

# Homework 3: Mathematics of Deep Learning (EN 580.745)

**Instructor:** René Vidal, Biomedical Engineering, Johns Hopkins University

**Due Date:** 12/17/2019, 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.
- You are free to use any of the well-known libraries (e.g. Tensorflow, PyTorch, Scikit-Learn etc.).
- Submit a zip file containing a single scanned PDF writeup along with your code and figures. For code, a Jupyter notebook is preferred. Typed writeups are nice, but clearly handwritten solutions are also fine.

**1. Variant of matrix factorization dropout.** Recall from lecture that dropout applied to matrix factorization,

$$\underset{U, V}{\text{minimize}} \quad \mathbb{E}_{z \sim \text{Ber}(\theta)^r} \|X - U \text{diag}(z) V^T\|_F^2, \quad (1)$$

introduces regularization on the product  $UV^T$  whose *lower convex envelope* is the squared nuclear norm  $\|UV^T\|_*^2$ . Furthermore, we showed that the dropout-induced regularizer is in fact equal to the squared nuclear norm, since the factors  $(U, V)$  can always be balanced.

In this exercise we will examine a variant of matrix factorization dropout that, perhaps surprisingly, induces squared nuclear norm regularization more directly. First, let  $U \in \mathbb{R}^{D \times r}$ ,  $V \in \mathbb{R}^{N \times r}$ . Assume without loss of generality that all rank-1 factors  $u_i v_i^T$  are unique<sup>1</sup>. We'll now define a distribution  $\mathcal{D}_{(U, V)}$  that samples normalized rank-1 factors  $(1/\|u_i v_i^T\|_F) u_i v_i^T$  with probability proportional to their contribution to  $UV^T$ .

Formally, let  $\mathcal{A} = \{\bar{u}\bar{v}^T \mid \|\bar{u}\bar{v}^T\|_F = 1\}$  be the set of unit norm rank-1 matrices, and define a probability mass  $p: \mathcal{A} \rightarrow [0, 1]$  for  $\mathcal{D}_{(U, V)}$  as follows

$$p(\bar{u}\bar{v}^T) = \begin{cases} \frac{\|u_i v_i^T\|_F}{\sum_{i=1}^r \|u_i v_i^T\|_F} & \text{if } \bar{u}\bar{v}^T = \frac{u_i v_i^T}{\|u_i v_i^T\|_F} \text{ for some } i \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

Now let  $\Theta(U, V) \triangleq \sum_{i=1}^r \|u_i v_i^T\|_F$ , and consider the stochastic optimization problem

$$\underset{U, V}{\text{minimize}} \quad \mathbb{E}_{\{\bar{u}_j \bar{v}_j^T\}_{j=1}^m \sim \mathcal{D}_{(U, V)}^m} \left\| X - \frac{\Theta(U, V)}{m} \sum_{j=1}^m \bar{u}_j \bar{v}_j^T \right\|_F^2, \quad (3)$$

where  $m \geq 2$ . According to this formulation, the data  $X$  are approximated by a random sample average of unit norm rank-one matrices, drawn independently from the distribution  $\mathcal{D}_{(U, V)}$ , which in turn is determined by the current factors  $(U, V)$ .

(a) **(15 points)** Compute the expectation in (3) to show that it is equivalent to

$$\underset{U, V}{\text{minimize}} \quad \|\tilde{X} - UV^T\|_F^2 + \frac{1}{m-1} \left( \sum_{i=1}^r \|u_i v_i^T\|_F \right)^2, \quad (4)$$

where  $\tilde{X} \triangleq \frac{m}{m-1} X$ .

(b) **(5 points)** As we've seen, the regularizer in (4) is in fact the variational form of the squared nuclear norm. We have therefore established a more direct connection between (3) and the squared nuclear norm, compared with (1). Despite this direct equivalence however, (3) is not as straightforward to optimize by SGD as is (1). Explain why not in 1-3 sentences.

<sup>1</sup>If not, we can consolidate the redundant factors without affecting  $UV^T$ .

**2. Implicit regularization for under-determined linear regression.** Consider the least-squares problem

$$\underset{\mathbf{x}}{\text{minimize}} f(\mathbf{x}) \triangleq \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2, \quad (5)$$

where  $\mathbf{A} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{x} \in \mathbb{R}^m$ ,  $\mathbf{y} \in \mathbb{R}^n$ , and  $\text{rank}(\mathbf{A}) = n < m$ . Furthermore, assume  $\mathbf{y} = \mathbf{A}\mathbf{x}^*$  for some  $\mathbf{x}^*$ .

(a) **(5 points)** Derive a closed-form expression for the solution(s) to (5).

(b) **(5 points)** Prove that (5) satisfies the following smoothness condition

$$f(\mathbf{w}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{w} - \mathbf{x} \rangle + \frac{1}{2\eta} \|\mathbf{w} - \mathbf{x}\|_2^2, \quad (6)$$

for some explicit  $\eta > 0$  depending on  $\mathbf{A}$ .

(c) **(5 points)** Prove that as a consequence, gradient descent with this fixed step size  $\eta$  converges to a critical point (hence a global minimizer) of (5). Recall that the gradient descent update can be written

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \eta \nabla f(\mathbf{x}^{(k)}). \quad (7)$$

(d) **(5 points)** Prove that if we initialize  $\mathbf{x}^{(0)} = \mathbf{0}$ , then gradient descent converges to the minimum  $\ell_2$ -norm solution

$$\hat{\mathbf{x}}_{\ell_2} \triangleq \arg \min_{\mathbf{x}} \|\mathbf{x}\|_2 \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{y}. \quad (8)$$

(e) **(10 points)** Now, suppose instead that we optimize a non-convex “factorized” variant of (5)

$$\underset{\mathbf{x}}{\text{minimize}} \frac{1}{2} \|\mathbf{A}(\mathbf{x} \odot \mathbf{x}) - \mathbf{y}\|_2^2, \quad (9)$$

where  $\odot$  denotes the entrywise product. Interestingly, Gunasekar and co-authors proved in [1, Corollary 2] that this seemingly small change in the optimization problem significantly affects the induced regularization. Specifically, they showed that gradient descent with very small step size and initialization very close to zero converges to the **minimum  $\ell_1$ -norm solution**

$$\hat{\mathbf{x}}_{\ell_1} \triangleq \arg \min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{y}. \quad (10)$$

Validate this result experimentally. Specifically, using MATLAB or Python, do the following.

- i) Generate  $\mathbf{A} \in \mathbb{R}^{20 \times 40}$  by sampling entries from the standard Gaussian distribution, and normalizing the columns to have unit  $\ell_2$  norm.
- ii) Generate a sparse  $\mathbf{x}^* \in \mathbb{R}^{40}$  with 10 non-zero entries whose locations are sampled uniformly at random. Generate the non-zero entries of  $\mathbf{x}^*$  by sampling from a Gaussian distribution with standard deviation  $1/\sqrt{10}$ , and taking absolute value so that all entries are non-negative. Generate  $\mathbf{y} = \mathbf{A}\mathbf{x}^*$ .
- iii) Implement gradient descent with a fixed step size  $\eta$  for problem (9).
- iv) Initialize  $\mathbf{x}^{(0)}$  near  $\mathbf{0}$ , e.g.  $\mathbf{x}^{(0)} = 10^{-3}\mathbf{1}$ , and run gradient descent for 300,000 steps with a small step size, e.g.  $\eta = 10^{-3}$ . Record the iterate  $\mathbf{x}^{(k)}$  every 10 steps.
- v) Compute the full Lasso solution path, e.g. using Scikit-Learn<sup>2</sup>. Recall that the Lasso problem is

$$\underset{\mathbf{x}}{\text{minimize}} g_\lambda(\mathbf{x}) \triangleq \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1. \quad (11)$$

The solution path then is the continuous curve  $\lambda \mapsto \hat{\mathbf{x}}(\lambda) \triangleq \arg \min_{\mathbf{x}} g_\lambda(\mathbf{x})$ . Note that in particular, as  $\lambda \rightarrow 0$ , the solution path converges to the minimum  $\ell_1$ -norm solution (10).

- vi) Plot the Lasso solution path<sup>3</sup> by plotting each coordinate  $(\hat{\mathbf{x}}(\lambda))_i$  ( $i = 1, \dots, m$ ) as a function of  $-\log_{10}(\lambda)$ . Similarly, in a separate subplot, plot the gradient descent iterates with iteration on the  $x$ -axis and each coordinate  $\mathbf{x}_i^{(k)}$  ( $i = 1, \dots, m$ ) as a separate line. Use a unique color for each coordinate so that they can be easily distinguished, and corresponding coordinates can be easily compared between the Lasso and

<sup>2</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.lasso\\_path.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.lasso_path.html)

<sup>3</sup>[https://scikit-learn.org/stable/auto\\_examples/linear\\_model/plot\\_lasso\\_coordinate\\_descent\\_path.html#sphx-glr-auto-examples-linear-model-plot-lasso-coordinate-descent-path-py](https://scikit-learn.org/stable/auto_examples/linear_model/plot_lasso_coordinate_descent_path.html#sphx-glr-auto-examples-linear-model-plot-lasso-coordinate-descent-path-py)

gradient descent. How does the Lasso solution path compare with that of gradient descent? Do the two paths seem to converge to the same ultimate solution?

(f) **(Bonus)** In [1], Gunasekar and co-authors are more interested in the more general matrix completion case

$$\underset{\mathbf{U}}{\text{minimize}} \frac{1}{2} \sum_{(i,j) \in \Omega} ((\mathbf{U}\mathbf{U}^\top)_{ij} - \mathbf{Y}_{ij})^2, \quad (12)$$

where  $\mathbf{Y} \in \mathbb{R}^{D \times D}$  is a symmetric and positive semi-definite low-rank matrix,  $\Omega$  is a set of indices for the “observed entries”, and  $\mathbf{U} \in \mathbb{R}^{D \times D}$ . The authors’ main conjecture is that gradient descent with small step size and near zero initialization converges to the **minimum nuclear norm solution**. Design and carry out an experiment to test this claim.

## References

- [1] S. Gunasekar, B. E. Woodworth, S. Bhojanapalli, B. Neyshabur, and N. Srebro. Implicit regularization in matrix factorization. In *Advances in Neural Information Processing Systems*, pages 6151–6159, 2017. 2, 3