## Homework 1: Mathematics of Deep Learning (EN 580.745)

Instructor: René Vidal, Biomedical Engineering, Johns Hopkins University

## Due Date: 10/24/2019, 11:59PM Eastern Time

## Instructions.

- You can discuss the problems with your peers at a high level, but you must write your own solutions.
- You can consult relevant background material, but don't seek out solutions to the problems themselves. Cite the outside material use.
- Submit a single PDF to Blackboard. Typed solutions are nice, but clearly handwritten solutions are also fine.

**1. Properties of the**  $\ell_1$  **norm.** Let  $x \in \mathbb{R}^n$  and recall the  $\ell_1$  norm  $||x||_1 \triangleq \sum_i |x_i|$ .

(a) (10 points) The subdifferential  $\partial f(x)$  of a convex function  $f : \mathbb{R}^n \to \mathbb{R}$  at x is defined to be

$$\partial f(\boldsymbol{x}) \triangleq \{ \boldsymbol{z} \in \mathbb{R}^n \mid f(\boldsymbol{y}) \ge f(\boldsymbol{x}) + \langle \boldsymbol{z}, \boldsymbol{y} - \boldsymbol{x} \rangle \text{ for all } \boldsymbol{y} \in \mathbb{R}^n \}.$$
(1)

Prove that

$$\partial \|\boldsymbol{x}\|_{1} = \{\operatorname{sign}(\boldsymbol{x}) + \boldsymbol{w} \mid \operatorname{supp}(\boldsymbol{w}) \subseteq \operatorname{supp}(\boldsymbol{x})^{c}, \max_{i} |w_{i}| \leq 1\},$$
(2)

where sign:  $\mathbb{R}^n \to \{-1, 0, 1\}$  denotes the sign function, and supp $(x) \subseteq [n]$  denotes the support of x, that is the set of indices where x is non-zero.

(b) (10 points) Define the proximal operator  $\operatorname{prox}_f \colon \mathbb{R}^n \to \mathbb{R}^n$  of a convex function f to be

$$\operatorname{prox}_{f}(\boldsymbol{x}) \triangleq \underset{\boldsymbol{a} \in \mathbb{R}^{n}}{\operatorname{arg\,min}} \ \frac{1}{2} \|\boldsymbol{a} - \boldsymbol{x}\|_{2}^{2} + f(\boldsymbol{a}).$$
(3)

Prove that  $\operatorname{prox}_{\tau \parallel \cdot \parallel_1}(\boldsymbol{x}) = \mathcal{S}_{\tau}(\boldsymbol{x})$ , where  $\mathcal{S}_{\tau}$  is the elementwise soft-thresholding operator

$$\mathcal{S}_{\tau}(x_i) \triangleq \operatorname{sign}(x_i) \max(|x_i| - \tau, 0).$$
(4)

**2.** Properties of the nuclear norm. Let  $X \in \mathbb{R}^{D \times N}$  be a matrix of rank r. Recall the nuclear norm  $||X||_* \triangleq \sum_{i=1}^r \sigma_i(X)$ , where  $\sigma_i(X)$  denotes the *i*th singular value of X. Let  $X = U\Sigma V^{\top}$  be the compact SVD, so that  $U \in \mathbb{R}^{D \times r}, \Sigma \in \mathbb{R}^{r \times r}$ , and  $V \in \mathbb{R}^{N \times r}$ . Recall also the spectral norm  $||X||_2 = \sigma_1(X)$ .

(a) (10 points) Prove that  $^{1}$ 

$$\partial \|\boldsymbol{X}\|_* = \{\boldsymbol{U}\boldsymbol{V}^\top + \boldsymbol{W} \mid \boldsymbol{U}^\top \boldsymbol{W} = \boldsymbol{0}, \, \boldsymbol{W}\boldsymbol{V} = \boldsymbol{0}, \, \|\boldsymbol{W}\|_2 \le 1\}.$$
(5)

- (b) (10 points) Let τ > 0. Prove that prox<sub>τ ||·||\*</sub>(X) = US<sub>τ</sub>(Σ)V<sup>T</sup>, where S<sub>τ</sub> is as in (4). Note that when generalizing (3) to matrices, the squared Frobenius norm <sup>1</sup>/<sub>2</sub> || A − X ||<sup>2</sup><sub>F</sub> is used in place of the squared ℓ<sub>2</sub> norm (||X||<sub>F</sub> ≜ (∑<sub>ij</sub> X<sup>2</sup><sub>ij</sub>)<sup>1/2</sup>).
- **3.** Low rank matrix factorization. Let  $Y \in \mathbb{R}^{D \times N}$  be of rank r. Consider the matrix factorization problem

minimize 
$$F(\boldsymbol{U}, \boldsymbol{V}) = \frac{1}{2} \| \boldsymbol{Y} - \boldsymbol{U} \boldsymbol{V}^{\top} \|_{F}^{2} + \frac{\tau}{2} (\| \boldsymbol{U} \|_{F}^{2} + \| \boldsymbol{V} \|_{F}^{2}),$$
 (6)

where  $U \in \mathbb{R}^{D \times r}$ ,  $V \in \mathbb{R}^{N \times r}$ . This is in fact a regularized variant of the problem studied in (Baldi & Hornik, 1989). It is also closely related to the nuclear norm proximal operator from Problem 2(b), as we will see.

(a) (10 points) Let  $Y = P \Sigma Q^{\top}$  be the compact SVD. Define

$$\widehat{\boldsymbol{U}} = \boldsymbol{P} \mathcal{S}_{\tau}(\Sigma)^{\frac{1}{2}} \qquad \widehat{\boldsymbol{V}} = \boldsymbol{Q} \mathcal{S}_{\tau}(\Sigma)^{\frac{1}{2}},\tag{7}$$

<sup>&</sup>lt;sup>1</sup>*Hint:* the Von-Neumann trace inequality could be useful.

where the square root of a diagonal matrix is applied elementwise. Prove that  $(\hat{U}, \hat{V})$  is a critical point of (6).

(b) (bonus) Prove that  $(\widehat{U}, \widehat{V})$  is in fact a global minimizer.

Problem (6) and the nuclear norm proximal operator are therefore closely related, in the sense that the global minimizer  $(\hat{U}, \hat{V})$  satisfies  $\hat{U}\hat{V}^{\top} = \text{prox}_{\tau \parallel \cdot \parallel_*}(Y)$ .