**1. Deriving the Residual Error for PCA:**

a. Prove that $\|\vec{x}_i - \sum_{j=1}^{k} z_{ij}\vec{v}_j\|^2 = \vec{x}_i^T\vec{x}_i - \sum_{j=1}^{k}\vec{v}_j^T\vec{x}_i\vec{x}_i^T\vec{v}_j$

We prove this algebraically:

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

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

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

Since $\vec{v}_i^T\vec{v}_j = 1$ iff $i=j$: $\vec{x}_i^T\vec{x}_i - 2\sum_{j=1}^{k} z_{ij}\vec{v}_j^T\vec{x}_i + \sum_{j=1}^{k}\vec{v}_j^T\vec{x}_i\vec{x}_i^T v_j$. Since $z_{ij} \in \mathbb{R}$:

$$\rightarrow \vec{x}_i^T\vec{x}_i - 2\sum_{j=1}^{k} z_{ij}\vec{v}_j^T\vec{x}_i + \sum_{j=1}^{k}\vec{v}_j^T\vec{x}_i\vec{x}_i^T\vec{v}_j.\text{ Hence, by simplification:}$$

$$\therefore \boxed{\|\vec{x}_i - \sum_{j=1}^{k} z_{ij}\vec{v}_j\|^2 = \vec{x}_i^T x_i - \sum_{j=1}^{k}\vec{v}_j^T\vec{x}_i\vec{x}_i^T\vec{v}_j}\text{ as desired.}$$

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

Since $\vec{v}_j^T Z \vec{v}_j = \lambda_j\vec{v}_j^T\vec{v}_j = \lambda_j$, we show the following algebraically:

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

$$= \frac{1}{n}\sum_{i=1}^{n}\vec{x}_i^T\vec{x}_i - \sum_{v=1}^{k}\sum_{i=1}^{n}\vec{v}_j \rightarrow \boxed{\frac{1}{n}\sum_{i=1}^{n}\vec{x}_i^T\vec{x}_i - \sum_{j=1}^{k}\lambda_j},\text{ as desired.}$$

c. If $k=d$, there is no truncation, so $J_D = \emptyset$. Use this to show that the error from only using $k < d$ terms is given by the following:

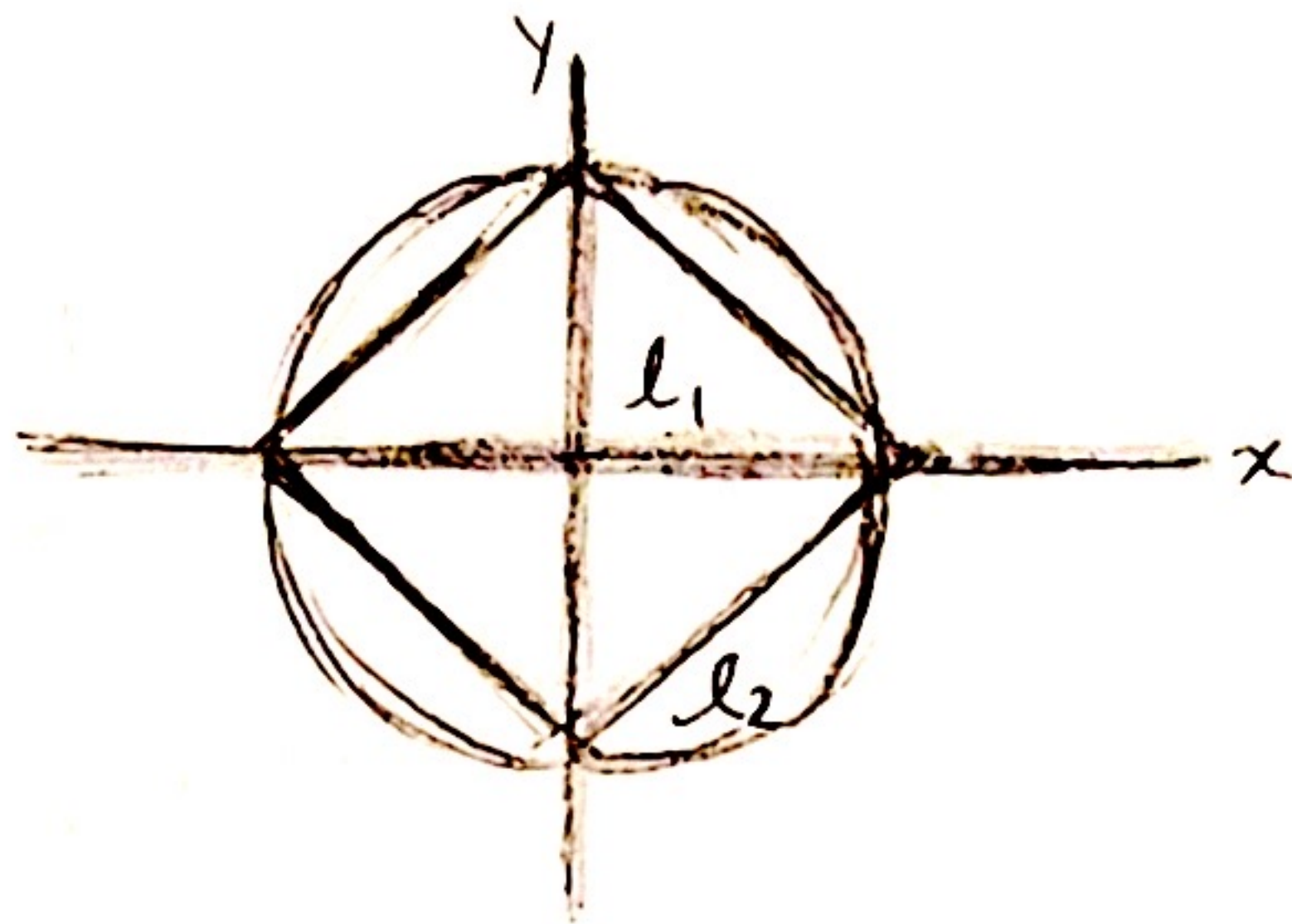$$J_k = \sum_{j=k+1}^{d}\lambda_j \quad (1).$$

Since we can partition the sum $\sum_{j=1}^{d}\lambda_j$ into $\sum_{j=1}^{k}\lambda_j$ and $\sum_{j=k+1}^{d}\lambda_j$, then:

$$J_k = \frac{1}{n}\sum_{i=1}^{n}\vec{x}_i^T\vec{x}_i - \sum_{j=1}^{d}\lambda_j + \sum_{j=k+1}^{d}\lambda_j = \boxed{\sum_{j=k+1}^{d}\lambda_j},\text{ as desired.}$$

## 2. $\ell_1$ - Regularization:

Consider the $\ell_1$ norm of a vector $x \in \mathbb{R}^n$ : $\|\vec{x}\|_1 = \sum |\vec{x}_i|$.

Draw the norm-ball $B_k = \{\vec{x} : \|\vec{x}\|_1 \leq k\}$ for $k=1$. On the same plot, draw the Euclidean norm-ball $A_k = \{\vec{x} : \|\vec{x}\|_2 \leq k\}$ for $k=1$ behind $B_k$.



Now, show that the optimization problem:

$$\text{minimize } f(x) \text{, subject to } \|\vec{x}\|_p \leq k$$

is equivalent to: minimize $f(x) + \lambda \|\vec{x}\|_p$. Then, argue why using $\ell_1$-regularization (adding a $\lambda \|\vec{x}\|_1$ term to the objective) will give sparser solutions than using $\ell_2$-regularization.

We re-write our original problem as $\inf_x \sup_{\lambda \geq 0} L(x, \lambda) = \inf_x \sup_{\lambda \geq 0} f(x) + \lambda(\|\vec{x}\|_p - k)$. In its dual, we can "flip" the infimum and supremum, such that:

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

Since the minimizing value of $f(x) + \lambda(\|\vec{x}\|_p - k)$ over $x$ is equivalent to the minimizing value of $f(x) + \lambda\|\vec{x}\|_p$, and $(-\lambda k)$ does not depend on $x$, we know that the optimizing $x$ will solve "minimize $f(x) + \lambda\|\vec{x}\|_p$" for some value $\lambda \geq 0$. Hence, in tandem with our plot, $\ell_1$-regularization is the projection of our actual optimal solution onto some well-defined $\ell_1$ norm-ball. As our $\ell_1$ ball has sharper edges, the probability of landing on an edge and not on the face [where both elements of the vector are non-zero] is infinitely larger than the $\ell_2$ ball. Specifically, this is due to the rotation invariance of the $\ell_2$, which does not hold for the $\ell_1$ ball. Furthermore, if we were to then generalize this to higher dimensions, the $\ell_1$-penalty would encourage more weights to be zero, compared to the $\ell_2$ ball, as desired.