## Exercise 1: Sparse Coding (20 + 20 P)

Let $x_1, \ldots, x_N \in \mathbb{R}^d$ be a dataset of $N$ examples. Let $z_i \in \mathbb{R}^h$ be the source associated to example $x_i$, and $W$ be a matrix of size $d \times h$ that linearly reconstructs the examples from the sources. We wish to minimize the objective:

$$J = \frac{1}{N} \sum_{i=1}^{N} \underbrace{\|x_i - W z_i\|^2}_{\text{reconstruction}} + \lambda \cdot \underbrace{\frac{1}{N} \sum_{i=1}^{N} \|z_i\|_1}_{\text{sparsity}} + \eta \underbrace{\|W\|_F^2}_{\text{regularization}}$$

with respect to the weights $W$ and the sources $z_1, \ldots, z_N$. The objective consists of three terms. The reconstruction term is the standard mean square error, the sparsity term consists of a standard $L_1$ penalty on the sources, and the last regularization term prevents the sparsity term from becoming ineffective.

(a) Show that for fixed sources, the optimal matrix $W$ is given in closed form as:

$$W = \Sigma_{XS}(\Sigma_{SS} + \eta I)^{-1}$$

where

$$\Sigma_{SS} = \frac{1}{N} \sum_{i=1}^{N} z_i z_i^\top \qquad \text{and} \qquad \Sigma_{XS} = \frac{1}{N} \sum_{i=1}^{N} x_i z_i^\top.$$

$$J(W) + \frac{\eta}{N} \sum_i (x_i - W z_i)^\top I (x_i - W z_i) \quad + \text{cst.} \quad + \eta \|W\|_F^2$$
$$= \frac{1}{N} \sum_i - 2 z_i W^\top z_i + W z_i z_i^\top W^\top + \text{cst} + \eta \|W\|_F^2$$

$$\frac{\partial J}{\partial W} = \frac{1}{N} \sum_i - 2 (x_i - W z_i) z_i^\top + \eta 2 W \overset{!}{=} 0$$

$$\iff W \left( \frac{1}{N} \sum_i z_i z_i^\top + \eta I \right) = \frac{1}{N} \sum_i x_i z_i^\top$$
$$\underbrace{\phantom{W \left( \frac{1}{N} \sum_i z_i z_i^\top + \eta I \right)}}_{\Sigma_{SS}} \qquad \underbrace{\phantom{\frac{1}{N} \sum_i x_i z_i^\top}}_{\Sigma_{XS}}$$

$$\iff W = \Sigma_{XS}(\Sigma_{SS} + \eta I)^{-1}$$

(b) We now consider the optimization of sources. Due to the 1-norm in the sparsity term, we cannot find a closed form solution. However, we consider a local relaxation of the optimization problem where the 1-norm of the sparsity term is linearized as

$$\|z_i\|_1 = q_i^\top z_i$$

with $q_i \in \{-1, 0, 1\}^d$ a constant vector. This relaxation makes the objective function quadratic with $z_i$.

Show that under this local relaxation, the solution of the optimization problem is given in closed form as:

$$z_i = (W^\top W)^{-1} (W^\top x_i - \lambda \cdot q_i / 2)$$

Although this solution is not the true minimum of $J$ (e.g. it is not sparse), it can serve as the end-point of some line-search method for finding good source vectors.

$$\frac{\partial J}{\partial z_i} = \frac{\lambda}{N} \cdot \left( - 2 W^\top (x_i - W z_i) \right) + \frac{\lambda}{N} \lambda q_i \overset{!}{=} 0$$

$$\iff W^\top W z_i = W^\top x_i - \frac{\lambda}{2} q_i$$

$$\iff z_i = (W^\top W)^{-1} (W^\top x_i - \lambda \cdot q_i / 2)$$

## Exercise 2: Auto-Encoders (20 P)

In this exercise, we would like to show an equivalence between linear autoencoders with tied weights (same parameters for the encoder and decoder) and PCA. We consider the special case of an autoencoder with a single hidden unit. In that case, the autoencoder consists of the two layers:

$$z_i = w^\top x_i \qquad \text{(encoder)}$$
$$\hat{x}_i = w z_i \qquad \text{(decoder)}$$

where $w \in \mathbb{R}^d$. We consider a dataset $x_1, \ldots, x_N$ assumed to be centered (i.e. $\sum_i x_i = 0$), and we define the objective that we would like to minimize to be the mean square error between the data and the reconstruction:

$$J(w) = \frac{1}{N} \sum_{i=1}^{N} \|x_i - \hat{x}_i\|^2$$

Furthermore, to make the objective closer to PCA, we can always rewrite the weight vector as $w = \alpha u$ where $u$ is a unit vector (of norm 1) and $\alpha$ some positive scalar, and search instead for the optimal parameters $u$ and $\alpha$.

(a) Show that the optimization problem can be equally rewritten as

$$\arg\min_{u, \alpha} J(w) = \arg\max_{u, \alpha} \; u^\top S u$$

where $S = \frac{1}{N} \sum_{i=1}^{N} x_i x_i^\top$, which is a common formulation of PCA.

$$J(w) = \frac{1}{N} \sum_i \|x_i - w w^\top x_i\|^2 = \frac{1}{N} \sum_i \|x_i - \alpha^2 u u^\top x_i\|^2$$

$$= \frac{1}{N} \sum_i \|x_i\|^2 - 2 \alpha^2 x_i^\top u u^\top x_i + \alpha^4 x_i^\top u u^\top u u^\top x_i$$
$$\underbrace{\phantom{u^\top u}}_{= 1}$$

$$= \frac{1}{N} \sum_i (x_i^\top x_i) \cdot u^\top x_i \; x_i^\top u + \text{cst.}$$

$$= (\alpha^4 - 2\alpha^2) \cdot \frac{u^\top \left[ \frac{1}{N} \sum_i x_i x_i^\top \right] u}{\underbrace{\phantom{u^\top \left[ \frac{1}{N} \sum_i x_i x_i^\top \right] u}}_{\geq 0 \; (S \text{ is p.s.d.})}}$$
$$\underbrace{\phantom{(\alpha^4 - 2\alpha^2)}}_{\leq 0 \text{ if } \ldots}$$

$$\implies \arg\min_{\alpha, u} J = \arg\max_{\alpha, u} \; u^\top S u$$