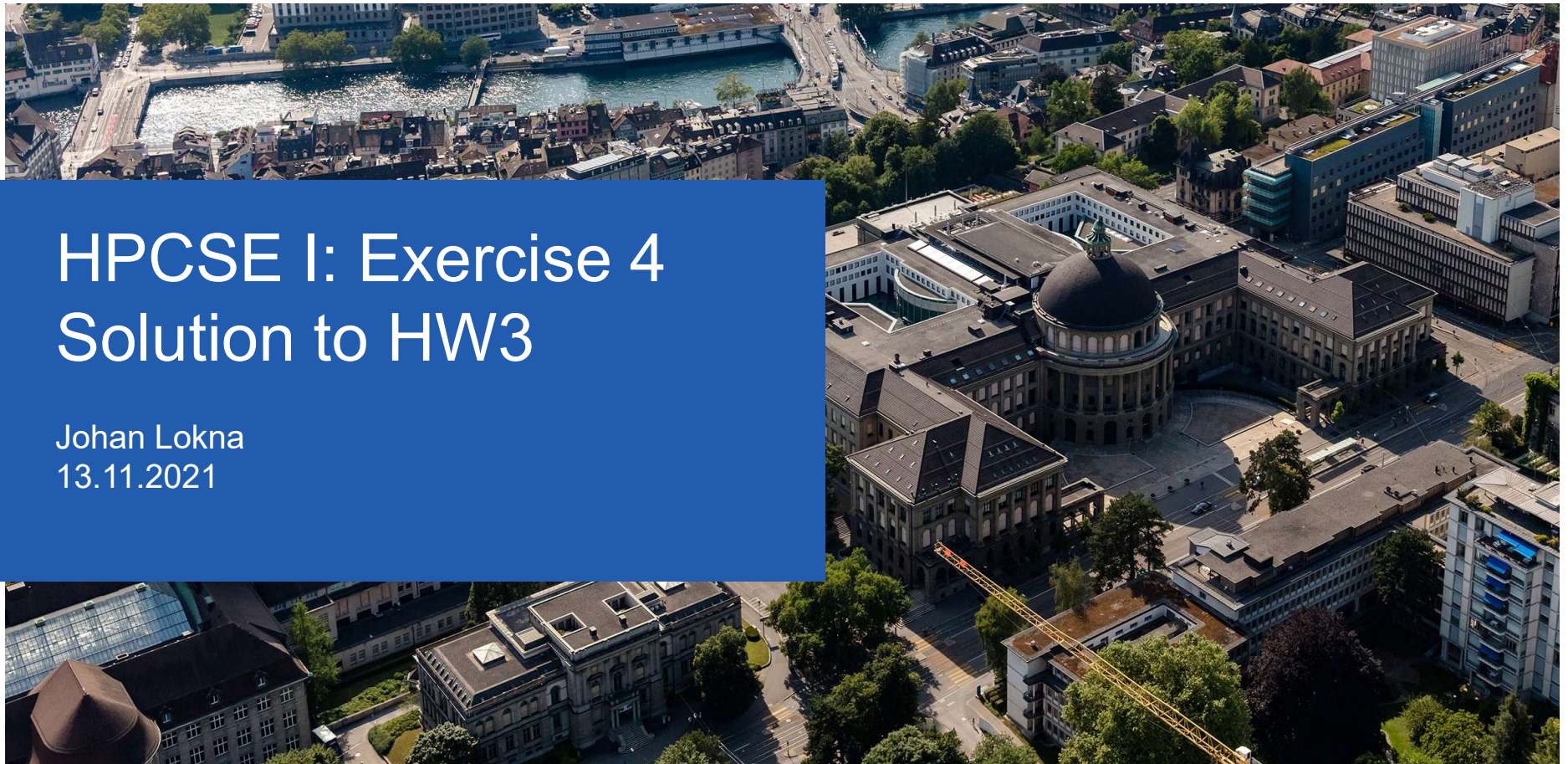


HPCSE I: Exercise 4 Solution to HW3

Johan Lokna
13.11.2021



Agenda

1. Solution to HW 3

- a) Question 1 – Perceptron
- b) Question 2 – Linear Auto-Encoders

2. Introduction to HW 4

Question 1a – Formulation

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):

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

[5 points] Implement the forward pass of the Perceptron in the file *perceptron.cpp*. Use a single call to BLAS in your implementation.

Question 1a – Solution

```
1. void Perceptron::forward(const scalar* data, scalar* out, const int batch_size) {  
2.  
3.     cblas_dgemm(CblasRowMajor, CblasNoTrans, CblasNoTrans,  
4.                 batch_size, out_dims_, in_dims_,  
5.                 1.0, data, in_dims_, weights_, out_dims_,  
6.                 0.0, out, out_dims_  
7.     );  
8.  
9. }
```

Question 1b – Formulation

- b) [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 .

$$\mathbf{O}^t = [\mathbf{o}_1^t \dots \mathbf{o}_n^t], \text{ such that } \mathbf{o}_i^t = \mathbf{X}\mathbf{w}_i^t \quad (3)$$

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:

3. Sangers's Rule:

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

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.

Question 1b – Solution

```
1. void SangersRule::backward(const scalar* data, int batch_size) {
2.
3.     forward(data, outputs_, batch_size);
4.
5.     cblas_dgemm(CblasRowMajor, CblasTrans, CblasNoTrans,
6.                 in_dims_, out_dims_, batch_size,
7.                 1.0, data, in_dims_, outputs_, out_dims_,
8.                 0.0, grad_, out_dims_
9.     );
10.
11.     for(int k = 0; k < out_dims_; ++k) {
12.
13.         cblas_dgemv(CblasRowMajor, CblasTrans,
14.                     batch_size, k + 1,
15.                     1.0, outputs_, out_dims_, outputs_ + k, out_dims_,
16.                     0.0, inference_, 1
17.         );
18.
19.         cblas_dgemv(CblasRowMajor, CblasNoTrans,
20.                     in_dims_, k + 1,
21.                     -1.0, weights_, out_dims_, inference_, 1,
22.                     1.0, grad_ + k, out_dims_
23.         );
24.
25.     }
26.
27. }
```


Question 1 – Grading

For each of the four TO DOs, the points are awarded accordingly:

- Correct method gives (**1 point**).
- Operating on the correct matrices and vectors gives (**1 point**).
- Using correct indexing gives (**1 point**).
- Passing the unit test gives (**1 point**).
- Using a single function call as specified in the description gives (**1 point**).

Hint: In order to test your implementation of each TO DO individually, you can copy the master solution for the other TO DOs and then run the unit test.

Question 2a – Formulation

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.

Question 2a – Solution & Grading

1) Formulate a predicate

2 Points

2) Prove the base case

2 Points

3) Assume the predicate holds for some value

1 Points

4) Induction step

4 Points

5) Conclusion

2 Points

$P(k) = 1 : \Leftrightarrow$ All linear neural networks with k layers are affine transformations

Using $k = 0$, we have the identity mapping, which is affine with $\mathbf{W} = \mathbf{I}$ and $\mathbf{b} = \mathbf{0}$.

Using $k = 1$, we only have one layer which is affine by definition.

For some $k \in \mathbb{N}$, assume $P(k) = 1$

$$\begin{aligned} & \left(\phi_{\mathbf{W}_{k+1}, \mathbf{b}_{k+1}}^{h_{k+1}, h_k} \circ \phi_{\mathbf{W}_k, \mathbf{b}_k}^{h_k, h_{k-1}} \circ \dots \circ \phi_{\mathbf{W}_1, \mathbf{b}_1}^{h_1, h_0} \right) (\mathbf{x}) \\ &= \phi_{\mathbf{W}_{k+1}, \mathbf{b}_{k+1}}^{h_{k+1}, h_k} \left(\left(\phi_{\mathbf{W}_k, \mathbf{b}_k}^{h_k, h_{k-1}} \circ \dots \circ \phi_{\mathbf{W}_1, \mathbf{b}_1}^{h_1, h_0} \right) (\mathbf{x}) \right) \\ &= \phi_{\mathbf{W}_{k+1}, \mathbf{b}_{k+1}}^{h_{k+1}, h_k} (\mathbf{W}\mathbf{x} + \mathbf{b}) \\ &= \mathbf{W}_{k+1}(\mathbf{W}\mathbf{x} + \mathbf{b}) + \mathbf{b}_{k+1} \\ &= \mathbf{W}_{k+1}\mathbf{W}\mathbf{x} + (\mathbf{W}_{k+1}\mathbf{b} + \mathbf{b}_{k+1}) = \mathbf{A}\mathbf{x} + \mathbf{c} \end{aligned}$$

$\forall k \in \mathbb{N}: P(k) = 1$

Question 2b – Formulation

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

The σ -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} ?

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

Question 2a – Solution & Grading

The first step of the proof is to transform the objective into matrix form:

$$\begin{aligned}\sum_{i=1}^N \left\| \left(\phi_{\mathbf{U},0}^{n,k} \circ \phi_{\mathbf{W},0}^{k,n} \right) (\mathbf{x}_i) - \mathbf{x}_i \right\|_2^2 &= \sum_{i=1}^N \left\| \mathbf{U} \mathbf{W} \mathbf{x}_i - \mathbf{x}_i \right\|_2^2 \\ &= \left\| \mathbf{X} \mathbf{W}^T \mathbf{U}^T - \mathbf{X} \right\|_F^2 \\ &= \left\| \tilde{\mathbf{X}}_{\mathbf{k}} - \mathbf{X} \right\|_F^2\end{aligned}$$

Correctly reforming the objective gives (**2 point**).

Question 2a – Solution & Grading

Then, we compute the rank of $\tilde{\mathbf{X}}_k$:

$$\begin{aligned} \text{Rank}(\tilde{\mathbf{X}}_k) &= \text{Rank}(\mathbf{X}\mathbf{U}^T\mathbf{W}^T) \\ &\leq \min\{\text{Rank}(X), \text{Rank}(U), \text{Rank}(W)\} \\ &\leq \min\{\min\{n, N\}, \min\{n, k\}, \min\{n, k\}\} \\ &= \min\{n, k, N\} = k \end{aligned}$$

Correctly computing an upper bound on the rank gives (**1 point**).

Question 2a – Solution & Grading

$$\begin{aligned}\tilde{\mathbf{X}}_k &= \mathbf{X} \mathbf{W}^T \mathbf{U}^T \\ &= \mathbf{X} \mathbf{V}_k \mathbf{V}_k^T \\ &= \left(\sum_{i=1}^{\min\{N,n\}} \sigma_i^2 \mathbf{u}_i \mathbf{v}_i^T \right) \mathbf{V}_k \mathbf{V}_k^T \\ &= \left(\sum_{i=1}^{\min\{N,n\}} \sigma_i^2 \mathbf{u}_i (\mathbf{v}_i^T \mathbf{V}_k) \right) \mathbf{V}_k^T \\ &= \left(\sum_{i=1}^k \sigma_i^2 \mathbf{u}_i \mathbf{e}_i^T \right) \mathbf{V}_k^T \\ &= \sum_{i=1}^k \sigma_i^2 \mathbf{u}_i (\mathbf{e}_i^T \mathbf{V}_k^T) \\ &= \sum_{i=1}^k \sigma_i^2 \mathbf{u}_i \mathbf{v}_i^T\end{aligned}$$

Proving that that using $\mathbf{U} = \mathbf{V}_k$ and $\mathbf{W} = \mathbf{V}_k^T$, results in $\tilde{\mathbf{X}}_k = \sum_{i=1}^k \sigma_i^2 \mathbf{u}_i \mathbf{v}_i^T$, gives (4 point).

Question 2a – Solution & Grading

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} .

Concluding using the Eckart-Young-Mirsky theorem gives (5 point).

Question 2a – Solution & Grading

There are several relations to PCA, here we point out a few:

1. One relation to PCA is that $\mathbf{v}_1, \dots, \mathbf{v}_k$ are the first k principal components of the data matrix. So a solution to PCA can be directly used to construct optimal weights to the linear auto-encoder.
2. Another is that the objective in equation (14) equals the minimal reconstruction loss if we allow for arbitrary pairs of linear encoding and decoding mappings.
3. A third relation is that both can be seen as mapping higher dimensional data onto a lower dimensional space while trying to maintain as much information as possible.

Pointing out one of the relations to PCA gives (3 point).