

# Some Notes on Deep CCA

Karl Stratos

## 1 CCA

Let  $\mathbf{X} \in \mathbb{R}^{d \times N}$  and  $\mathbf{Y} \in \mathbb{R}^{d' \times N}$  denote two views of  $N$  samples. Let  $\mu_{\mathbf{X}}$  and  $\mu_{\mathbf{Y}}$  denote the column means of  $\mathbf{X}$  and  $\mathbf{Y}$  and assume the following centered, scaled matrices

$$\bar{\mathbf{X}} = \frac{1}{\sqrt{N}} (\mathbf{X} - \mu_{\mathbf{X}}) \quad \bar{\mathbf{Y}} = \frac{1}{\sqrt{N}} (\mathbf{Y} - \mu_{\mathbf{Y}})$$

where matrix-vector subtraction is understood as columnwise. Then the cross-covariance and covariance matrices (omitting regularization) are estimated by

$$\mathbf{C}_{\mathbf{XY}} = \bar{\mathbf{X}} \bar{\mathbf{Y}}^\top \quad \mathbf{C}_{\mathbf{X}} = \bar{\mathbf{X}} \bar{\mathbf{X}}^\top \quad \mathbf{C}_{\mathbf{Y}} = \bar{\mathbf{Y}} \bar{\mathbf{Y}}^\top$$

Given  $m \leq \min(d, d')$ , CCA calculates  $m$ -dimensional representations of the two views  $\underline{\mathbf{X}} \in \mathbb{R}^{m \times N}$  and  $\underline{\mathbf{Y}} \in \mathbb{R}^{m \times N}$  by

$$\underline{\mathbf{X}} = \mathbf{A}^\top \bar{\mathbf{X}} \quad \underline{\mathbf{Y}} = \mathbf{B}^\top \bar{\mathbf{Y}}$$

where  $\mathbf{A} = \mathbf{C}_{\mathbf{X}}^{-1/2} U_m$  and  $\mathbf{B} = \mathbf{C}_{\mathbf{Y}}^{-1/2} V_m$ . The columns of  $U_m \in \mathbb{R}^{d \times m}$  and  $V_m \in \mathbb{R}^{d' \times m}$  are the left/right singular vectors of

$$\boldsymbol{\Omega} = \mathbf{C}_{\mathbf{X}}^{-1/2} \mathbf{C}_{\mathbf{XY}} \mathbf{C}_{\mathbf{Y}}^{-1/2}$$

corresponding to the  $m$  largest singular values  $1 \geq \sigma_1(\boldsymbol{\Omega}) \geq \dots \geq \sigma_m(\boldsymbol{\Omega})$ : each  $\sigma_i(\boldsymbol{\Omega})$  is precisely the  $i$ -th canonical correlation defined in the original formulation of CCA (i.e., as an iterative optimization problem). Check that these representations satisfy

$$\underline{\mathbf{X}} \underline{\mathbf{X}}^\top = \underline{\mathbf{Y}} \underline{\mathbf{Y}}^\top = I_m \quad \underline{\mathbf{X}} \underline{\mathbf{Y}}^\top = \text{diag}(\sigma_1(\boldsymbol{\Omega}) \dots \sigma_m(\boldsymbol{\Omega}))$$

Because of the variational characterization of singular vectors (Theorem A.2), we can equivalently state that  $U_m, V_m$  are a solution of

$$\max_{\substack{\hat{U} \in \mathbb{R}^{d \times m}: \hat{U}^\top \hat{U} = I_m \\ \hat{V} \in \mathbb{R}^{d' \times m}: \hat{V}^\top \hat{V} = I_m}} \text{tr}(\hat{U}^\top \boldsymbol{\Omega} \hat{V}) = \sum_{i=1}^m \sigma_i(\boldsymbol{\Omega})$$

such that the columns ordered corresponding to nonincreasing singular values. To see why the maximum is achieved at  $\sum_{i=1}^m \sigma_i(\boldsymbol{\Omega})$ , see the proof of the theorem.

## 2 Deep CCA

In deep CCA [1], we assume that data matrices are already transformed by differentiable parameters to the target dimension  $m$  before considering the CCA objective. To be concrete, assume parameters  $W_1 \in \mathbb{R}^{m \times d}$  and  $W_2 \in \mathbb{R}^{m \times d'}$  yielding  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{m \times N}$  by

$$\mathbf{X} = \tanh(W_1 \widetilde{\mathbf{X}}) \quad \mathbf{Y} = \tanh(W_2 \widetilde{\mathbf{Y}})$$

where in this case we change our notation and write  $\widetilde{\mathbf{X}} \in \mathbb{R}^{d \times N}$  and  $\widetilde{\mathbf{Y}} \in \mathbb{R}^{d' \times N}$  to denote the original data matrices. As before, let  $\mathbf{\Omega} \in \mathbb{R}^{m \times m}$  denote the correlation matrix resulting from these matrices (the operations involved in calculating  $\mathbf{\Omega}$  are all fully differentiable). Now we consider the CCA objective and note that the optimum

$$\max_{\substack{\widehat{U} \in \mathbb{R}^{d \times m}: \widehat{U}^\top \widehat{U} = I_m \\ \widehat{V} \in \mathbb{R}^{d' \times m}: \widehat{V}^\top \widehat{V} = I_m}} \text{tr}(\widehat{U}^\top \mathbf{\Omega} \widehat{V}) = \sum_{i=1}^m \sigma_i(\mathbf{\Omega}) = \|\mathbf{\Omega}\|_1$$

coincides with the nuclear norm of  $\mathbf{\Omega}$  because this is the sum of *all* singular values of  $\mathbf{\Omega}$  without truncation. A main reason to prefer the nuclear norm  $\|\mathbf{\Omega}\|_1 := \text{tr}((\mathbf{\Omega}^\top \mathbf{\Omega})^{1/2})$  for differentiability is that it is given by the trace operator which has many nice differentiable properties (e.g., see [here](#)). The matrix gradients

$$\frac{\partial \|\mathbf{\Omega}\|_1}{\partial \mathbf{X}} \in \mathbb{R}^{m \times N} \quad \frac{\partial \|\mathbf{\Omega}\|_1}{\partial \mathbf{Y}} \in \mathbb{R}^{m \times N}$$

are given by an SVD of  $\mathbf{\Omega}$ . So this can be viewed as a final node in the computation graph with  $\widetilde{\mathbf{X}}, \widetilde{\mathbf{Y}}$  as parents whose forward and backward passes are calculated by an SVD.

So all we're doing is updating parameters  $W_1$  and  $W_2$  by taking gradient steps on  $\|\mathbf{\Omega}\|_1$ , which is a function of these parameters. There is no explicit CCA calculation during training! Rather, they are trained to maximize the objective of an implicit trivial CCA (i.e., involving no dimensionality reduction) on top of the representations they induce. Thus to calculate correlations, we need to explicitly perform a trivial CCA after training to obtain actual CCA projection matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times m}$ . We would then calculate the sum of canonical correlations on the training data as follows:

1. Apply the trained  $W_1$  and  $W_2$  to obtain  $m$ -dimensional data  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{m \times N}$ .
2. Apply the CCA projection matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times m}$  and return  $\text{tr}(\mathbf{A}^\top \mathbf{X} \mathbf{Y}^\top \mathbf{B})$ .

One final note: the deep CCA objective does not decompose over samples and invalidates the correctness of stochastic gradient descent. This issue is examined more closely in follow-up works [4, 3].

## A Matrix Facts

### A.1 Trace

**Theorem A.1** (Von Neumann, 7.4.11 in [2]). Let  $A, B \in \mathbb{R}^{m \times n}$  where  $m \leq n$ . Let  $\sigma_1(A) \geq \dots \geq \sigma_m(A)$  denote the  $m$  largest singular values of  $A$  in nonincreasing order (likewise for  $B$ ). Then

$$\text{tr}(AB^\top) \leq \sum_{i=1}^m \sigma_i(A)\sigma_i(B)$$

### A.2 SVD

**Theorem A.2.** Let  $A \in \mathbb{R}^{m \times n}$  where  $m \leq n$ . For any  $k \leq m$ ,

$$(U_k, V_k) \in \underset{\substack{\hat{U} \in \mathbb{R}^{m \times k}: \hat{U}^\top \hat{U} = I_k \\ \hat{V} \in \mathbb{R}^{n \times k}: \hat{V}^\top \hat{V} = I_k}}{\text{arg max}} \text{tr}(\hat{U}^\top A \hat{V})$$

where the columns of  $U_k \in \mathbb{R}^{m \times k}$  and  $V_k \in \mathbb{R}^{n \times k}$  are the left and right singular vectors of  $A$  corresponding to the largest  $k$  singular values.

*Proof.* For any orthonormal  $\hat{U} \in \mathbb{R}^{m \times k}$  and  $\hat{V} \in \mathbb{R}^{n \times k}$ ,

$$\text{tr}(\hat{U}^\top A \hat{V}) = \text{tr}(A \hat{V} \hat{U}^\top) \leq \sum_{i=1}^m \sigma_i(A) \sigma_i(\hat{V} \hat{U}^\top) = \sum_{i=1}^k \sigma_i(A)$$

To see the last equality, note that  $\sigma_i(\hat{V} \hat{U}^\top) = \lambda_i(\hat{U} \hat{U}^\top)^{1/2}$  where  $\hat{U} \hat{U}^\top$  is the projection operator onto a  $k$ -dimensional subspace of  $\mathbb{R}^m$ . This upper bound is reached by taking  $\hat{U} = U_k$  and  $\hat{V} = V_k$ .  $\square$

## References

- [1] Andrew, G., Arora, R., Bilmes, J., and Livescu, K. (2013). Deep canonical correlation analysis. In *International Conference on Machine Learning*, pages 1247–1255.
- [2] Horn, R. A. and Johnson, C. R. (2013). *Matrix analysis*.
- [3] Wang, W., Arora, R., Livescu, K., and Bilmes, J. (2015a). On deep multi-view representation learning. In *International Conference on Machine Learning*, pages 1083–1092.
- [4] Wang, W., Arora, R., Livescu, K., and Srebro, N. (2015b). Stochastic optimization for deep cca via nonlinear orthogonal iterations. In *Communication, Control, and Computing (Allerton), 2015 53rd Annual Allerton Conference on*, pages 688–695. IEEE.