

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.

- (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.
- (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.

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.

- (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 λ .

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 $\sum_{i=1}^d \mathbf{S}_{ii}$ is an upper bound to the eigenvalue λ_1 .
- (b) *State* the conditions on the data for which the upper bound is tight.
- (c) *Show* that $\max_{i=1}^d \mathbf{S}_{ii}$ is a lower bound to the eigenvalue λ_1 .
- (d) *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 $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ 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.

Solution:

First the objective function can be written as

$$\begin{aligned} J(\boldsymbol{\theta}) &= \sum_{k=1}^N \|\boldsymbol{\theta} - \mathbf{x}_k\|^2 \\ &= \sum_{k=1}^N \boldsymbol{\theta}^\top \boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \mathbf{x}_k + \mathbf{x}_k^\top \mathbf{x}_k \end{aligned}$$

not related to \mathbf{x}_k

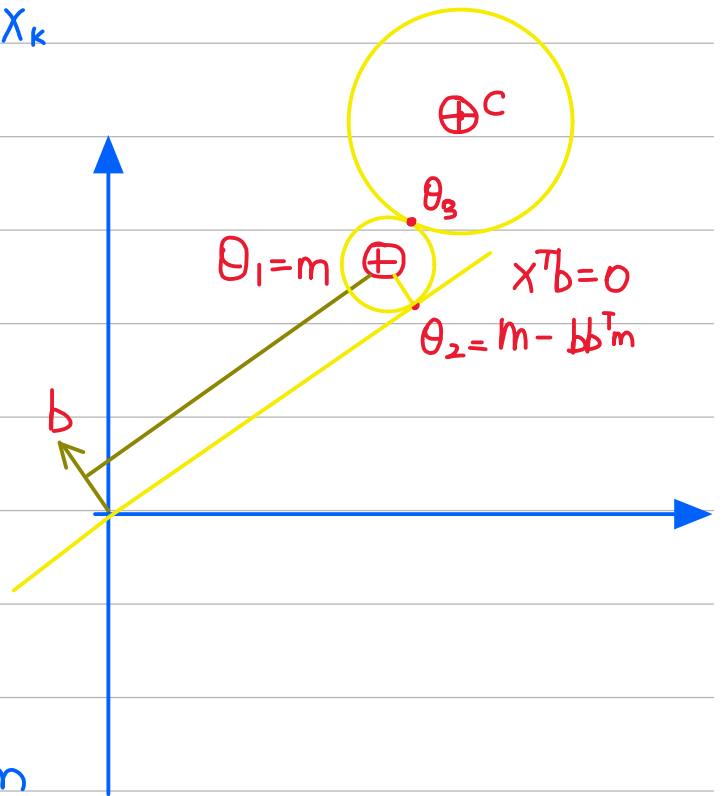
$$\underset{\boldsymbol{\theta}}{\operatorname{argmin}} \sum_{k=1}^N \boldsymbol{\theta}^\top \boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \mathbf{x}_k$$

$$\underset{\boldsymbol{\theta}}{\operatorname{argmin}} N \boldsymbol{\theta}^\top \boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \sum_{k=1}^N \mathbf{x}_k$$

$$\underset{\boldsymbol{\theta}}{\operatorname{argmin}} \boldsymbol{\theta}^\top \boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \mathbf{m}$$

$$\underset{\boldsymbol{\theta}}{\operatorname{argmin}} \boldsymbol{\theta}^\top \boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \mathbf{m} + \mathbf{m}^\top \mathbf{m}$$

$$\underset{\boldsymbol{\theta}}{\operatorname{argmin}} \|\boldsymbol{\theta} - \mathbf{m}\|^2 \quad \boldsymbol{\theta}^* = \mathbf{m}$$



So when this objective is constrained by $\boldsymbol{\theta}^\top \mathbf{b} = 0$ using Lagrange multiplier we can form a Lagrange function:

日期: /

$$\max_{\lambda} \min_{\theta} \mathcal{L}(\theta, \lambda) = \frac{1}{2} \|\theta - m\|^2 + \lambda \theta^T b$$

$$\left\{ \begin{array}{l} \frac{\partial \mathcal{L}}{\partial \theta^T} = (\theta - m) + \lambda b = 0 \\ \frac{\partial \mathcal{L}}{\partial \lambda} = \theta^T b = 0 \end{array} \right.$$



$$\theta - m + \lambda b = 0 \quad \text{①}$$

$$\therefore \theta^* = m - \lambda b$$

Then we need to find the expression of λ .

let ① be left multiplied by b^T .

$$b^T \theta - b^T m + \lambda b^T b = 0$$

Since $b^T \theta = 0$ and b is a unit vector, i.e. $b^T b = 1$

$$\therefore \lambda = b^T m$$

$$\therefore \theta_2^* = m - b^T m b$$

$$= m - b b^T m$$

(b) Using the same method, find the parameter θ that minimizes $J(\theta)$ subject to $\|\theta - c\|^2 = 1$, where c is a vector in \mathbb{R}^d different from m . Give a geometrical interpretation to your solution.

Solution:

The Lagrange function is

$$\mathcal{L}(\theta, \lambda) = \frac{1}{2} \|\theta - m\|^2 + \frac{1}{2} \lambda (\|\theta - c\|^2 - 1)$$

日期: /

$$\frac{\partial L}{\partial \theta} = (\theta - m) + \lambda(\theta - c) = 0$$

$$\therefore (\theta - c) + \lambda(\theta - c) = m - c$$

$$(1 + \lambda)(\theta - c) = m - c \longrightarrow \theta = c + \frac{m - c}{1 + \lambda}$$

$$\therefore \quad \Downarrow \text{ take 2nd Norm}$$

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

$$\Downarrow$$
$$(1 + \lambda)^2 = \|m - c\|^2$$

$$\therefore \lambda = -1 \pm \|m - c\|$$

$$\therefore \theta^* = c \pm \frac{m - c}{\|m - c\|}$$

We can compare the objective value:

$$J\left(c + \frac{m - c}{\|m - c\|}\right) = \left\| c + \frac{m - c}{\|m - c\|} - m \right\|^2 \longrightarrow \left\| -(m - c) + \frac{m - c}{\|m - c\|} \right\|^2$$

$$J\left(c - \frac{m - c}{\|m - c\|}\right) = \left\| c - \frac{m - c}{\|m - c\|} - m \right\|^2 \longrightarrow \left\| -(m - c) - \frac{m - c}{\|m - c\|} \right\|^2$$

$$\left\| (m - c) \times \left(\frac{1}{\|m - c\|} - 1 \right) \right\|^2 \leq \left\| (m - c) \times \left(-\frac{1}{\|m - c\|} - 1 \right) \right\|^2$$

$$\Downarrow$$
$$\theta_3^* = c + \frac{m - c}{\|m - c\|}$$

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.

Solution:

We can derive the function as.

$$\begin{aligned} & \frac{1}{N} \sum_{k=1}^N \left[\mathbf{u}^\top \mathbf{x}_k - \frac{1}{N} \left(\sum_{l=1}^N \mathbf{u}^\top \mathbf{x}_l \right) \right]^2 \\ &= \frac{1}{N} \sum_{k=1}^N \left[\mathbf{u}^\top \mathbf{x}_k - \mathbf{u}^\top \mathbf{m} \right]^2 \\ &= \frac{1}{N} \sum_{k=1}^N (\mathbf{u}^\top \mathbf{x}_k - \mathbf{u}^\top \mathbf{m})^\top (\mathbf{u}^\top \mathbf{x}_k - \mathbf{u}^\top \mathbf{m}) \\ &= \frac{1}{N} \sum_{k=1}^N \left(\mathbf{u}^\top (\mathbf{x}_k - \mathbf{m}) \right)^\top \left(\mathbf{u}^\top (\mathbf{x}_k - \mathbf{m}) \right) \\ &= \frac{1}{N} \sum_{k=1}^N (\mathbf{x}_k - \mathbf{m})^\top \mathbf{u} \mathbf{u}^\top (\mathbf{x}_k - \mathbf{m}) \\ &= \frac{1}{N} \sum_{k=1}^N \mathbf{u}^\top (\mathbf{x}_k - \mathbf{m}) (\mathbf{x}_k - \mathbf{m})^\top \mathbf{u} \end{aligned}$$

↓ delete $\frac{1}{N}$

$$\mathbf{u}^\top \left(\sum_{k=1}^N (\mathbf{x}_k - \mathbf{m}) (\mathbf{x}_k - \mathbf{m})^\top \right) \mathbf{u}$$

↓

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

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

$$Su = \lambda u$$

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

Solution:

We can form the Lagrange function as:

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

$$\frac{\partial \mathcal{L}}{\partial u^T} = Su - \lambda u = 0$$

$$\therefore Su = \lambda u \quad \textcircled{1}$$

So we have proofed that it can be solved by solving an Eigenvalue problem.

Then we left multiply u^T on both sides.

$$u^T S u = \lambda u^T u = \lambda \|u\|^2 = \lambda$$

$$\therefore \arg\max_u u^T S u \iff \arg\max_{\lambda} \lambda$$

\therefore the optimal results is the eigenvector u^* that correspond to the largest eigenvalue λ_1 .

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

Let λ_1 denote the largest eigenvalue of the matrix 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 S .

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

Solution:

$$\lambda_i = u_i^T S u_i = \sum_{k=1}^N (u_i^T (x_k - m))^2 \geq 0$$

$$\sum_i S_{ii} = \text{trace}(S)$$

From Matrix cookbook, Page 6, formular (12)

we have:

$$\sum_i S_{ii} = \text{trace}(S) = \sum_i \lambda_i \geq \lambda_1$$

\therefore the largest eigenvalue is upperbounded by $\sum_i S_{ii}$

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

Solution:

When all other eigenvalues $\lambda_2, \lambda_3 \dots \lambda_d$ are 0.

Then

$$\sum_i S_{ii} = \lambda_1$$

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

Solution:

↙ every element ≥ 0

$$\lambda = \max_{\|u\|=1} u^T S u \geq \max_{u \in \{e_1, \dots, e_d\}} u^T S u = \max_{i=1 \dots d} e_i^T S e_i$$

$$= \{e_1^T S e_1, e_2^T S e_2 \dots e_d^T S e_d\}$$

$$\therefore \lambda \geq \max_{i=1}^d S_{ii}$$

[e_1, \dots, e_d are canonical coordinate vectors.]

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

Solution:

If we want $\lambda = \max_{i=1}^d S_{ii} = \max_{i=1}^d e_i^T S e_i$

then $u \in \{e_1, \dots, e_d\}$

So that it can pick out the largest diagonal element

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 .

Solution:

$$\begin{aligned} \mathcal{E}_k(\mathbf{w}^{(T+1)}) &= \left| \frac{\mathbf{w}^{(T+1)\top} \mathbf{u}_k}{\mathbf{w}^{(T+1)\top} \mathbf{u}_1} \right| \\ &= \left| \frac{(\mathbf{S}\mathbf{w}^{(T)})^\top \mathbf{u}_k}{\|\mathbf{S}\mathbf{w}^{(T)}\|} \right| = \left| \frac{(\mathbf{S}\mathbf{w}^{(T)})^\top \mathbf{u}_k}{(\mathbf{S}\mathbf{w}^{(T)})^\top \mathbf{u}_1} \right| \end{aligned}$$

Since $\mathbf{S} = \mathbf{S}^\top$

$$= \left| \frac{\mathbf{w}^{(T)\top} \mathbf{S} \mathbf{u}_k}{\mathbf{w}^{(T)\top} \mathbf{S} \mathbf{u}_1} \right| = \left| \frac{\mathbf{w}^{(T)\top} \lambda_k \mathbf{u}_k}{\mathbf{w}^{(T)\top} \lambda_1 \mathbf{u}_1} \right|$$

$$= \left| \frac{\mathbf{w}^{(T)\top} \mathbf{u}_k}{\mathbf{w}^{(T)\top} \mathbf{u}_1} \right| \cdot \left| \frac{\lambda_k}{\lambda_1} \right|$$

$$= \mathcal{E}_k(\mathbf{w}^{(T)}) \left| \frac{\lambda_k}{\lambda_1} \right| = \mathcal{E}_k(\mathbf{w}^{(T-1)}) \left| \frac{\lambda_k}{\lambda_1} \right|^2$$

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