

# EECS208 Written HW4

**Issued:** Nov. 4. **Due:** Nov. 14, 11:59 PM via Gradescope

**Reading:** Chapters 7 of *High-Dim Data Analysis with Low-Dim Models*.

## Problem 1 (Complete Dictionary Learning via $\ell^4$ Norm Maximization)

*Exercise 7.2 of High-Dim Data Analysis with Low-Dim Models.* In this exercise, we derive and practice an algorithm to solve  $\ell^4$  norm maximization problem (7.3.14) for complete dictionary learning.

1. Derive the gradient  $\varphi(\mathbf{A}) = \|\mathbf{A}\mathbf{Y}\|_4^4$  with respect to  $\mathbf{A}$ .
2. Derive a project gradient ascent for maximizing  $\varphi(\mathbf{A})$ :

$$\mathbf{A}_{k+1} = \mathcal{P}_{\mathcal{O}(n)}[\mathbf{A}_k + \gamma \cdot \nabla\varphi(\mathbf{A}_k)]. \quad (0.1)$$

Hint: you may directly apply the following claim (without proof):

**Claim 0.1 (Projection onto Orthogonal Group)**  $\forall \mathbf{A} \in \mathbb{R}^{n \times n}$ , the orthogonal matrix which has minimum Frobenius norm with  $\mathbf{A}$  is the following

$$\mathcal{P}_{\mathcal{O}(n;\mathbb{R})}(\mathbf{A}) = \arg \min_{\mathbf{M} \in \mathcal{O}(n;\mathbb{R})} \|\mathbf{M} - \mathbf{A}\|_F^2 = \mathbf{U}\mathbf{V}^\top, \quad (0.2)$$

where  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top = \text{SVD}(\mathbf{A})$ .

If you can prove Claim 0.1, you will **receive extra credits** for this problem.

3. Let  $\mathbf{Y} \in \mathbb{R}^{n \times p}$  and suppose each entry of  $\mathbf{Y}$  is iid drawn from a standard Gaussian distribution  $\mathcal{N}(0, 1)$  with probability  $\theta$ , and equals to 0 with probability  $1 - \theta$ :

$$\bar{y}_{i,j} = \begin{cases} z \sim \mathcal{N}(0, 1) & \text{with probability } \theta, \\ 0, & \text{otherwise.} \end{cases} \quad (0.3)$$

and  $\mathbf{A}_0$  is a randomly initialized orthogonal matrix. Conduct simulation of the algorithm and play with different step size  $\gamma$  of the gradient ascent, what is the relationship between the step size  $\gamma$  and the convergence speed? What happens if you make the step size to be infinite? That is,

$$\mathbf{A}_{k+1} = \mathcal{P}_{\mathcal{O}(n)}[\nabla\varphi(\mathbf{A}_k)]. \quad (0.4)$$

## Problem 2 (Low-rank Regularization via the $\log \det(\cdot)$ Function)

*Exercise 7.4 of High-Dim Data Analysis with Low-Dim Models.* When a matrix  $\mathbf{X} \in \mathbb{R}^{n \times n}$  is symmetric and positive semi-definite, the nuclear norm  $\|\mathbf{X}\|_*$  is the same as its trace of the matrix. In this exercise, we try to study the connection of the convex nuclear norm (or the trace norm) with another popular smooth but nonconvex surrogate for minimizing  $\text{rank}(\mathbf{X})$  is to minimize the quantity

$$\min_{\mathbf{X} \in \mathcal{C}} f(\mathbf{X}) \doteq \log \det(\mathbf{X} + \delta \mathbf{I}), \quad (0.5)$$

where  $\delta > 0$  is a small regularization constant and  $\mathbf{X}$  belongs to some constraint set  $\mathcal{C}$ . To see how this objective is related to the trace norm:

1. First, show that  $\nabla_{\mathbf{X}} f(\mathbf{X}) = (\mathbf{X} + \delta \mathbf{I})^{-1}$ .
2. Second, the first-order expansion of  $f(\mathbf{X})$  around a point  $\mathbf{X}_k$  is given by:

$$f(\mathbf{X}) \approx f(\mathbf{X}_k) + \text{tr}((\mathbf{X}_k + \delta \mathbf{I})^{-1}(\mathbf{X} - \mathbf{X}_k)) + o(\|\mathbf{X} - \mathbf{X}_k\|). \quad (0.6)$$

Then to minimize  $f(\mathbf{X})$ , we can use a greedy descent algorithm with the iteration

$$\mathbf{X}_{k+1} = \arg \min_{\mathbf{X} \in \mathcal{C}} \text{tr}((\mathbf{X}_k + \delta \mathbf{I})^{-1} \mathbf{X}). \quad (0.7)$$

Argue that when  $\mathbf{X}_k$  is initialized around  $\mathbf{X}_o = \mathbf{I}$ , then the above iteration becomes minimizing the trace norm  $\mathbf{X}_{k+1} = \arg \min_{\mathbf{X} \in \mathcal{C}} \text{tr}(\mathbf{X})$ .