

## Exercise Sheet 4

### Exercise 1: Lagrange Multipliers (10 + 10 P)

Let  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$  be a dataset of  $N$  data points. We consider the objective function

$$J(\boldsymbol{\theta}) = \sum_{k=1}^N \|\boldsymbol{\theta} - \mathbf{x}_k\|^2$$

to be minimized with respect to the parameter  $\boldsymbol{\theta} \in \mathbb{R}^d$ . In absence of constraints, the parameter  $\boldsymbol{\theta}$  that minimizes this objective is given by the empirical mean  $\mathbf{m} = \frac{1}{N} \sum_{k=1}^N \mathbf{x}_k$ . However, this is generally not the case when the parameter  $\boldsymbol{\theta}$  is constrained.

- Using the method of Lagrange multipliers, *find* the parameter  $\boldsymbol{\theta}$  that minimizes  $J(\boldsymbol{\theta})$  subject to the constraint  $\boldsymbol{\theta}^\top \mathbf{b} = 0$ , with  $\mathbf{b}$  some unit vector in  $\mathbb{R}^d$ . Give a geometrical interpretation to your solution.
- Using the same method, *find* the parameter  $\boldsymbol{\theta}$  that minimizes  $J(\boldsymbol{\theta})$  subject to  $\|\boldsymbol{\theta} - \mathbf{c}\|^2 = 1$ , where  $\mathbf{c}$  is a vector in  $\mathbb{R}^d$  different from  $\mathbf{m}$ . Give a geometrical interpretation to your solution.

### Exercise 2: Principal Component Analysis (10 + 10 P)

We consider a dataset  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ . Principal component analysis searches for a unit vector  $\mathbf{u} \in \mathbb{R}^d$  such that projecting the data on that vector produces a distribution with maximum variance. Such vector can be found by solving the optimization problem:

$$\arg \max_{\mathbf{u}} \frac{1}{N} \sum_{k=1}^N \left[ \mathbf{u}^\top \mathbf{x}_k - \frac{1}{N} \left( \sum_{l=1}^N \mathbf{u}^\top \mathbf{x}_l \right) \right]^2 \quad \text{with} \quad \|\mathbf{u}\|^2 = 1$$

- Show that the problem above can be rewritten as

$$\arg \max_{\mathbf{u}} \mathbf{u}^\top \mathbf{S} \mathbf{u} \quad \text{with} \quad \|\mathbf{u}\|^2 = 1$$

where  $\mathbf{S} = \sum_{k=1}^N (\mathbf{x}_k - \mathbf{m})(\mathbf{x}_k - \mathbf{m})^\top$  is the scatter matrix, and  $\mathbf{m} = \frac{1}{N} \sum_{k=1}^N \mathbf{x}_k$  is the empirical mean.

- Show using the method of Lagrange multipliers that the problem above can be reformulated as solving the eigenvalue problem

$$\mathbf{S} \mathbf{u} = \lambda \mathbf{u}$$

and retaining the eigenvector  $\mathbf{u}$  associated to the highest eigenvalue  $\lambda$ .

### Exercise 3: Bounds on Eigenvalues (5 + 5 + 5 + 5 P)

Let  $\lambda_1$  denote the largest eigenvalue of the matrix  $\mathbf{S}$ . The eigenvalue  $\lambda_1$  quantifies the variance of the data when projected on the first principal component. Because its computation can be expensive, we study how the latter can be bounded with the diagonal elements of the matrix  $\mathbf{S}$ .

- Show that  $\sum_{i=1}^d \mathbf{S}_{ii}$  is an upper bound to the eigenvalue  $\lambda_1$ .
- State the conditions on the data for which the upper bound is tight.
- Show that  $\max_{i=1}^d \mathbf{S}_{ii}$  is a lower bound to the eigenvalue  $\lambda_1$ .
- State the conditions on the data for which the lower bound is tight.

**Exercise 4: Iterative PCA (10 P)**

When performing principal component analysis, computing the full eigendecomposition of the scatter matrix  $\mathbf{S}$  is typically slow, and we are often only interested in the first principal components. An efficient procedure to find the first principal component is *power iteration*. It starts with a random unit vector  $\mathbf{w}^{(0)} \in \mathbb{R}^d$ , and iteratively applies the parameter update

$$\mathbf{w}^{(t+1)} = \mathbf{S}\mathbf{w}^{(t)} / \|\mathbf{S}\mathbf{w}^{(t)}\|$$

until some convergence criterion is met. Here, we would like to show the exponential convergence of power iteration. For this, we look at the error terms

$$\mathcal{E}_k(\mathbf{w}) = \left| \frac{\mathbf{w}^\top \mathbf{u}_k}{\mathbf{w}^\top \mathbf{u}_1} \right| \quad \text{with } k = 2, \dots, d,$$

and observe that they should all converge to zero as  $\mathbf{w}$  approaches the eigenvector  $\mathbf{u}_1$  and becomes orthogonal to other eigenvectors.

- (a) Show that  $\mathcal{E}_k(\mathbf{w}^{(T)}) = |\lambda_k/\lambda_1|^T \cdot \mathcal{E}_k(\mathbf{w}^{(0)})$ , i.e. the convergence of the algorithm is exponential with the number of time steps  $T$ .

**Exercise 5: Programming (30 P)**

Download the programming files on ISIS and follow the instructions.

### Exercise 1: Lagrange Multipliers (10 + 10 P)

Let  $x_1, \dots, x_N \in \mathbb{R}^d$  be a dataset of  $N$  data points. We consider the objective function

$$J(\theta) = \sum_{k=1}^N \|\theta - x_k\|^2$$

to be minimized with respect to the parameter  $\theta \in \mathbb{R}^d$ . In absence of constraints, the parameter  $\theta$  that minimizes this objective is given by the empirical mean  $m = \frac{1}{N} \sum_{k=1}^N x_k$ . However, this is generally not the case when the parameter  $\theta$  is constrained.

- a) Using the method of Lagrange multipliers, find the parameter  $\theta$  that minimizes  $J(\theta)$  subject to the constraint  $\theta^\top b = 0$ , with  $b$  some unit vector in  $\mathbb{R}^d$ . Give a geometrical interpretation to your solution.

a)

$$\mathcal{L}(\theta, \lambda) = J(\theta) + \lambda(\theta^\top b)$$

$$= \sum_{k=1}^N \|\theta - x_k\|^2 + \lambda(\theta^\top b)$$

$$= \sum_{k=1}^N [\theta^2 - 2\theta x_k + x_k^2] + \lambda \cdot \theta^\top b$$

$$= N\theta^2 - 2 \cdot \sum_{k=1}^N \theta x_k + \sum_{k=1}^N x_k^2 + \lambda \theta^\top b$$

$$\nabla_{\theta} \mathcal{L} = 2N\theta - 2 \sum_{k=1}^N x_k + \lambda b \stackrel{!}{=} 0 \quad \Leftrightarrow$$

$$\theta^\top b = 0$$

$$\Leftrightarrow \left(m - \frac{\lambda b}{2N}\right)^\top b = 0$$

$$\Leftrightarrow \left(m^\top - \frac{\lambda}{2N} b^\top\right) b = 0$$

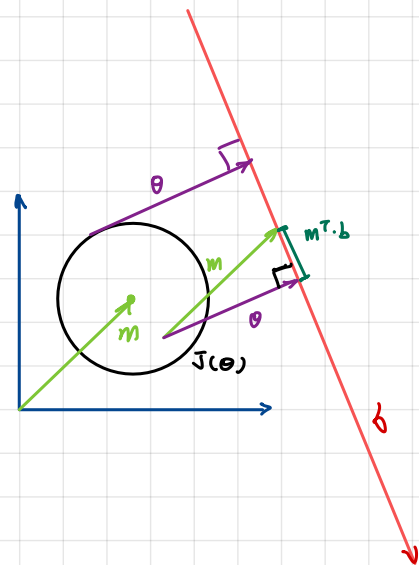
$$\Leftrightarrow \lambda b^\top \cdot b = 2N m^\top \cdot b$$

$$\Leftrightarrow \lambda = \frac{2N \cdot m^\top \cdot b}{b^\top \cdot b}$$

$$\Rightarrow \theta^* = m - \frac{\lambda b}{2N} = m - \frac{m^\top \cdot b \cdot b}{b^\top \cdot b} = m - \frac{m^\top \cdot b \cdot b}{\|b\|^2 \stackrel{!}{=} 1} = m - (b^\top m) b$$

$\theta$  is orthogonal to  $b$  by  $\theta^\top b = 0$ .

Geometrically the optimal  $\theta$  is obtained by subtracting the projection of  $m$  onto the direction of  $b$



The function  $\theta \mapsto \frac{1}{2} \|\theta - m\|^2$  is strictly convex  
 $\rightarrow$  Found solution is global minimum.

$$\begin{aligned} \underset{\theta}{\operatorname{argmin}} \sum_{k=1}^N \|\theta - x_k\|^2 &= \underset{\theta}{\operatorname{argmin}} \sum_k \|\theta - 2\theta^\top x_k \\ &= \underset{\theta}{\operatorname{argmin}} \sum_k (\|\theta\|^2 - 2\theta^\top m) \\ &= \underset{\theta}{\operatorname{argmin}} \sum_k \|\theta - m\|^2 \\ &= \underset{\theta}{\operatorname{argmin}} \frac{1}{2} \|\theta - m\|^2 \end{aligned}$$

$$\begin{aligned} \theta &= \frac{\sum_{k=1}^N x_k - \lambda b}{2N} \\ &= m - \frac{\lambda b}{2N} \\ \frac{b^\top \theta}{\|b\|} &= b^\top m - \frac{\lambda \|b\|}{2N} \Rightarrow \lambda = \frac{2N b^\top m}{\|b\| \stackrel{!}{=} 1} = 2N b^\top m \end{aligned}$$

→ 2

$$\begin{aligned} \mathcal{L}(\theta, \lambda) &= J(\theta) + \lambda (\|\theta - c\|^2 - 1) \\ &= \sum_{k=1}^N \|\theta - x_k\|^2 + \lambda (\|\theta - c\|^2 - 1) \\ &= \sum_{k=1}^N (\theta - x_k)^T (\theta - x_k) + \lambda (\theta - c)^T (\theta - c) - \lambda \\ \nabla_{\theta} \mathcal{L} &= \sum_{k=1}^N 2(\theta - x_k) + 2\lambda(\theta - c) = \vec{0} \\ \Leftrightarrow \quad 2\theta \cdot N - \sum_{k=1}^N 2x_k + 2\lambda\theta - 2\lambda c &= \vec{0} \\ \Leftrightarrow \quad (2N + 2\lambda)\theta &= \sum_{k=1}^N 2x_k + 2\lambda c \\ \Leftrightarrow \quad \theta &= \frac{1}{2N + 2\lambda} \left( \sum_{k=1}^N 2x_k + 2\lambda c \right) \\ &= \frac{1}{N + \lambda} \cdot (N \cdot m + \lambda c) \end{aligned}$$

$$\|\theta - c\|^2 = 1 \quad \Leftrightarrow (\theta - c)^T (\theta - c) = \frac{1}{(N + \lambda)^2} (N \cdot m + \lambda c - cN - \lambda)^T (N \cdot m - cN) = 1$$

$$\Leftrightarrow (Nm - cN)^T (Nm - cN) = (N + \lambda)^2$$

$$\Leftrightarrow N^2 (m - c)^T (m - c) = (N + \lambda)^2$$

$$\Leftrightarrow (N+\lambda)^2 = N^2 \|m-c\|^2$$

$$\Leftrightarrow \lambda = -N \pm \sqrt{N^2 \|m-c\|^2}$$

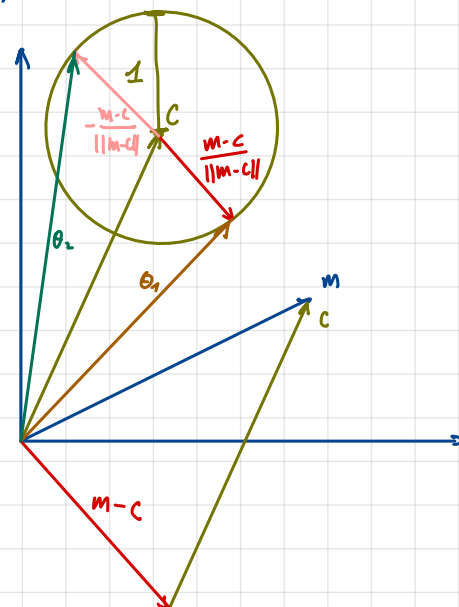
$$\begin{aligned}\theta^* &= \frac{1}{N+\lambda} \cdot (N \cdot m + \lambda c) = \frac{1}{N \cdot \cancel{N} \pm N \cdot \|m - c\|} \cdot (m \cdot N - cN \pm cN \cdot \|m - c\|) \\ &= \frac{N}{\pm N \cdot \|m - c\|} \cdot (m - c \pm c \cdot \|m - c\|) \\ &= \frac{1}{\pm \|m - c\|} (m - c \pm c \|m - c\|) \\ &= \frac{m - c}{\pm \|m - c\|} + c\end{aligned}$$

$\theta$  must lie on a circle of radius 1 centered at  $c$  by  $\|\theta - c\|^2 = 1$ .

The optimal  $\theta$  can be obtained by moving 1 unit from  $c$  along the unit vector  $m - c$  or the opposite direction  $-(m - c)$ .

To find the optimal solution, we observe that:

$$J\left(c + \frac{m-c}{\|m-c\|}\right) \leq J\left(c - \frac{m-c}{\|m-c\|}\right)$$



**Exercise 2: Principal Component Analysis (10 + 10 P)**

We consider a dataset  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ . Principal component analysis searches for a unit vector  $\mathbf{u} \in \mathbb{R}^d$  such that projecting the data on that vector produces a distribution with maximum variance. Such vector can be found by solving the optimization problem:

$$\arg \max_{\mathbf{u}} \frac{1}{N} \sum_{k=1}^N \left[ \mathbf{u}^\top \mathbf{x}_k - \frac{1}{N} \left( \sum_{l=1}^N \mathbf{u}^\top \mathbf{x}_l \right) \right]^2 \quad \text{with} \quad \|\mathbf{u}\|^2 = 1$$

(a) Show that the problem above can be rewritten as

$$\arg \max_{\mathbf{u}} \mathbf{u}^\top \mathbf{S} \mathbf{u} \quad \text{with} \quad \|\mathbf{u}\|^2 = 1$$

where  $\mathbf{S} = \sum_{k=1}^N (\mathbf{x}_k - \mathbf{m})(\mathbf{x}_k - \mathbf{m})^\top$  is the scatter matrix, and  $\mathbf{m} = \frac{1}{N} \sum_{k=1}^N \mathbf{x}_k$  is the empirical mean.

$$\begin{aligned} & \arg \max_{\mathbf{u}} \frac{1}{N} \sum_{k=1}^N \left[ \mathbf{u}^\top \mathbf{x}_k - \frac{1}{N} \left( \sum_{l=1}^N \mathbf{u}^\top \mathbf{x}_l \right) \right]^2 \quad \text{with} \quad \|\mathbf{u}\|^2 = 1 \\ &= \arg \max_{\mathbf{u}} \frac{1}{N} \sum_{k=1}^N \left[ \mathbf{u}^\top \mathbf{x}_k - \mathbf{u}^\top \mathbf{m} \right]^2 \\ &= \arg \max_{\mathbf{u}} \frac{1}{N} \sum_{k=1}^N \left[ \mathbf{u}^\top (\mathbf{x}_k - \mathbf{m}) \right]^2 \\ &= \arg \max_{\mathbf{u}} \frac{1}{N} \sum_{k=1}^N \left[ \mathbf{u}^\top (\mathbf{x}_k - \mathbf{m}) (\mathbf{x}_k - \mathbf{m})^\top \mathbf{u} \right] \\ &= \arg \max_{\mathbf{u}} \frac{1}{N} \left[ \mathbf{u}^\top \sum_{k=1}^N (\mathbf{x}_k - \mathbf{m}) (\mathbf{x}_k - \mathbf{m})^\top \mathbf{u} \right] \\ &= \arg \max_{\mathbf{u}} \frac{1}{N} \left[ \mathbf{u}^\top \mathbf{S} \mathbf{u} \right] \\ &\Leftrightarrow \arg \max_{\mathbf{u}} \mathbf{u}^\top \mathbf{S} \mathbf{u} \quad \text{with} \quad \|\mathbf{u}\|^2 = 1 \end{aligned}$$

(b) Show using the method of Lagrange multipliers that the problem above can be reformulated as solving the eigenvalue problem

$$\mathbf{S} \mathbf{u} = \lambda \mathbf{u}$$

and retaining the eigenvector  $\mathbf{u}$  associated to the highest eigenvalue  $\lambda$ .

$$\mathcal{L}(\mathbf{u}, \lambda) = \mathbf{u}^\top \mathbf{S} \mathbf{u} - \lambda (\|\mathbf{u}\|^2 - 1)$$

$$= \mathbf{u}^\top \mathbf{S} \mathbf{u} - \lambda \mathbf{u}^\top \mathbf{u} + \lambda$$

$$\nabla_{\mathbf{u}} \mathcal{L} = 2\mathbf{S} \mathbf{u} - 2\lambda \mathbf{u} = 0 \quad \Leftrightarrow \mathbf{S} \mathbf{u} = \lambda \mathbf{u} \quad (\mathbf{S} \text{ is symmetric})$$

$$\arg \max_{\mathbf{u}} \mathbf{u}^\top \mathbf{S} \mathbf{u} = \lambda \cdot \mathbf{u}^\top \cdot \mathbf{u} = \lambda \|\mathbf{u}\| = \lambda$$

$\Rightarrow$  To find the maximum of  $\mathbf{u}^\top \mathbf{S} \mathbf{u}$  we should choose the highest eigenvalue  $\lambda$

### Exercise 3: Bounds on Eigenvalues (5 + 5 + 5 + 5 P)

Let  $\lambda_1$  denote the largest eigenvalue of the matrix  $S$ . The eigenvalue  $\lambda_1$  quantifies the variance of the data when projected on the first principal component. Because its computation can be expensive, we study how the latter can be bounded with the diagonal elements of the matrix  $S$ .

(a) Show that  $\sum_{i=1}^d S_{ii}$  is an upper bound to the eigenvalue  $\lambda_1$ .

$$\sum_{i=1}^d S_{ii} = \sum_{i=1}^d \lambda_i = \lambda_1 + \sum_{i=2}^d \lambda_i$$

$$\text{So } \sum_{i=1}^d S_{ii} \geq \lambda_1 \text{ holds if } \sum_{i=2}^d \lambda_i \geq 0$$

To prove  $\sum_{i=2}^d \lambda_i \geq 0$ , we could show the matrix  $S$  is semi-positive definite.

But first, we must show  $S$  is symmetric:

$$\begin{aligned} S &= \sum_{k=1}^N (X_k - m)(X_k - m)^T \\ S^T &= \sum_{k=1}^N [(X_k - m)(X_k - m)^T]^T \\ &= \sum_{k=1}^N [(X_k - m)^T]^T (X_k - m)^T \\ &= \sum_{k=1}^N [(X_k - m)(X_k - m)^T] = S \\ \Rightarrow S \text{ is symmetric} \end{aligned}$$

And  $S$  is semi-positive definite if  $\vec{x}^T S \vec{x} \geq 0$  for all  $\vec{x} \in \mathbb{R}^d$

$$\begin{aligned} \Rightarrow \sum_{k=1}^N \vec{x}^T (X_k - m)(X_k - m)^T \vec{x} &= \sum_{k=1}^N \vec{x}^T (X_k - m) \vec{x}^T (X_k - m)^T \\ &= \sum_{k=1}^N \vec{x}^T \vec{x} (X_k - m)^T (X_k - m) \\ &= \sum_{k=1}^N \|\vec{x}\|^2 \cdot \|X_k - m\|^2 \geq 0 \end{aligned}$$

$\Rightarrow S$  is semi-positive definite

$$\Rightarrow \sum_{i=2}^d \lambda_i \geq 0$$

$$\Rightarrow \sum_{i=1}^d S_{ii} \geq \lambda_1$$

(b) State the conditions on the data for which the upper bound is tight.

- } the upper bound is tight if  $\sum_{i=1}^d S_{ii} = \lambda_1$

$$\text{But } \sum_{i=1}^d S_{ii} = \sum_{i=1}^d \lambda_i = \lambda_1 + \sum_{i=2}^d \lambda_i \geq \lambda_1$$

The equal sign holds only if  $\sum_{i=2}^d \lambda_i = 0$

So the condition is  $\sum_{i=2}^d \lambda_i = 0$

$$\lambda_2 = \dots = \lambda_d = 0$$

$\rightarrow$  data lies on a one-dimensional subspace

(c) Show that  $\max_{i=1}^d S_{ii}$  is a lower bound to the eigenvalue  $\lambda_1$ .

$\lambda_1$  is the largest eigenvalue of  $S$

$$\text{So } \lambda \vec{x} = S \vec{x} \Rightarrow \lambda_1 \vec{x}^T \vec{x} = \vec{x}^T S \vec{x}$$

$$\begin{aligned} \lambda_1 &= \max_{u: \|u\|=1} u^T S u > \max_{u \in \{e_1, \dots, e_d\}} u^T S u \\ &= \max_{i=1}^d e_i^T S e_i \\ &= \max_{i=1}^d S_{ii} \end{aligned}$$

Now we assume  $\|\vec{x}\| = \vec{x}^T \cdot \vec{x} = 1$

$$\Rightarrow \lambda_1 = \max \vec{x}^T S \vec{x}$$

We choose a unit vector  $\vec{e}_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \end{pmatrix}$

$$\Rightarrow \lambda_1 = \max \vec{x}^T S \vec{x} \geq \max \vec{e}_i^T S \vec{e}_i = \max S_{ii}$$

$$\Rightarrow \lambda_1 \geq \max S_{ii}$$

(d) State the conditions on the data for which the lower bound is tight.

-5

the lower bound is tight if  $\lambda_1 = \max S_{ii}$

$$\sum_{i=1}^d S_{ii} = \sum_{i=1}^d \lambda_i = \lambda_1 + \sum_{i=2}^d \lambda_i = S_{11} + \sum_{i=2}^d S_{ii} \quad (\text{assume } \max S_{ii} = S_{11})$$

$$\text{if } \lambda_1 = \max S_{ii} = S_{11} \Rightarrow \sum_{i=2}^d \lambda_i = \sum_{i=2}^d S_{ii}$$

$$u_1 \in \{e_1, \dots, e_d\}$$

where  $e_1, \dots, e_d$  are the canonical unit vector

#### Exercise 4: Iterative PCA (10 P)

When performing principal component analysis, computing the full eigendecomposition of the scatter matrix  $S$  is typically slow, and we are often only interested in the first principal components. An efficient procedure to find the first principal component is *power iteration*. It starts with a random unit vector  $w^{(0)} \in \mathbb{R}^d$ , and iteratively applies the parameter update

$$w^{(t+1)} = Sw^{(t)} / \|Sw^{(t)}\|$$

until some convergence criterion is met. Here, we would like to show the exponential convergence of power iteration. For this, we look at the error terms

$$\mathcal{E}_k(w) = \left| \frac{w^\top u_k}{w^\top u_1} \right| \quad \text{with } k = 2, \dots, d,$$

and observe that they should all converge to zero as  $w$  approaches the eigenvector  $u_1$  and becomes orthogonal to other eigenvectors.

- (a) Show that  $\mathcal{E}_k(w^{(T)}) = |\lambda_k/\lambda_1|^T \cdot \mathcal{E}_k(w^{(0)})$ , i.e. the convergence of the algorithm is exponential with the number of time steps  $T$ .

from lecture 23

$$\begin{aligned} Sw^{(t)} &= S \left( \sum_{i=1}^d \langle w^{(t)}, u_i \rangle u_i \right) = \sum_{i=1}^d \langle w^{(t)}, u_i \rangle Su_i \\ &= \sum_{i=1}^d \langle w^{(t)}, u_i \rangle \lambda_i u_i \end{aligned}$$

$$w^{(t+1)} = \frac{Sw^{(t)}}{\|Sw^{(t)}\|} = \frac{\sum_{i=1}^d \langle w^{(t)}, u_i \rangle \lambda_i u_i}{\left\| \sum_{i=1}^d \langle w^{(t)}, u_i \rangle \lambda_i u_i \right\|}$$

$$w^{(t+1)\top} u_1 = \frac{\langle w^{(t)}, u_1 \rangle \lambda_1 u_1^\top u_1}{\|Sw^{(t)}\|} = \frac{\langle w^{(t)}, u_1 \rangle \lambda_1}{\|Sw^{(t)}\|}$$

$$w^{(t+1)\top} u_k = \frac{\langle w^{(t)}, u_k \rangle \lambda_k}{\|Sw^{(t)}\|} u_k^\top u_k = \frac{\langle w^{(t)}, u_k \rangle \lambda_k}{\|Sw^{(t)}\|}$$

$$\begin{aligned} \mathcal{E}_k(w^{(t+1)}) &= \left| \frac{w^{(t+1)\top} u_k}{w^{(t+1)\top} u_1} \right| = \frac{\langle w^{(t)}, u_k \rangle \lambda_k}{\langle w^{(t)}, u_1 \rangle \lambda_1} \\ &= \frac{\lambda_k}{\lambda_1} \cdot \left| \frac{w^{(t)\top} u_k}{w^{(t)\top} u_1} \right| \\ &= \frac{\lambda_k}{\lambda_1} \cdot \mathcal{E}_k(w^{(t)}) \\ &= \left( \frac{\lambda_k}{\lambda_1} \right)^t \mathcal{E}_k(w^{(0)}) \\ &= \left( \frac{\lambda_k}{\lambda_1} \right)^{1+T} \mathcal{E}_k(w^{(0)}) \end{aligned}$$

$$\Rightarrow \mathcal{E}_k(w^{(T)}) = \left( \frac{\lambda_k}{\lambda_1} \right)^T \mathcal{E}_k(w^{(0)})$$