)\text{P}(\tilde{x}[i] \neq y) \quad (34)$$

$$= \frac{\sigma^2}{2} + 0 \quad (35)$$

$$= \frac{\sigma^2}{2}. \quad (36)$$

By Theorem 2.3 in the notes on linear regression the linear MMSE estimator equals

$$\hat{y}(\tilde{x}) = \tilde{x}^T \frac{16}{15\sigma^2} \begin{bmatrix} 1 & -\frac{1}{4} \\ -\frac{1}{4} & 1 \end{bmatrix} \begin{bmatrix} \frac{\sigma^2}{2} \\ \frac{\sigma^2}{2} \end{bmatrix} \quad (37)$$

$$= \frac{2(\tilde{x}[1] + \tilde{x}[2])}{5}. \quad (38)$$

b. By Theorem 2.3 in the notes on linear regression the corresponding MSE equals

$$\text{MSE} = \text{Var}(\tilde{y}) - \Sigma_{\tilde{x}\tilde{y}}^T \Sigma_{\tilde{x}}^{-1} \Sigma_{\tilde{x}\tilde{y}} \quad (39)$$

$$= \sigma^2 - \begin{bmatrix} \frac{\sigma^2}{2} \\ \frac{\sigma^2}{2} \end{bmatrix}^T \frac{16}{15\sigma^2} \begin{bmatrix} 1 & -\frac{1}{4} \\ -\frac{1}{4} & 1 \end{bmatrix} \begin{bmatrix} \frac{\sigma^2}{2} \\ \frac{\sigma^2}{2} \end{bmatrix} \quad (40)$$

$$= \sigma^2 - \frac{2\sigma^2}{5} \quad (41)$$

$$= \frac{3\sigma^2}{5}. \quad (42)$$

c. Note that conditioned on  $\tilde{x}[1] \neq \tilde{x}[2]$ , there are only three possible events that are equally likely:  $\{\tilde{x}[1] = \tilde{y}, \tilde{x}[2] = \tilde{z}_2\}$ ,  $\{\tilde{x}[1] = \tilde{z}_1, \tilde{x}[2] = \tilde{y}\}$ , and  $\{\tilde{x}[1] = \tilde{z}_1, \tilde{x}[2] = \tilde{z}_2\}$ . By iterated expectation

$$\text{Var}(\tilde{x}[i] \mid \tilde{x}[1] \neq \tilde{x}[2]) = \text{E}(\tilde{x}[i]^2 \mid \tilde{x}[1] \neq \tilde{x}[2]) \quad (43)$$

$$= \text{E}(\tilde{x}[i]^2 \mid \tilde{x}[1] = \tilde{y}, \tilde{x}[2] = \tilde{z}_2)\text{P}(\tilde{x}[1] = \tilde{y}, \tilde{x}[2] = \tilde{z}_2) \quad (44)$$

$$+ \text{E}(\tilde{x}[i]^2 \mid \tilde{x}[1] = \tilde{z}_1, \tilde{x}[2] = \tilde{y})\text{P}(\tilde{x}[1] = \tilde{z}_1, \tilde{x}[2] = \tilde{y}) \quad (45)$$

$$+ \text{E}(\tilde{x}[i]^2 \mid \tilde{x}[1] = \tilde{z}_1, \tilde{x}[2] = \tilde{z}_2)\text{P}(\tilde{x}[1] = \tilde{z}_1, \tilde{x}[2] = \tilde{z}_2) \quad (46)$$

$$= \text{E}(\tilde{y}^2)\text{P}(\tilde{x}[1] = \tilde{y}, \tilde{x}[2] = \tilde{z}_2) \quad (47)$$

$$+ \text{E}(\tilde{z}_1^2)\text{P}(\tilde{x}[1] = \tilde{z}_1, \tilde{x}[2] = \tilde{y}) \quad (48)$$

$$+ \text{E}(\tilde{z}_1^2)\text{P}(\tilde{x}[1] = \tilde{z}_1, \tilde{x}[2] = \tilde{z}_2) \quad (49)$$

$$= \sigma^2. \quad (50)$$

Again, by iterated expectation

$$\text{Cov}(\tilde{x}[1]\tilde{x}[2] \mid \tilde{x}[1] \neq \tilde{x}[2]) = \text{E}(\tilde{x}[1]\tilde{x}[2] \mid \tilde{x}[1] \neq \tilde{x}[2]) \quad (51)$$

$$= \text{E}(\tilde{y}\tilde{z}_2)\text{P}(\tilde{x}[1] = \tilde{y}, \tilde{x}[2] = \tilde{z}_2) \quad (52)$$

$$+ \text{E}(\tilde{z}_1\tilde{y})\text{P}(\tilde{x}[1] = \tilde{z}_1, \tilde{x}[2] = \tilde{y}) \quad (53)$$

$$+ \text{E}(\tilde{z}_1\tilde{z}_2)\text{P}(\tilde{x}[1] = \tilde{z}_1, \tilde{x}[2] = \tilde{z}_2) \quad (54)$$

$$= 0. \quad (55)$$

The covariance matrix conditioned on  $\tilde{x}[1] \neq \tilde{x}[2]$  equals

$$\Sigma_{\tilde{x} \mid \tilde{x}[1] \neq \tilde{x}[2]} = \sigma^2 I. \quad (56)$$

and its inverse equals

$$\Sigma_{\tilde{x} \mid \tilde{x}[1] \neq \tilde{x}[2]}^{-1} = \frac{1}{\sigma^2} I. \quad (57)$$

The crosscovariance equals

$$\text{Cov}(\tilde{x}[1]\tilde{y} \mid \tilde{x}[1] \neq \tilde{x}[2]) = \text{E}(\tilde{x}[1]\tilde{y} \mid \tilde{x}[1] \neq \tilde{x}[2]) \quad (58)$$

$$= \text{E}(\tilde{y}^2)\text{P}(\tilde{x}[1] = \tilde{y}, \tilde{x}[2] = \tilde{z}_2) \quad (59)$$

$$+ \text{E}(\tilde{z}_1\tilde{y})\text{P}(\tilde{x}[1] = \tilde{z}_1, \tilde{x}[2] = \tilde{y}) \quad (60)$$

$$+ \text{E}(\tilde{z}_1\tilde{y})\text{P}(\tilde{x}[1] = \tilde{z}_1, \tilde{x}[2] = \tilde{z}_2) \quad (61)$$

$$= \frac{\sigma^2}{3}. \quad (62)$$

By Theorem 2.3 in the notes on linear regression the linear MMSE estimator equals

$$\hat{y}(\tilde{x}) = \frac{1}{\sigma^2} \tilde{x}^T \begin{bmatrix} \frac{\sigma^2}{3} \\ \frac{\sigma^2}{3} \end{bmatrix} \quad (63)$$

$$= \frac{\tilde{x}[1] + \tilde{x}[2]}{3}, \quad (64)$$

which is our estimate if  $\tilde{x}[1] \neq \tilde{x}[2]$ . The corresponding MSE conditioned on  $\tilde{x}[1] \neq \tilde{x}[2]$

$$\text{E}(\tilde{y} - \mid \tilde{x}[1] \neq \tilde{x}[2]) = \text{Var}(\tilde{y} \mid \tilde{x}[1] \neq \tilde{x}[2]) - \Sigma_{\tilde{x}\tilde{y} \mid \tilde{x}[1] \neq \tilde{x}[2]}^T \Sigma_{\tilde{x} \mid \tilde{x}[1] \neq \tilde{x}[2]}^{-1} \Sigma_{\tilde{x}\tilde{y} \mid \tilde{x}[1] \neq \tilde{x}[2]} \quad (65)$$

$$= \sigma^2 - \frac{1}{\sigma^2} \begin{bmatrix} \frac{\sigma^2}{3} \\ \frac{\sigma^2}{3} \end{bmatrix}^T \begin{bmatrix} \frac{\sigma^2}{3} \\ \frac{\sigma^2}{3} \end{bmatrix} \quad (66)$$

$$= \sigma^2 - \frac{2\sigma^2}{9} \quad (67)$$

$$= \frac{7\sigma^2}{9}. \quad (68)$$

Finally, the MSE equals

$$\text{E}((\tilde{y} - \hat{y}_{\text{nl}}(\tilde{x}))^2) = \text{E}((\tilde{y} - \hat{y}_{\text{nl}}(\tilde{x}))^2 \mid \tilde{x}[1] = \tilde{x}[2]) \text{P}(\tilde{x}[1] = \tilde{x}[2]) \quad (69)$$

$$+ \text{E}((\tilde{y} - \hat{y}_{\text{nl}}(\tilde{x}))^2 \mid \tilde{x}[1] \neq \tilde{x}[2]) \text{P}(\tilde{x}[1] \neq \tilde{x}[2])$$

$$= 0 + \frac{3}{4} \text{E}(\tilde{y} - \hat{y}_{\text{nl}}(\tilde{x}) \mid \tilde{x}[1] \neq \tilde{x}[2]) \quad (70)$$

$$= \frac{7\sigma^2}{12}. \quad (71)$$

5. *Linear regression with dimensionality reduction.*

- No, let  $\beta^*$  be a solution. There exists at least one nonzero vector  $\beta^{\text{null}}$  in the null space of  $X^T$  because it is low rank. As a result,  $\beta^* + \alpha\beta^{\text{null}}$  is also a solution for any value of  $\alpha \in \mathbb{R}$ .
- By Corollary 7.2 in the notes on PCA, the optimal projection is obtained by computing the eigendecomposition of the covariance matrix  $XX^T$  and using the first  $r$  principal directions. Since  $XX^T = US^2U^T$ , we set  $P := U^T$ .
- The  $\ell_2$  norm is preserved exactly. Since  $X = USV^T$ , we can express each data point as  $x_i = Uw_i$  where  $w_i := (SV^T)_i$  is the  $i$ th column of  $SV^T$ . Then we have

$$\sum_{i=1}^n \|U^T x_i\|_2^2 = \sum_{i=1}^n x_i^T U U^T x_i \quad (72)$$

$$= \sum_{i=1}^n w_i^T w_i \quad (73)$$

$$= \sum_{i=1}^n w_i^T U^T U w_i \quad (74)$$

$$= \sum_{i=1}^n x_i^T x_i \quad (75)$$

$$= \sum_{i=1}^n \|x_i\|_2^2, \quad (76)$$

because  $U^T U = I$ .

- The matrix  $U^T X$  contains the points of reduced dimensionality as its columns. We fit the model by solving

$$\min_{\beta \in \mathbb{R}^r} \|y - X^T U \beta\|_2. \quad (77)$$

$U^T X$  is full rank and has SVD  $SV^T$ , so by Eq. (54) in the lecture notes on linear regression the solution equals

$$\beta_{\text{LS}} = S^{-1} V^T y \quad (78)$$

$$= S^{-1} V^T (X^T \beta_{\text{true}} + z) \quad (79)$$

$$= U^T \beta_{\text{true}} + S^{-1} V^T z. \quad (80)$$

- No. We obtain  $U^T \beta_{\text{true}}$  instead. Since  $U^T$  is a fat matrix, we cannot recover  $\beta_{\text{true}}$  from  $U^T \beta_{\text{true}}$ . However, we can still predict new values of  $y$  by first projecting the features  $x_{\text{test}}$  using  $U^T$  and setting

$$y_{\text{pred}} := \langle U^T x_{\text{test}}, \beta_{\text{LS}} \rangle. \quad (81)$$

6. *Linear regression with orthogonal features.*

- Since  $X$  has orthonormal rows,  $XX^T = I$ , so the OLS estimator is  $Xy$ . Each OLS coefficient is equal to the inner product of a row of  $X$  and the response vector.

b. The ridge-regression estimator equals

$$\beta_{\text{RR}} = (XX^T + \lambda I)^{-1}Xy \quad (82)$$

$$= \frac{Xy}{1 + \lambda} \quad (83)$$

$$= \frac{\beta_{\text{OLS}}}{1 + \lambda}. \quad (84)$$

c. The coefficients equal

$$\beta_{\text{RR}} = \frac{X\tilde{y}}{1 + \lambda} \quad (85)$$

$$= \frac{XX^T\tilde{\beta} + X\tilde{z}}{1 + \lambda} \quad (86)$$

$$= \frac{\tilde{\beta} + X\tilde{z}}{1 + \lambda} \quad (87)$$

$$= \tilde{\beta} - \frac{\lambda\tilde{\beta}}{1 + \lambda} + \frac{X\tilde{z}}{1 + \lambda} \quad (88)$$

which implies

$$\mathbb{E}(\|\tilde{\beta} - \tilde{\beta}_{\text{RR}}\|_2^2) = \mathbb{E}\left(\left(\frac{\lambda\tilde{\beta}}{1 + \lambda} - \frac{X\tilde{z}}{1 + \lambda}\right)^T \left(\frac{\lambda\tilde{\beta}}{1 + \lambda} - \frac{X\tilde{z}}{1 + \lambda}\right)\right) \quad (89)$$

$$= \frac{\lambda^2\mathbb{E}(\|\tilde{\beta}\|_2^2)}{(1 + \lambda)^2} + \frac{\mathbb{E}(\tilde{z}^T X^T X \tilde{z})}{(1 + \lambda)^2}. \quad (90)$$

The cross terms are zero by independence and because both random vectors have zero mean.  $X\tilde{z}$  is a  $p$ -dimensional random vector with covariance matrix  $X\sigma^2 I X^T = \sigma^2 I$ , so  $\mathbb{E}(\tilde{z}^T X^T X \tilde{z})$  is equal to  $\sigma^2 p$  by Theorem 4.5 in the notes on linear regression. We have

$$\mathbb{E}(\|\tilde{\beta} - \tilde{\beta}_{\text{RR}}\|_2^2) = \frac{\lambda^2}{(1 + \lambda)^2} + \frac{\sigma^2 p}{(1 + \lambda)^2} := f(\lambda). \quad (91)$$

The derivative with respect to  $\lambda$  equals

$$f'(\lambda) = \frac{2\lambda}{(1 + \lambda)^3} - \frac{2\sigma^2 p}{(1 + \lambda)^3}, \quad (92)$$

which is zero at  $\lambda = \sigma^2 p$ , negative if  $\lambda$  is smaller and positive if  $\lambda$  is larger. The optimal  $\lambda$  therefore equals  $\sigma^2 p$ . It is proportional to the noise variance and to the number of features.