

**Set 3 - Principal Component Analysis and Sanger's rule**

Issued: October 29, 2021

Hand in (optional): November 12, 2021 12:00am

**Question 1: Principal Component Analysis with Sanger's Rule (25 points)**

In this question we want to investigate methods for computing the Eigen values and vectors of a dataset's co-variance matrix:

$$\mathbf{C} := \frac{1}{N-1} \mathbf{X}^T \mathbf{X}, \text{ with } \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \dots \\ \mathbf{x}_N^T \end{bmatrix} \in \mathbb{R}^{N \times n} \quad (1)$$

All methods are based on the Perceptron, but use different rules for updating the weights. Importantly, the updates strongly depend on the output from the forward-pass of the Perceptron, as defined in equation (2):

$$\mathbf{O} = \mathbf{XW}, \text{ with } \mathbf{W} \in \mathbb{R}^{n \times k} \quad (2)$$

Here, the weight matrix,  $\mathbf{W}$ , stores the  $k$  weight vectors as column vectors. We want to define our update procedure in such a manner that the  $i$ -th column vector,  $\mathbf{w}_i$ , converges to the  $i$ -th Eigen vector with unit length of the co-variance matrix.

- [5 points]** Implement the forward pass of the Perceptron in the file *perceptron.cpp*. Use a single call to BLAS in your implementation.
- [20 points]** In order to find Eigen values and vectors, we need to consider the evolution of our weights. We define  $\mathbf{O}^t$  to be the output matrix from step  $t$ .

$$\mathbf{O}^t = \begin{bmatrix} \mathbf{o}_1^t & \dots & \mathbf{o}_n^t \end{bmatrix}, \text{ such that } \mathbf{o}_i^t = \mathbf{Xw}_i^t \quad (3)$$

As indicated, we see that our weights evolve over time and  $\mathbf{w}_i^t$  is the  $i$ -th weight vector at step  $t$ . We then need to define the gradient update rules, which is done in the following manner:

- Hebb's Rule:

$$\nabla \mathbf{w}_i^t \propto \mathbf{X}^T \mathbf{o}_i^t \quad (4)$$

2. Oja's Rule:

$$\begin{aligned}\nabla \mathbf{w}_i^t &\propto \mathbf{X}^T \mathbf{o}_i^t - \left( (\mathbf{o}_i^t)^T \mathbf{o}_i^t \right) \mathbf{w}_i \\ &= \mathbf{X}^T \mathbf{o}_i^t - \|\mathbf{o}_i^t\|^2 \mathbf{w}_i\end{aligned}\tag{5}$$

3. Sangers's Rule:

$$\begin{aligned}\nabla \mathbf{w}_i^t &\propto \mathbf{X}^T \mathbf{o}_i^t - \sum_{j=1}^i \left( (\mathbf{o}_i^t)^T \mathbf{o}_j^t \right) \mathbf{w}_j \\ &= \mathbf{X}^T \mathbf{o}_i^t - \begin{bmatrix} \mathbf{w}_1 & \dots & \mathbf{w}_i \end{bmatrix} \begin{bmatrix} \mathbf{o}_1^T \\ \dots \\ \mathbf{o}_i^T \end{bmatrix} \mathbf{o}_i\end{aligned}\tag{6}$$

We only specify the gradient as being proportional to the given formulas. The reason is that the length of our update is equal to a specified learning rate,  $\eta$ . The update step is then given by equation (7):

$$\mathbf{w}_i^{t+1} = \mathbf{w}_i^t + \eta \cdot \frac{\nabla \mathbf{w}_i^t}{\|\nabla \mathbf{w}_i^t\|}\tag{7}$$

The Hebb and Oja rules have already been defined in *hebb.cpp* and *oja.cpp* respectively and can be used for inspiration. Your task is to implement Sanger's update rule in *sanger.cpp* as given by equation (6). Also in this task, each TO DO corresponds to a single function call; either to BLAS or another class method.

*Hint:* You can test your implementation using the unit test - just run "make unit\_test". If you pass all tests, your implementation is very likely to be correct. Also, in order to investigate the different methods closer, you can create plots by running "make plots". The plots display the evolution of the Eigen values and the final Eigen vectors. They can therefore both be used for debugging and gaining more insight into the different methods.

*Hint:* Reference for [BLAS](#) and [C-BLAS](#).

## Question 2: Linear Auto Encoders (25 points)

This question will cover the fundamental theory of linear neural networks and more specifically linear auto-encoders.

A linear neural network is, just as any conventional feed-forward neural network, a concatenation of structured mappings between layers; in each layer, we first apply an affine transformation before employing an activation function to each individual output. In contrast to other neural networks, in a linear neural network we can combine the two steps into a single affine transformation. Hence, we can write a general linear neural network as given in equation (8) with the definition of a linear layer as given in equation (9).

$$f : \mathbb{R}^n \longrightarrow \mathbb{R}^m, \mathbf{x} \mapsto \left( \phi_{\mathbf{W}_{L+1}, \mathbf{b}_{L+1}}^{m, h_L} \circ \dots \circ \phi_{\mathbf{W}_2, \mathbf{b}_2}^{h_2, h_1} \circ \phi_{\mathbf{W}_1, \mathbf{b}_1}^{h_1, n} \right) (\mathbf{x}) \quad (8)$$

$$\phi_{\mathbf{W}, \mathbf{b}}^{k, l} : \mathbb{R}^l \longrightarrow \mathbb{R}^k, \mathbf{x} \mapsto \mathbf{W}\mathbf{x} + \mathbf{b}, \text{ with } \mathbf{W} \in \mathbb{R}^{k \times l}, \mathbf{b} \in \mathbb{R}^k \quad (9)$$

a) [10 points] Prove that any linear neural network is an affine mapping.

*Hint:* Prove by induction on the number of layers that for any linear neural network  $f : \mathbb{R}^n \longrightarrow \mathbb{R}^m$ , there exist  $\mathbf{W} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$  such that  $\forall \mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}$ .

b) [15 points] In the linear auto-encoder we consider a linear neural network with a single hidden layer of size  $k$ . Let  $\mathbf{X} \in \mathbb{R}^{N \times n}$  be the data matrix with the individual data points as row vectors. Furthermore, we generally have the following decomposition of  $\mathbf{X}$ , also known as the SVD-decomposition:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \dots \\ \mathbf{x}_N^T \end{bmatrix} = \sum_{i=1}^{\min\{N, n\}} \sigma_i^2 \mathbf{u}_i \mathbf{v}_i^T \quad (10)$$

The  $\sigma$ -values are ordered  $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_{\min\{N, n\}}^2 \geq 0$ , and the  $\mathbf{u}$ - and  $\mathbf{v}$ -vectors are orthonormal, meaning that  $\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$  and  $\mathbf{v}_i^T \mathbf{v}_j = \delta_{ij}$ . Lastly, we define  $\mathbf{V}_k := [\mathbf{v}_1 \dots \mathbf{v}_k]$  with  $k \leq \min\{N, n\}$ .

Based on these definitions prove that  $(\mathbf{V}_k, \mathbf{V}_k^T)$  is an optimizer of the objective of the linear auto-encoder as given in equation (11). How does this result relate to the PCA of  $\mathbf{X}$ ?

$$(\mathbf{V}_k, \mathbf{V}_k^T) \in \arg \min_{U \in \mathbb{R}^{n \times k}, W \in \mathbb{R}^{k \times n}} \sum_{i=1}^N \left\| \left( \phi_{\mathbf{U}, \mathbf{0}}^{n, k} \circ \phi_{\mathbf{W}, \mathbf{0}}^{k, n} \right) (\mathbf{x}_i) - \mathbf{x}_i \right\|^2 \quad (11)$$

*Hint:* Use the Eckart-Young-Mirsky theorem which states that  $\sum_{i=1}^k \sigma_i^2 \mathbf{u}_i \mathbf{v}_i^T$  is the best rank  $k$  approximation of  $\mathbf{X}$ . Also, the following identities might prove useful:

- $A \in \mathbb{R}^{m \times n} \rightarrow \text{Rank}(A) \leq \min\{m, n\}$
- $\text{Rank}(AB) \leq \min\{\text{Rank}(A), \text{Rank}(B)\}$

## References

- [1] Oja, Erkki Simplified neuron model as a principal component analyzer, *Journal of mathematical biology*, 1982.
- [2] Oja, Erkki, Principal components, minor components, and linear neural networks, *Neural networks*, 1992.
- [3] Terence, D. Sanger, Optimal unsupervised learning in a single-layer linear feedforward neural network, *Neural Networks*, 1989.
- [4] Bruno, A. Olshausen, Linear Hebbian learning and PCA.
- [5] Baldi, Pierre and Hornik, Kurt, Neural networks and principal component analysis: Learning from examples without local minima, *Neural networks*, 1989.