

Exercise Sheet 3 (theory part)

Exercise 1: Principal Component Analysis (15 + 15 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:

$$\arg \max_{\mathbf{u}} \frac{1}{N} \sum_{i=1}^N [\mathbf{u}^\top (\mathbf{x}_i - \mathbf{m})]^2 \quad \text{with} \quad \|\mathbf{u}\|^2 = 1 \quad \text{and} \quad \mathbf{m} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i.$$

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

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

where $\Sigma = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mathbf{m})(\mathbf{x}_i - \mathbf{m})^\top$ is the covariance matrix.

$$\begin{aligned} & \frac{1}{N} \sum_{i=1}^N [\mathbf{u}^\top (\mathbf{x}_i - \mathbf{m})]^2 \\ &= \frac{1}{N} \sum_{i=1}^N \mathbf{u}^\top (\mathbf{x}_i - \mathbf{m})(\mathbf{x}_i - \mathbf{m})^\top \mathbf{u} \\ &= \mathbf{u}^\top \left(\underbrace{\frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mathbf{m})(\mathbf{x}_i - \mathbf{m})^\top}_{\Sigma} \right) \mathbf{u} \end{aligned}$$

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

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

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

We apply the method of Lagrange multipliers:

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

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}} = 2\Sigma \mathbf{u} - 2\lambda \mathbf{u} \stackrel{\text{def}}{=} \mathbf{0} \quad \Rightarrow \quad \Sigma \mathbf{u} = \lambda \mathbf{u}$$

Hence, the solution \mathbf{u} is an eigenvector of Σ . To find which one it is, we multiply on both sides by \mathbf{u}^\top :

$$\mathbf{u}^\top \Sigma \mathbf{u} = \underbrace{\mathbf{u}^\top (\lambda \mathbf{u})}_{\lambda}$$

The left hand side is the objective to be optimized and the right hand side is the eigenvalue. Maximizing the objective is therefore achieved by choosing the eigenvector with maximum associated eigenvalue.

Exercise 2: Bounds on Eigenvalues (10 + 10 P)

Let λ_1 denote the largest eigenvalue of the matrix Σ . The eigenvalue λ_1 measures the variance of the data when projected on the first principal component. We study how the latter can be bounded with the diagonal elements of the matrix Σ .

Preliminary observation for all eigenvalues/eigenvectors $(\lambda_i, \mathbf{u}_i)$:

$$\lambda_i = \mathbf{u}_i^\top \Sigma \mathbf{u}_i = \sum_k (\mathbf{u}_i^\top (\mathbf{x}_k - \mathbf{m}))^2 \geq 0$$

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

$$\sum_{i=1}^d \Sigma_{ii} = \sum_{i=1}^d \lambda_i \geq \lambda_1$$

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

Let $\mathbf{e}_1, \dots, \mathbf{e}_d$ be the collection of canonical coordinate vectors and observe that $\|\mathbf{e}_i\| = 1$.

$$\lambda_1 = \max_{\mathbf{u}: \|\mathbf{u}\|=1} \mathbf{u}^\top \Sigma \mathbf{u} \geq \max_{\mathbf{u} \in \{\mathbf{e}_1, \dots, \mathbf{e}_d\}} \mathbf{u}^\top \Sigma \mathbf{u} = \max_{i=1}^d \mathbf{e}_i^\top \Sigma \mathbf{e}_i = \max_{i=1}^d \Sigma_{ii}$$

Exercise 3: Iterative PCA (5 + 10 + 5 P)

When performing principal component analysis, computing the full eigendecomposition of the covariance matrix Σ 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

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

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

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

will demonstrate that they all converge to zero as \mathbf{w} approaches the eigenvector \mathbf{u}_1 and becomes orthogonal to other eigenvectors. To demonstrate this, we proceed in three steps:

(a) Show that $\mathcal{E}_k(\mathbf{w}^{(t+1)}) = \left| \frac{\mathbf{w}^{(t)\top} \Sigma \mathbf{u}_k}{\mathbf{w}^{(t)\top} \Sigma \mathbf{u}_1} \right|$.

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

(b) Starting from the result above, show that $\mathcal{E}_k(\mathbf{w}^{(t+1)}) = |\lambda_k / \lambda_1| \cdot \mathcal{E}_k(\mathbf{w}^{(t)})$ (Hint: to show this, it is useful to recall that the covariance is linked to its eigenvectors and eigenvalues through the equation $\Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i$.)

$$\begin{aligned}
\mathcal{E}_k(\mathbf{w}^{(t+1)}) &= \left| \frac{\mathbf{w}^{(t)\top} \Sigma \mathbf{u}_k}{\mathbf{w}^{(t)\top} \Sigma \mathbf{u}_1} \right| \\
&= \left| \frac{\mathbf{w}^{(t)\top} \lambda_k \mathbf{u}_k}{\mathbf{w}^{(t)\top} \lambda_1 \mathbf{u}_1} \right| \\
&= \left| \frac{\mathbf{w}^{(t)\top} \mathbf{u}_k}{\mathbf{w}^{(t)\top} \mathbf{u}_1} \right| \cdot \left| \frac{\lambda_k}{\lambda_1} \right| \\
&= \mathcal{E}_k(\mathbf{w}^{(t)}) \cdot \left| \frac{\lambda_k}{\lambda_1} \right|
\end{aligned}$$

(c) Starting from the result above, 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 .

$$\mathcal{E}_k(\mathbf{w}^{(T)}) = \mathcal{E}_k(\mathbf{w}^{(T-1)}) \cdot \left| \frac{\lambda_k}{\lambda_1} \right| = \dots = \mathcal{E}_k(\mathbf{w}^{(0)}) \cdot \left| \frac{\lambda_k}{\lambda_1} \right|^T$$