

Feel free to work with other students, but make sure you write up the homework and code on your own (no copying homework or code; no pair programming). Feel free to ask students or instructors for help debugging code or whatever else, though.

**1 (Murphy 12.5 - Deriving the Residual Error for PCA)** It may be helpful to reference section 12.2.2 of Murphy.

(a) Prove that

$$\left\| \mathbf{x}_i - \sum_{j=1}^k z_{ij} \mathbf{v}_j \right\|^2 = \mathbf{x}_i^\top \mathbf{x}_i - \sum_{j=1}^k \mathbf{v}_j^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v}_j.$$

Hint: first consider the case when  $k = 2$ . Use the fact that  $\mathbf{v}_i^\top \mathbf{v}_j$  is 1 if  $i = j$  and 0 otherwise. Recall that  $z_{ij} = \mathbf{x}_i^\top \mathbf{v}_j$ .

(b) Now show that

$$J_k = \frac{1}{n} \sum_{i=1}^n \left( \mathbf{x}_i^\top \mathbf{x}_i - \sum_{j=1}^k \mathbf{v}_j^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v}_j \right) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^\top \mathbf{x}_i - \sum_{j=1}^k \lambda_j.$$

Hint: recall that  $\mathbf{v}_j^\top \Sigma \mathbf{v}_j = \lambda_j \mathbf{v}_j^\top \mathbf{v}_j = \lambda_j$ .

(c) If  $k = d$  there is no truncation, so  $J_d = 0$ . Use this to show that the error from only using  $k < d$  terms is given by

$$J_k = \sum_{j=k+1}^d \lambda_j.$$

Hint: partition the sum  $\sum_{j=1}^d \lambda_j$  into  $\sum_{j=1}^k \lambda_j$  and  $\sum_{j=k+1}^d \lambda_j$ .

For this problem, we have

- $\mathbf{v}_j \in \mathbb{R}^D$  denotes the  $j$ th principal direction
- $\mathbf{x}_i \in \mathbb{R}^D$  denotes the  $i$ th high-d representation
- $\mathbf{z}_i \in \mathbb{R}^K$  denotes the  $i$ th low-d representation

To begin (a), just expand out the expression

a)  $\|X_i - \sum_{j=1}^K z_{ij} V_j\|^2$  (Recall  $z_{ij} = X_i^T V_j = V_j^T X_i$ )

$$= X_i^T X_i - \underbrace{X_i^T \sum_{j=1}^K z_{ij} V_j}_{\text{These two evaluate to scalars and are the transpose of each other, so are equal}} - \underbrace{\sum_{j=1}^K V_j^T z_{ij} X_i}_{\text{These two evaluate to scalars and are the transpose of each other, so are equal}} + \underbrace{\left[ \sum_{j=1}^K z_{ij} V_j \right]^T \sum_{j=1}^K z_{ij} V_j}_{\left[ z_{ij} V_j \right]^T = V_j^T z_{ij}^T}$$

$$= X_i^T X_i - 2 \sum_{j=1}^K \underbrace{z_{ij}}_{\text{This is a number}} \underbrace{V_j^T X_i}_{\text{So this must be also, take transpose}} + \sum_{j=1}^K \sum_{j=1}^K \underbrace{V_j^T z_{ij}^T z_{ij} V_j}_{= z_{ij}^T V_j^T V_j z_{ij} = z_{ij}^T \delta_{jj} z_{ij} = z_{ij} \delta_{jj} z_{ij}}$$

$$= X_i^T X_i - 2 \sum_{j=1}^K V_j^T X_i X_i^T V_j + \sum_{j=1}^K V_j^T X_i X_i^T V_j$$

$$= \boxed{X_i^T X_i - \sum_{j=1}^K V_j^T X_i X_i^T V_j}$$

← This term is really  $\sum z_{ij}^2$

b) Recall  $\Sigma = \frac{1}{N} \sum_{i=1}^N X_i X_i^T$  is the empirical covariance matrix

$$\rightarrow J_K = \frac{1}{n} \sum_{i=1}^n \left( X_i^T X_i - \sum_{j=1}^K V_j^T X_i X_i^T V_j \right) = \frac{1}{n} \sum_{i=1}^n X_i^T X_i - \sum_{j=1}^K V_j^T \Sigma V_j$$

$$= \boxed{\frac{1}{n} \sum_{i=1}^n X_i^T X_i - \sum_{j=1}^K \lambda_j}$$

←  $V_j$  are eigenvectors of  $\Sigma$  w/  $\lambda = \lambda_j$   
 $\rightarrow V_j^T \Sigma V_j = \lambda_j V_j^T V_j = \lambda_j$  (assuming  $\|V_j\|=1$ )

c)  $J_{K=d} = \frac{1}{n} \sum_{i=1}^n X_i^T X_i - \sum_{j=1}^{K'} \lambda_j - \sum_{j=K'+1}^d \lambda_j = 0$  because we include all dimensions

So  $J_{K=K'} = J_{K=K'} - J_{K=d} = \boxed{\sum_{j=K'+1}^d \lambda_j}$



2 ( $\ell_1$ -Regularization) Consider the  $\ell_1$  norm of a vector  $\mathbf{x} \in \mathbb{R}^n$ :

$$\|\mathbf{x}\|_1 = \sum_i |x_i|.$$

Draw the norm-ball  $B_k = \{\mathbf{x} : \|\mathbf{x}\|_1 \leq k\}$  for  $k = 1$ . On the same graph, draw the Euclidean norm-ball  $A_k = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq k\}$  for  $k = 1$  behind the first plot. (Do not need to write any code, draw the graph by hand).

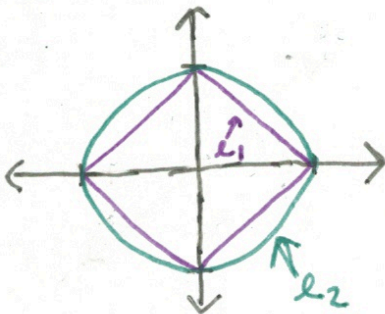
Show that the optimization problem

$$\begin{aligned} &\text{minimize: } f(\mathbf{x}) \\ &\text{subj. to: } \|\mathbf{x}\|_p \leq k \end{aligned}$$

is equivalent to

$$\text{minimize: } f(\mathbf{x}) + \lambda \|\mathbf{x}\|_p$$

(hint: create the Lagrangian). With this knowledge, and the plots given above, argue why using  $\ell_1$  regularization (adding a  $\lambda \|\mathbf{x}\|_1$  term to the objective) will give sparser solutions than using  $\ell_2$  regularization for suitably large  $\lambda$ .



It is clearest to see in two dimensions.

$$\|\vec{x}\|_1 = x + y = 1$$

$$\|\vec{x}\|_2 = \sqrt{x^2 + y^2} = 1$$

This extends into  $n$ -dimensions with  $n$ -dimensional hypercubes and hyperspheres.

Consider representing the minimization problem using Lagrange multipliers:

$$L(\vec{x}, \lambda) = f(\vec{x}) + \lambda (\underbrace{\|\vec{x}\|_p - 1}_{\text{this is zero}})$$

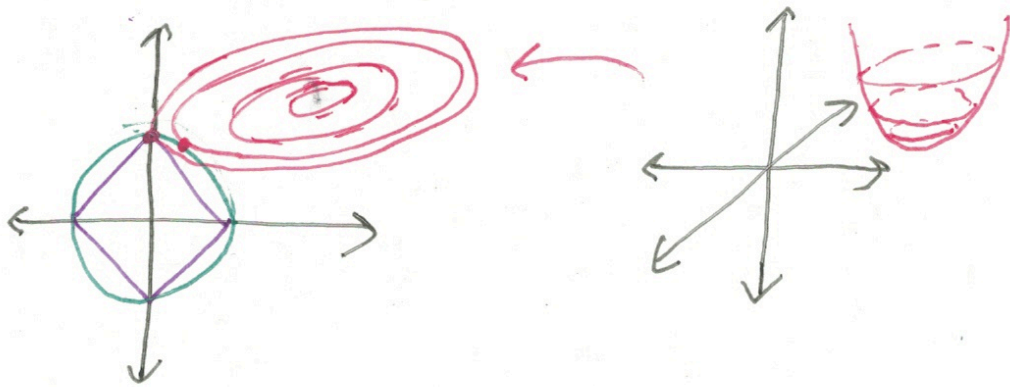
Taking partials of  $L$  yields

$$\frac{\partial L}{\partial x_i} = \frac{\partial f(\vec{x})}{\partial x_i} + \frac{\partial}{\partial x_i} [\lambda \|\vec{x}\|_p]$$

Therefore minimizing the constrained cost function is equivalent to minimizing

$$f(\vec{x}) + \lambda \|\vec{x}\|_p$$

Now looking at the plot,  $L_1$  will give sparser solutions because



We are looking for the minimum value on our convex cost function that intersects with our  $p$ -norm ball.

The  $L_1$  has sharper corners, so we are more likely to hit at the corner vs. at exactly the top of a circle. These corner hits correspond to zero values for one of our fit parameters, i.e. sparser solutions.