

1 Intro

2 Optimization Problems

2.1 Objective Functions

I would like to develop rigorous probabilistic models for each of these objective functions below. The place to start is probably Bach and Jordan (2005).

2.1.1 Classical Two-View CCA

$$\begin{aligned} & \text{minimize} \quad \|\mathbf{X}_1 \Phi_1 - \mathbf{X}_2 \Phi_2\|_F^2 \\ & \text{subject to} \quad \Phi_1^\top \mathbf{X}_1^\top \mathbf{X}_1 \Phi_1 = \mathbf{I}_k \\ & \quad \quad \quad \Phi_2^\top \mathbf{X}_2^\top \mathbf{X}_2 \Phi_2 = \mathbf{I}_k, \end{aligned}$$

where Φ_1 and Φ_2 are the optimization variables.

2.1.2 Complete Multiview Graph

$$\begin{aligned} & \text{minimize} \quad \sum_{i=1}^m \sum_{j=i+1}^m \|\mathbf{X}_i \Phi_i - \mathbf{X}_j \Phi_j\|_F^2 \\ & \text{subject to} \quad \Phi_i^\top \mathbf{X}_i^\top \mathbf{X}_i \Phi_i = \mathbf{I}_k \quad i = 1, \dots, m, \end{aligned}$$

where Φ_i are the optimization variables.

2.1.3 (Mean Field Variational?) Approximation of Complete Graph

$$\begin{aligned} & \text{minimize} \quad \sum_{i=1}^m \|\mathbf{X}_i \Phi_i - \Psi\|_F^2 \\ & \text{subject to} \quad \Phi_i^\top \mathbf{X}_i^\top \mathbf{X}_i \Phi_i = \mathbf{I}_k \quad i = 1, \dots, m, \end{aligned}$$

where Φ_i and Ψ are the optimization variables.

2.1.4 Arbitrary Dependency Graph

$$\begin{aligned} & \text{minimize} \quad \sum_{(i,j) \in E} \|\mathbf{X}_i \Phi_i - \mathbf{X}_j \Phi_j\|_F^2 \\ & \text{subject to} \quad \Phi_i^\top \mathbf{X}_i^\top \mathbf{X}_i \Phi_i = \mathbf{I}_k \quad i = 1, \dots, m, \end{aligned}$$

where E is the set of edges in the dependency graph between the views, and Φ_i are the optimization variables.

2.2 Constraints

The classical formulation of generalized eigenvalue problems including CCA is to have quadratic equality constraints $\mathbf{X}^\top \mathbf{A} \mathbf{X} = \mathbf{I}$ on the parameters \mathbf{X} for some quadratic form \mathbf{A} . The purpose of these constraints is to force the parameters to counteract/reflect the scale of the data so that the model is robust to each view having dramatically different scale.

2.2.1 Issues with Equality Constraints

For both the classical and projected gradient algorithms for estimating CCA parameters, to satisfy the equality constraints, we are forced to orthogonalize some basis either via Gram-Schmidt orthogonalization or matrix square-roots. For batch classical and projected gradient optimization techniques, this causes scaling issues for both computation and storage. For the stochastic projected gradient algorithm, the projections onto a non-convex feasible region (a quadratic surface) result in lack of convergence guarantees and very likely slower local convergence when this is attainable. Additionally, if we want to use the transformations as dimensionality reduction on unseen data, we might consider satisfaction of equality constraints to be overfitting.

2.2.2 Alternatives

If we want to improve convergence guarantees and reduce computation, but still account for and be robust to differences in scale between the two data views, it may be beneficial to penalized the parameters for straying from the equality constraints rather than requiring exact constraint satisfaction. In fact, when using the stochastic projected gradient algorithm, violating the equality constraints by a little bit is unavoidable because we use minibatches for projecting the parameters, but the feasible region is defined in terms of the full dataset.

If we define some penalty function $R(\mathbf{X}) = \|\mathbf{X}^\top \mathbf{A} \mathbf{X} - \mathbf{I}\|_F^2$ and use this penalty function as a proximal function or regularizer, it may result in more efficient computation and better convergence properties for parameter estimation and better generalization ability for the dimensionality reduction defined in terms of the parameters.

3 Algorithms

3.1 Stochastic m -View AppGrad

The algorithm in this section is derived from a similar algorithm for classical two-view CCA given in Ma et al. (2015).

Algorithm 1 Stochastic m -view, rank- k AppGrad

```

1: Input: batch size  $n$ , step sizes  $\eta_{1t}, \dots, \eta_{mt}$ , max iter  $T$ , num components  $k$ 
2: Output: Canonical bases  $(\Phi_1, \dots, \Phi_m)$ 
3: Initialization:  $t = 0$ ,  $\tilde{\Phi}_1^{(0)} \in \mathbb{R}^{p_1 \times k}, \dots, \tilde{\Phi}_m^{(0)} \in \mathbb{R}^{p_m \times k}$  each with entries sampled individually from  $\mathcal{N}(0, 1)$ ,
    $\Phi_j^{(0)} = \mathbf{0}$  for  $j = 1, \dots, m$ 
4: while  $t < T$  and not converged do
5:   Receive minibatch  $\mathbf{X}_1^{(t)} \in \mathbb{R}^{n \times p_1}, \dots, \mathbf{X}_m^{(t)} \in \mathbb{R}^{n \times p_m}$ 
6:   for  $i = 1, \dots, m$  do
7:      $\tilde{\Phi}_j^{(t+1)} \leftarrow \tilde{\Phi}_j^{(t)} - \eta_{jt} \left( \mathbf{X}_i^{(t)} \right)^\top \left( (m-1) \mathbf{X}_i^{(t)} \tilde{\Phi}_i^{(t)} - \sum_{j \neq i} \mathbf{X}_j^{(t)} \Phi_j^{(t)} \right)$ 
8:     SVD:  $\left( \tilde{\Phi}_i^{(t+1)} \right)^\top \left( \frac{1}{n} \mathbf{X}_i^\top \mathbf{X}_i \right) \tilde{\Phi}_i^{(t+1)} = \mathbf{U}_i^\top \mathbf{D}_i \mathbf{U}_i$ 
9:      $\Phi_i^{(t+1)} \leftarrow \tilde{\Phi}_i^{(t+1)} \mathbf{U}_i^\top \mathbf{D}_i^{-\frac{1}{2}} \mathbf{U}_i$ 
10:   end for
11:    $t \leftarrow t + 1$ 
12: end while
```

Minibatch Management: It should be noted that the minibatches are managed by fixed-length queues to deal with different sampling rates for different views, so the pseudocode for Algorithm 1 is not entirely faithful to my implementation. I would like to eventually give more rigorous treatment of this later, possibly using some Markov chain theory.

Parameter Updates: It should be noted that I am not making naïve gradient updates. This is overlooked in the pseudocode for clarity and relevance reasons. I am actually applying AdaGrad scaling with optional dual averaging and shrinkage-and-thresholding to the singular values of the parameter matrix and gradient matrix (Duchi et al., 2011).

3.1.1 Why gradients?

Gradient-based updates allow us to easily derive stochastic and mini-batch algorithms and make incremental updates where we have complete control over the scale and influence of the incremental updates. These abilities facilitate a filter-like perspective of CCA and similar problems, and since we can control the scale/influence of each incremental update, we can control the rate of change of each parameter as a function of many quantities, including the update rates of other parameters. We can use tools from online learning and online optimization to approach these problems efficiently and with strong theoretical guarantees.

3.2 GenELinK and CCALin

Algorithms 2 and 3 in this section are taken directly from Ge et al. (2016). Algorithms 4 and 5 are inspired by Ge et al. (2016).

Algorithm 2 GenELinK

```

1: Input:  $T, k$ , symmetric matrix  $\mathbf{A}$ , PSD matrix  $\mathbf{B}$ , subroutine  $\text{GS}_{\mathbf{B}}(\cdot)$  that performs Gram-Schmidt process, with
   inner product  $\langle \cdot, \cdot \rangle_{\mathbf{B}}$ 
2: Output: top  $k$  eigen-space  $\mathbf{W} \in \mathbb{R}^{d \times k}$ 
3:  $\tilde{\mathbf{W}}^{(0)} \leftarrow$  random  $d \times k$  matrix with each entry i.i.d. from  $\mathcal{N}(0, 1)$ 
4:  $\mathbf{W}^{(0)} \leftarrow \text{GS}_{\mathbf{B}}(\tilde{\mathbf{W}}^{(0)})$ 
5: for  $t = 0, \dots, T - 1$  do
6:    $\mathbf{\Gamma}^{(t)} \leftarrow \left( (\mathbf{W}^{(t)})^{\top} \mathbf{B} \mathbf{W}^{(t)} \right)^{-1} \left( (\mathbf{W}^{(t)})^{\top} \mathbf{A} \mathbf{W}^{(t)} \right)$ 
7:    $\tilde{\mathbf{W}}^{(t+1)} \leftarrow \arg \min_{\mathbf{W}} \text{tr} \left( \frac{1}{2} \mathbf{W}^{\top} \mathbf{B} \mathbf{W} - \mathbf{W}^{\top} \mathbf{A} \mathbf{W} \right)$ 
8:   {Use an optimization subroutine with initialization  $\mathbf{W}^{(t)} \mathbf{\Gamma}^{(t)}$ }
9:    $\mathbf{W}^{(t+1)} \leftarrow \text{GS}_{\mathbf{B}}(\tilde{\mathbf{W}}^{(t+1)})$ 
10: end for
11: return  $\mathbf{W}^{(T)}$ 

```

Algorithm 3 CCALin

```

1: Input:  $T, k$ , data matrix  $\mathbf{X} \in \mathbb{R}^{n \times d_1}$ ,  $\mathbf{Y} \in \mathbb{R}^{n \times d_2}$ , subroutine  $\text{GS}_{\mathbf{M}}(\cdot)$  that performs Gram-Schmidt process,
   with inner product  $\langle \cdot, \cdot \rangle_{\mathbf{M}}$  where  $\mathbf{M}$  is any PSD matrix
2: Output: top  $k$  canonical subspace  $\mathbf{W}_x \in \mathbb{R}^{d_1 \times k}$ ,  $\mathbf{W}_y \in \mathbb{R}^{d_2 \times k}$ 
3:  $\mathbf{S}_{xx} \leftarrow \mathbf{X}^{\top} \mathbf{X} / n$ ,  $\mathbf{S}_{yy} \leftarrow \mathbf{Y}^{\top} \mathbf{Y} / n$ ,  $\mathbf{S}_{xy} \leftarrow \mathbf{X}^{\top} \mathbf{Y} / n$ 
4:  $\mathbf{A} \leftarrow \begin{pmatrix} 0 & \mathbf{S}_{xy} \\ \mathbf{S}_{xy}^{\top} & 0 \end{pmatrix}$ ,  $\mathbf{B} \leftarrow \begin{pmatrix} \mathbf{S}_{xx} & 0 \\ 0 & \mathbf{S}_{yy} \end{pmatrix}$ 
5:  $\begin{pmatrix} \tilde{\mathbf{W}}_x \in \mathbb{R}^{d_1 \times k} \\ \tilde{\mathbf{W}}_y \in \mathbb{R}^{d_2 \times k} \end{pmatrix} \leftarrow \text{GenELinK}(\mathbf{A}, \mathbf{B}, T, k)$ 
6:  $\mathbf{U} \leftarrow 2k \times k$  random Gaussian matrix
7:  $\tilde{\mathbf{W}}_x \leftarrow \tilde{\mathbf{W}}_x \mathbf{U}$ ,  $\tilde{\mathbf{W}}_y \leftarrow \tilde{\mathbf{W}}_y \mathbf{U}$ 
8:  $\mathbf{W}_x \leftarrow \text{GS}_{\mathbf{S}_{xx}}(\tilde{\mathbf{W}}_x)$ ,  $\mathbf{W}_y \leftarrow \text{GS}_{\mathbf{S}_{yy}}(\tilde{\mathbf{W}}_y)$ 

```

3.2.1 GenELinK and CCALin Filters

For transforming the CCA algorithm with GenELinK subroutine into an adaptive filtering algorithm, we need to ensure that two criteria are satisfied:

- There is an updated set of canonical bases available for filtering unseen datapoints after each minibatch is processed.
- State from previous parameter updates are maintained across minibatches.

Algorithm 4 GenELinK Filter

```
1: Initialization:  $\tilde{\mathbf{W}}^{(0)} \leftarrow$  random  $d \times k$  matrix with each entry i.i.d. from  $\mathcal{N}(0, 1)$ ,  $t \leftarrow 0$ 
2: while True do
3:   Receive symmetric matrix  $\mathbf{A}^{(t)}$ , PSD matrix  $\mathbf{B}^{(t)}$ , subroutine  $\text{GS}_{\mathbf{B}^{(t)}}(\cdot)$  that performs Gram-Schmidt process,
   with inner product  $\langle \cdot, \cdot \rangle_{\mathbf{B}^{(t)}}$ 
4:    $\mathbf{W}^{(t)} \leftarrow \text{GS}_{\mathbf{B}^{(t)}}(\tilde{\mathbf{W}}^{(t)})$ 
5:    $\mathbf{\Gamma}^{(t)} \leftarrow \left( (\mathbf{W}^{(t)})^\top \mathbf{B}_t \mathbf{W}^{(t)} \right)^{-1} \left( (\mathbf{W}^{(t)})^\top \mathbf{A}_t \mathbf{W}^{(t)} \right)$ 
6:    $\tilde{\mathbf{W}}^{(t+1)} \leftarrow \arg \min_{\mathbf{W}} \text{tr} \left( \frac{1}{2} \mathbf{W}^\top \mathbf{B}^{(t)} \mathbf{W} - \mathbf{W}^\top \mathbf{A}^{(t)} \mathbf{W} \right)$ 
7:   {Use an optimization subroutine with initialization  $\mathbf{W}^{(t)} \mathbf{\Gamma}^{(t)}$ }
8:    $\mathbf{W}^{(t+1)} \leftarrow \text{GS}_{\mathbf{B}^{(t)}}(\tilde{\mathbf{W}}^{(t+1)})$ 
9:   yield  $\mathbf{W}^{(t+1)}$ 
10:   $t \leftarrow t + 1$ 
11: end while
```

Algorithm 5 CCALin Filter

```
1: Initialization:  $t \leftarrow 0$ ,  $\text{GF}(\cdot, \cdot, \cdot)$  an instance of GenELinK Filter
2: while True do
3:   Receive  $\mathbf{X}^{(t)} \in \mathbb{R}^{n \times d_1}$ ,  $\mathbf{Y}^{(t)} \in \mathbb{R}^{n \times d_2}$ 
4:    $\mathbf{S}_{xx}^{(t)} \leftarrow (\mathbf{X}^{(t)})^\top \mathbf{X}^{(t)} / n$ ,  $\mathbf{S}_{yy}^{(t)} \leftarrow (\mathbf{Y}^{(t)})^\top \mathbf{Y}^{(t)} / n$ ,  $\mathbf{S}_{xy}^{(t)} \leftarrow (\mathbf{X}^{(t)})^\top \mathbf{Y}^{(t)} / n$ 
5:    $\mathbf{A}^{(t)} \leftarrow \begin{pmatrix} 0 & \mathbf{S}_{xy}^{(t)} \\ (\mathbf{S}_{xy}^{(t)})^\top & 0 \end{pmatrix}$ ,  $\mathbf{B}^{(t)} \leftarrow \begin{pmatrix} \mathbf{S}_{xx}^{(t)} & 0 \\ 0 & \mathbf{S}_{yy}^{(t)} \end{pmatrix}$ 
6:    $\begin{pmatrix} \bar{\mathbf{W}}_x \in \mathbb{R}^{d_1 \times k} \\ \bar{\mathbf{W}}_y \in \mathbb{R}^{d_2 \times k} \end{pmatrix} \leftarrow \text{GF}(\mathbf{A}^{(t)}, \mathbf{B}^{(t)}, \text{GS}_{\mathbf{B}^{(t)}})$ 
7:    $\mathbf{U}^{(t)} \leftarrow 2k \times k$  random Gaussian matrix
8:    $\tilde{\mathbf{W}}_x^{(t)} \leftarrow \bar{\mathbf{W}}_x^{(t)} \mathbf{U}^{(t)}$ ,  $\tilde{\mathbf{W}}_y^{(t)} \leftarrow \bar{\mathbf{W}}_y^{(t)} \mathbf{U}^{(t)}$ 
9:    $\mathbf{W}_x^{(t)} \leftarrow \text{GS}_{\mathbf{S}_{xx}^{(t)}}(\tilde{\mathbf{W}}_x^{(t)})$ ,  $\mathbf{W}_y^{(t)} \leftarrow \text{GS}_{\mathbf{S}_{yy}^{(t)}}(\tilde{\mathbf{W}}_y^{(t)})$ 
10:  yield  $\mathbf{W}_x^{(t)}$ ,  $\mathbf{W}_y^{(t)}$ 
11:   $t \leftarrow t + 1$ 
12: end while
```

To satisfy these criteria, it will likely be necessary to make a call to the GenELinK subroutine or something like it for each minibatch. There are at least two clear options aside from what is described in Algorithm 4.

- Perform the full GenELinK subroutine on each minibatch and initialize each round’s GenELinK with the previous round’s canonical bases.
- For each minibatch, perform only a single stochastic gradient step on the quadratic program inside the GenELinK subroutine before giving the output back to the CCA algorithm

3.2.2 m -View CCALin

To derive the m -view variation of CCALin, it should be sufficient to receive m datasets $\mathbf{X}_1, \dots, \mathbf{X}_m$ (or minibatches for the filtering scenario) and change the generation of \mathbf{A} and \mathbf{B} to the following

$$\mathbf{A} \leftarrow \begin{pmatrix} 0 & \mathbf{S}_{x_1 x_2} & \cdots & \mathbf{S}_{x_1 x_m} \\ (\mathbf{S}_{x_1 x_2})^\top & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{S}_{x_{m-1} x_m} \\ (\mathbf{S}_{x_1 x_m})^\top & \cdots & (\mathbf{S}_{x_{m-1} x_m})^\top & 0 \end{pmatrix}$$

$$\mathbf{B} \leftarrow \begin{pmatrix} \mathbf{S}_{x_1} & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \mathbf{S}_{x_m} \end{pmatrix}$$

where $\mathbf{S}_{x_i x_j} \leftarrow \mathbf{X}_i^\top \mathbf{X}_j / n$ and $\mathbf{S}_{x_i} \leftarrow \mathbf{X}_i^\top \mathbf{X}_i / n$.

References

- Bach, F. R. and Jordan, M. I. (2005). A probabilistic interpretation of canonical correlation analysis.
- Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *The Journal of Machine Learning Research*, 12:2121–2159.
- Ge, R., Jin, C., Kakade, S. M., Netrapalli, P., and Sidford, A. (2016). Efficient algorithms for large-scale generalized eigenvector computation and canonical correlation analysis. *arXiv preprint arXiv:1604.03930*.
- Ma, Z., Lu, Y., and Foster, D. (2015). Finding linear structure in large datasets with scalable canonical correlation analysis. *arXiv preprint arXiv:1506.08170*.