

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

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

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

(b) Now show that

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

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

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

a) we have $\left\| \bar{\mathbf{x}}_i - \sum_{j=1}^k z_{ij} \bar{\mathbf{v}}_j \right\|^2$

$$= \left(\bar{\mathbf{x}}_i - \sum_{j=1}^k z_{ij} \bar{\mathbf{v}}_j \right)^T \left(\bar{\mathbf{x}}_i - \sum_{j=1}^k z_{ij} \bar{\mathbf{v}}_j \right)$$

$$= \left(\bar{\mathbf{x}}_i^T - \sum_{j=1}^k z_{ij} \bar{\mathbf{v}}_j^T \right) \left(\bar{\mathbf{x}}_i - \sum_{j=1}^k z_{ij} \bar{\mathbf{v}}_j \right) \text{ multiply out to get:}$$

$$= \bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_i - \bar{\mathbf{x}}_i^T \sum_{j=1}^k z_{ij} \bar{\mathbf{v}}_j - \bar{\mathbf{x}}_i \sum_{j=1}^k z_{ij} \bar{\mathbf{v}}_j^T + \sum_{j=1}^k z_{ij} \bar{\mathbf{v}}_j^T \sum_{j=1}^k z_{ij} \bar{\mathbf{v}}_j$$

$$= \bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_i - \underbrace{\sum_{j=1}^k z_{ij}^2 - \sum_{j=1}^k z_{ij}^2}_{\text{Because } z_{ij} = \bar{\mathbf{x}}_i^T \bar{\mathbf{v}}_j} + \sum_{j=1}^k z_{ij} \bar{\mathbf{v}}_j^T \sum_{j=1}^k z_{ij} \bar{\mathbf{v}}_j$$

$$= \bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_i - 2 \sum_{j=1}^k z_{ij}^2 + \sum_{j=1}^k z_{ij}^2 \text{ Because } \mathbf{v}_i^T \mathbf{v}_j \text{ is 1 if } i=j \text{ \& 0 otherwise}$$

$$= \bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_i - \sum_{j=1}^k z_{ij}^2 = \boxed{\bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_i - \sum_{j=1}^k \bar{\mathbf{v}}_j^T \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \bar{\mathbf{v}}_j} \text{ as desired}$$

b) We have $J_k = \frac{1}{n} \sum_{i=1}^n [\bar{x}_i^T \bar{x}_i - \sum_{j=1}^k (\bar{v}_j^T \bar{x}_i x_i^T \bar{v}_j)]$

$$= \frac{1}{n} \sum_{i=1}^n \bar{x}_i^T \bar{x}_i - \sum_{j=1}^k \left[\bar{v}_j^T \left(\frac{1}{n} \sum_{i=1}^n \Sigma_i \right) \bar{v}_j \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \bar{x}_i^T \bar{x}_i - \sum_{j=1}^k \bar{v}_j^T \Sigma \bar{v}_j$$

↑ avg of covariance matrices

Because $\bar{v}_j^T \Sigma \bar{v}_j = \lambda_j \bar{v}_j^T \bar{v}_j = \lambda_j$

$$= \boxed{\frac{1}{n} \sum_{i=1}^n \bar{x}_i^T \bar{x}_i - \sum_{j=1}^k \lambda_j} \text{ as desired.}$$

c) From part b, we know that

$$J_d = \frac{1}{n} \sum_{i=1}^n \bar{x}_i^T \bar{x}_i - \sum_{j=1}^d \lambda_j = 0$$

$$\underbrace{\frac{1}{n} \sum_{i=1}^n \bar{x}_i^T \bar{x}_i - \sum_{j=1}^k \lambda_j}_{= J_k} - \sum_{j=k+1}^d \lambda_j = 0$$

$$\boxed{J_k = \sum_{j=k+1}^d \lambda_j} \text{ as desired}$$

2 (ℓ_1 -Regularization) Consider the ℓ_1 norm of a vector $x \in \mathbb{R}^n$:

$$\|x\|_1 = \sum_i |x_i|.$$

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

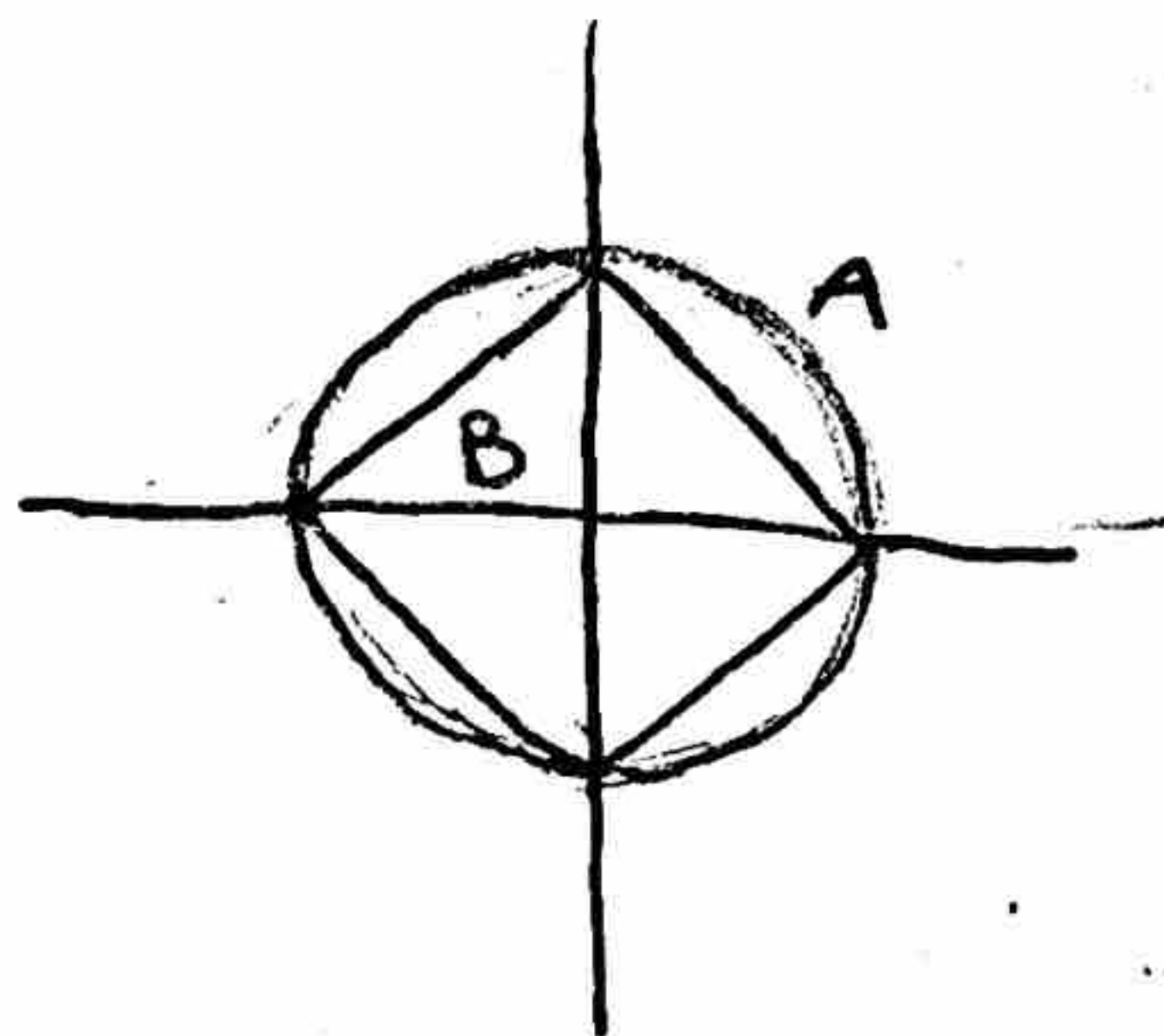
Show that the optimization problem

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

is equivalent to

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

(hint: create the Lagrangian). With this knowledge, and the plots given above, argue why using ℓ_1 regularization (adding a $\lambda \|x\|_1$ term to the objective) will give sparser solutions than using ℓ_2 regularization for suitably large λ .



We have the Lagrangian:

$$L(x, \lambda) = f(x) + \lambda (\|x\|_p - k)$$

$$= f(x) + \lambda \|x\|_p - \lambda k$$

does not depend on x so doesn't matter

So minimizing $f(\vec{x})$ subject to $\|\vec{x}\|_p \leq k$ is equivalent to minimizing $f(\vec{x}) + \lambda \|\vec{x}\|_p$

Checked solution to confirm intuition:

Because ℓ_1 regularization has edges and corners, when the solution is projected onto the surface, it is much more likely^{than a sphere} (infinitely more likely in fact) to land on a corner where one of the parameters will be sent to zero, yielding sparser solutions.