

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.

- (a) 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.

$$\begin{aligned} \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta}) &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \sum_{k=1}^N \|\boldsymbol{\theta}\|^2 - 2\boldsymbol{\theta}^\top \mathbf{x}_k + \|\mathbf{x}_k\|^2 = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} N \cdot \|\boldsymbol{\theta}\|^2 - 2\boldsymbol{\theta}^\top \mathbf{m} \cdot N \\ &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \|\boldsymbol{\theta}\|^2 - 2\boldsymbol{\theta}^\top \mathbf{m} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \frac{1}{2} [\|\boldsymbol{\theta}\|^2 - 2\boldsymbol{\theta}^\top \mathbf{m} + \|\mathbf{m}\|^2] \\ &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \frac{1}{2} [\boldsymbol{\theta} - \mathbf{m}]^2 \end{aligned}$$

$$\mathcal{L}(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda \cdot g(\mathbf{x}) = \frac{1}{2} [\boldsymbol{\theta} - \mathbf{m}]^2 + \lambda \cdot \boldsymbol{\theta}^\top \mathbf{b}$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \mathbf{m} + \lambda \cdot \mathbf{b} = 0 \quad \Rightarrow \quad \boldsymbol{\theta} \mathbf{b}^\top - \mathbf{b}^\top \mathbf{m} + \lambda \mathbf{b}^\top \mathbf{b} = 0 \quad \Rightarrow \quad \lambda = \mathbf{b}^\top \mathbf{m}$$

$$\Rightarrow \boldsymbol{\theta} - \mathbf{m} + \mathbf{b}^\top \mathbf{m} = 0$$

$$\Rightarrow \boldsymbol{\theta} = \mathbf{m} - \mathbf{b}^\top \mathbf{m} \mathbf{b}$$

① 先得得到 λ 和别的参数的关系

② 把 λ 代回得到 $\boldsymbol{\theta}$ 的结果

- (b) 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.

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} [\boldsymbol{\theta} - \mathbf{m}]^2 + \frac{1}{2} \lambda \cdot (\|\boldsymbol{\theta} - \mathbf{c}\|^2 - 1)$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \mathbf{m} + (\boldsymbol{\theta} - \mathbf{c}) \cdot \lambda = 0 \quad \Rightarrow \quad (1 + \lambda)(\boldsymbol{\theta} - \mathbf{c}) = \mathbf{m} - \mathbf{c}$$

$$\Rightarrow (1 + \lambda)^2 \|\boldsymbol{\theta} - \mathbf{c}\|^2 = \|\mathbf{m} - \mathbf{c}\|^2$$

$$\Rightarrow \|1 + \lambda\| = \pm \|\mathbf{m} - \mathbf{c}\|$$

$$\Rightarrow (\boldsymbol{\theta} - \mathbf{c}) = \pm \frac{\mathbf{m} - \mathbf{c}}{\|\mathbf{m} - \mathbf{c}\|}$$

$$\Rightarrow \boldsymbol{\theta} = \mathbf{c} \pm \frac{\mathbf{m} - \mathbf{c}}{\|\mathbf{m} - \mathbf{c}\|}$$

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

We consider a dataset $x_1, \dots, x_N \in \mathbb{R}^d$. Principal component analysis searches for a unit vector $u \in \mathbb{R}^d$ such that projecting the data on that vector produces a distribution with maximum variance. Such vector can be found by solving the optimization problem:

$$\arg \max_u \frac{1}{N} \sum_{k=1}^N \left[u^\top x_k - \frac{1}{N} \left(\sum_{l=1}^N u^\top x_l \right) \right]^2 \quad \text{with} \quad \|u\|^2 = 1$$

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

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

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

$$\begin{aligned} & \arg \max_u \frac{1}{N} \sum_{k=1}^N \left[u^\top x_k - \frac{1}{N} \left(\sum_{l=1}^N u^\top x_l \right) \right]^2 \\ &= \arg \max_u \frac{1}{N} \sum_{k=1}^N \left[u^\top \left(x_k - \frac{1}{N} \sum_{l=1}^N x_l \right) \right]^2 \\ &= \arg \max_u \frac{1}{N} \sum_{k=1}^N \left[u^\top (x_k - m) \right]^2 \\ &= \arg \max_u u^\top \sum_{k=1}^N (x_k - m)^2 u \\ &= \arg \max_u u^\top \sum_{k=1}^N (x_k - m)(x_k - m)^\top u \end{aligned}$$

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

$$S u = \lambda u$$

and retaining the eigenvector u associated to the highest eigenvalue λ .

$$\begin{aligned} J(u) &= u^\top \sum_{k=1}^N (x_k - m)(x_k - m)^\top u = u^\top S u \\ g(u) &= \|u\|^2 - 1 \end{aligned}$$

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

$$\nabla_u \mathcal{L}(u) = 2S u - 2\lambda u = 0 \quad \Rightarrow \quad S u = \lambda u$$

$$u^\top S u = u^\top \lambda u = u^\top u \lambda = \lambda$$

$$\max u^\top S u = \max \lambda$$

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

Let λ_1 denote the largest eigenvalue of the matrix \mathbf{S} . The eigenvalue λ_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} .

- (a) Show that $\sum_{i=1}^d S_{ii}$ is an upper bound to the eigenvalue λ_1 .

$$\sum_{i=1}^d S_{ii} = \text{Tr}(\mathbf{S}) = \sum_{i=1}^d \lambda_i \geq \lambda_1$$

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

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

$$\Rightarrow \text{Var}(u_i^T \mathbf{x}) = 0 \quad \forall i \in \{2, \dots, d\}$$

- (c) Show that $\max_{i=1}^d S_{ii}$ is a lower bound to the eigenvalue λ_1 .

$$\lambda_1 = u_1^T \mathbf{S} u_1 = \max_u u^T \mathbf{S} u \geq \max_{i=1}^d e_i^T \mathbf{S} e_i = \max_{i=1}^d S_{ii}$$

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

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

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 . (*Hint: to show this, it is useful to rewrite the scatter matrix in terms of eigenvalues and eigenvectors, i.e. $\mathbf{S} = \sum_{i=1}^d \mathbf{u}_i \mathbf{u}_i^\top \lambda_i$.*)

$$\begin{aligned} \mathcal{E}_k(\mathbf{w}^{(T)}) &= \left| \frac{\mathbf{w}^{(T)\top} \mathbf{u}_k}{\mathbf{w}^{(T)\top} \mathbf{u}_1} \right| = \left| \frac{\mathbf{S} \mathbf{w}^{(T-1)\top} \cdot \mathbf{u}_k}{\mathbf{S} \mathbf{w}^{(T-1)\top} \cdot \mathbf{u}_1} \right| = \left| \frac{\mathbf{w}^{(T-1)\top} \cdot \sum_{i=1}^d \mathbf{u}_i \mathbf{u}_i^\top \lambda_i \cdot \mathbf{u}_k}{\mathbf{w}^{(T-1)\top} \cdot \sum_{i=1}^d \mathbf{u}_i \mathbf{u}_i^\top \lambda_i \cdot \mathbf{u}_1} \right| \\ &= \left| \frac{\mathbf{w}^{(T-1)\top} \sum_{i=1}^d \mathbf{u}_i \delta_{i,k} \lambda_i}{\mathbf{w}^{(T-1)\top} \sum_{i=1}^d \mathbf{u}_i \delta_{i,1} \lambda_i} \right| \\ &= \left| \frac{\mathbf{w}^{(T-1)\top} \cdot \mathbf{u}_k \cdot \lambda_k}{\mathbf{w}^{(T-1)\top} \cdot \mathbf{u}_1 \cdot \lambda_1} \right| \\ &= \mathcal{E}_k(\mathbf{w}^{(T-1)}) \cdot \left| \frac{\lambda_k}{\lambda_1} \right| \end{aligned}$$

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