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 $\mathbf{v}_j^\top \mathbf{\Sigma} \mathbf{v}_j = \lambda_j \mathbf{v}_j^\top \mathbf{v}_j = \lambda_j$.

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

■

(a)

$$\left\| \vec{x}_i - \sum_{j=1}^{k} z_{ij} \vec{v}_j \right\|_2^2 = \left( \vec{x}_i - \sum_{j=1}^{k} z_{ij} \vec{v}_j \right)^T \left( \vec{x}_i - \sum_{j=1}^{k} z_{ij} \vec{v}_j \right)$$

$$= \vec{x}_i^T \vec{x}_i - 2 \sum_{j=1}^{k} z_{ij} \vec{v}_j^T \vec{x}_i + \left( \sum_{j=1}^{k} z_{ij} \vec{v}_j \right)^T \left( \sum_{j=1}^{k} z_{ij} \vec{v}_j \right)$$

$$= \vec{x}_i^T \vec{x}_i - 2 \sum_{j=1}^{k} z_{ij} \vec{v}_j^T \vec{x}_i + \sum_{j=1}^{k} \vec{v}_j^T z_{ij} z_{ij} \vec{v}_j$$

$$= \vec{x}_i^T \vec{x}_i - 2 \sum_{j=1}^{k} z_{ij} \vec{v}_j^T \vec{x}_i + \sum_{j=1}^{k} \vec{v}_j^T \vec{x}_i \vec{x}_i^T \vec{v}_j$$

$$= \vec{x}_i^T \vec{x}_i - \sum_{j=1}^{k} \vec{v}_j^T \vec{x}_i \vec{x}_i^T \vec{v}_j \quad , \text{ as desired.}$$

(b)

$$\bar{J}_k = \frac{1}{n} \sum_{i=1}^{n} \left( \vec{x}_i^T \vec{x}_i - \sum_{j=1}^{b} \vec{v}_j^T \vec{x}_i \vec{x}_i^T \vec{v}_j \right)$$

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

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

(c) Because $\vec{v}_d = 0$, $\sum_{j=1}^{d} \lambda_j = \frac{1}{n} \sum_{i=1}^{n} \vec{x}_i^T \vec{x}_i$. So

$$\bar{J}_k = \frac{1}{n} \sum_{i=1}^{n} \vec{x}_i^T \vec{x}_i - \sum_{j=1}^{d} \lambda_j + \sum_{j=k+1}^{d} \lambda_j = \sum_{j=k+1}^{d} \lambda_j.$$

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

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

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).
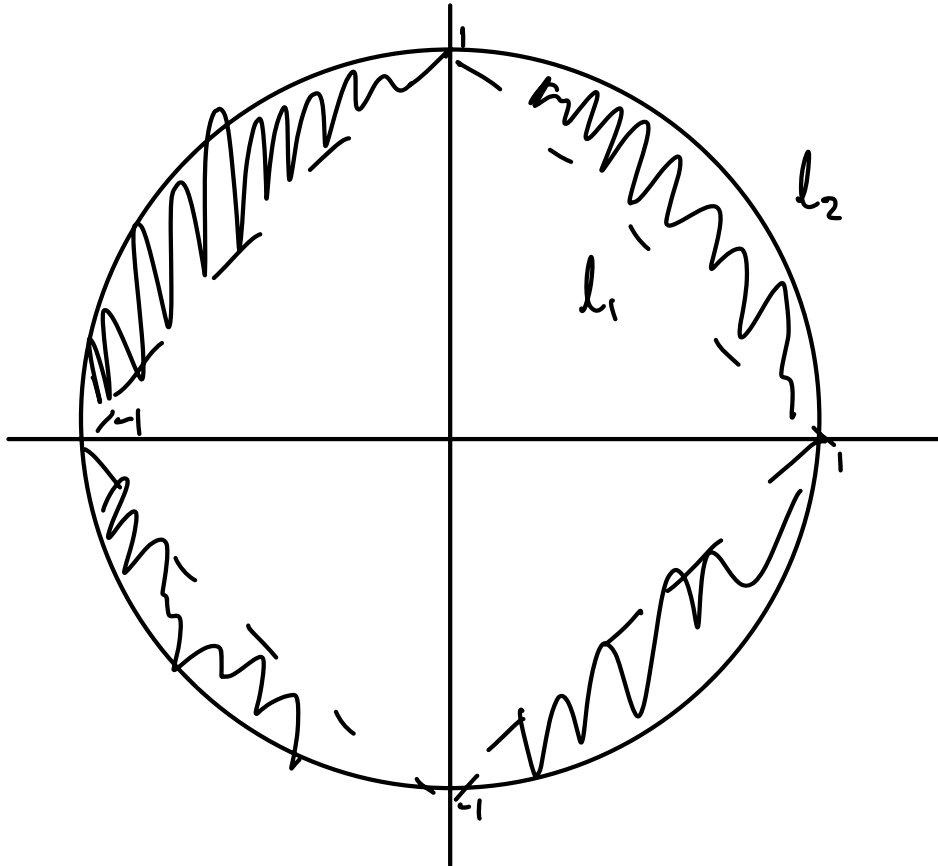
Show that the optimization problem

    minimize: $f(\mathbf{x})$
      subj. to: $\|\mathbf{x}\|_p \leq k$

is equivalent to

    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 $\lambda$.



2

Minimize: $f(\vec{x})$

Subject to: $\|\vec{x}_\rho\| \leq k$.

$$\inf_{x} \sup_{\lambda \geq 0} L(\vec{x}, \lambda) = \inf_{x} \sup_{\lambda \geq 0} f(x) + \lambda \left(\|\vec{x}\|_p - k\right).$$

$$\sup_{\lambda \geq 0} \inf_{\vec{x}} f(\vec{x}) + \lambda \left(\|\vec{x}\|_p - k\right) = \sup_{\lambda \geq 0} g(\lambda).$$

reduces to the problem of

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

It is clear this quantity is minimized when more weights are zero.

We seek to maximize $P(\theta|D) = \dfrac{P(D|\theta)\,P(\theta)}{P(D)}$

$$= \log P(D|\theta) + \log P(\theta) - \log P(D),$$

or minimize $-\log P(D|\theta) - \log P(\theta)$.

Given a prior $\theta_i \sim \mathrm{Lap}(0, b)$.

$$-\log P(\theta) = -\log \prod_i \exp\left(-\frac{|\theta_i|}{b}\right) + C$$

$$= \frac{1}{b}\sum_i |\theta_i| + Z$$

$$= \lambda \|\theta\|_1 + Z.$$

Our original problem is equivalent to

■

Minimizing $-\log P(D|\theta) + \lambda \|\theta\|_1$, as desired.