

# A Appendix

## A.1 Proofs

*Proof of Proposition 1.* First, we use a property of elliptically contoured distributions [Cambanis et al., 1981, Corollary 5] to obtain

$$\begin{aligned}\mathbb{E}[\mathbf{C}^\top X \mid \mathbf{B}^\top X] &= \mathbf{a} + \text{Cov}(\mathbf{C}^\top X, \mathbf{B}^\top X) \text{Var}^{-1}(\mathbf{B}^\top X) [\mathbf{B}^\top X - \mathbb{E}(\mathbf{B}^\top X)] \\ &= \mathbf{a} + \mathbf{C}^\top \text{Var}(X) \mathbf{B} (\mathbf{B}^\top \text{Var}(X) \mathbf{B})^{-1} \mathbf{B}^\top (X - \mathbb{E}X)\end{aligned}$$

for some constant  $\mathbf{a}$ . From condition (1) and the law of total covariance,

$$\begin{aligned}\text{Cov}(S, \mathbf{C}^\top X) &= \text{Cov}(f_S(\mathbf{B}^\top X, \epsilon_S), \mathbf{C}^\top X) \\ &= \mathbb{E}[\text{Cov}(f_S(\mathbf{B}^\top X, \epsilon_S), \mathbf{C}^\top X \mid \mathbf{B}^\top X, \epsilon_S)] + \\ &\quad \text{Cov}[\mathbb{E}(f_S(\mathbf{B}^\top X, \epsilon_S) \mid \mathbf{B}^\top X, \epsilon_S), \mathbb{E}(\mathbf{C}^\top X \mid \mathbf{B}^\top X, \epsilon_S)] \\ &= \text{Cov}(f_S(\mathbf{B}^\top X, \epsilon_S), \mathbb{E}(\mathbf{C}^\top X \mid \mathbf{B}^\top X)) \\ &= \text{Cov}(f_S(\mathbf{B}^\top X, \epsilon_S), X) \mathbf{B} (\mathbf{B}^\top \text{Var}(X) \mathbf{B})^{-1} \mathbf{B}^\top \text{Var}(X) \mathbf{C}.\end{aligned}$$

Thus, we have  $\text{Cov}(S, \mathbf{C}^\top X) = \mathbf{0}$  if  $\mathbf{B}^\top \text{Var}(X) \mathbf{C} = \mathbf{0}$  which implies that the columns of  $\mathbf{C}$  lie in the nullspace of  $\text{Var}(X) \mathbf{B}$ .  $\square$

*Proof of Theorem 1.* We first show that the basis (8) of the classifier hypothesis RKHS is orthonormal, and then compute the canonical angles using the basis (8) and an orthonormal basis of  $\mathcal{F}$ . Denote by  $\xi_i := (\gamma_i - \rho_i \sigma_i) \mathbf{Q} \mathbf{M} \mathbf{\Lambda}^{-1/2} \mathbf{U}_i + \rho_i \mathbf{A} \mathbf{T} \mathbf{\Omega}^{-1/2} \mathbf{V}_i$  the  $i$ -th basis function in (8), we have

$$\begin{aligned}\langle \xi_i, \xi_j \rangle &= (\gamma_i - \rho_i \sigma_i)(\gamma_j - \rho_j \sigma_j) \mathbf{U}_i^\top \mathbf{U}_j + \rho_i \rho_j \mathbf{V}_i^\top \mathbf{V}_j + 2\rho_i \sigma_i (\gamma_i - \rho_i \sigma_i) \mathbb{1}_{i=j} \\ &= [(\gamma_i - \rho_i \sigma_i)^2 + \rho_i^2 + 2\rho_i \sigma_i (\gamma_i - \rho_i \sigma_i)] \mathbb{1}_{i=j} \\ &= \mathbb{1}_{i=j},\end{aligned}$$

where the first equality follows from the orthonormal basis (7). This shows that (8) is an orthonormal basis. Using the orthonormal basis  $\{\psi_i := \phi \mathbf{Q} \mathbf{M} \mathbf{\Lambda}_i^{-1/2}\}_{i=1}^r$  of  $\mathcal{F}$ , we can use the SVD to compute the canonical angles (see e.g., Algorithm 6.4.3 in [Golub and Van Loan, 2013]) as

$$[\xi_1, \dots, \xi_d]^\top [\psi_1, \dots, \psi_r] = \text{diag}(\gamma_i - \rho_i \sigma_i) \mathbf{U}^\top + \text{diag}(\rho_i) \mathbf{\Sigma} \mathbf{U}^\top = \mathbf{I}_d \text{diag}(\gamma_i) \mathbf{U}^\top. \quad (1)$$

Here,  $\text{diag}(d_i)$  denotes the diagonal matrix with diagonal elements  $d_i$ . Note that the last term in (1) is the (thin) SVD, and the singular values  $\gamma_i$  are the canonical angles between  $\mathcal{M}$  and  $\mathcal{F}$ . Finally, we relate the canonical angles to the operator norm in (9). Recall that the orthogonal projector can be expressed as the tensor product  $\mathcal{P}_{\mathcal{F}} = \sum_{i=1}^r \psi_i \otimes \psi_i$ , and  $\mathcal{P}_{\mathcal{F}} h = \sum_{i=1}^r \langle h, \psi_i \rangle \psi_i$ . We have

$$\begin{aligned}\|\mathcal{P}_{\mathcal{F}} - \mathcal{P}_{\mathcal{M}}\| &= \|(\mathcal{P}_{\mathcal{F}} + \mathcal{P}_{\mathcal{F}})(\mathcal{P}_{\mathcal{F}} - \mathcal{P}_{\mathcal{M}})(\mathcal{P}_{\mathcal{M}} + \mathcal{P}_{\mathcal{M}})\| \\ &= \|\mathcal{P}_{\mathcal{F}} \mathcal{P}_{\mathcal{M}} - \mathcal{P}_{\mathcal{F}} \mathcal{P}_{\mathcal{M}}\| \\ &= \sup_{h \in \mathcal{H}_{\kappa, n}: \|h\| \leq 1} (\|\mathcal{P}_{\mathcal{F}} \mathcal{P}_{\mathcal{M}} h\| + \|\mathcal{P}_{\mathcal{F}} \mathcal{P}_{\mathcal{M}} h\|),\end{aligned}$$

where  $\mathcal{F}$  and  $\mathcal{M}$  represent respectively the orthogonal complements of  $\mathcal{F}$  and  $\mathcal{M}$ . It can be shown that  $\mathcal{P}_{\mathcal{F}} \mathcal{P}_{\mathcal{M}}$  and  $\mathcal{P}_{\mathcal{F}} \mathcal{P}_{\mathcal{M}}$  have the same nonzero singular values which are the sines of the principal angles between  $\mathcal{F}$  and  $\mathcal{M}$  (see e.g., p.249 of Stewart, 2001). From (1), these principal angles are  $\arccos(\gamma_i)$ . Thus, we obtain  $\|\mathcal{P}_{\mathcal{F}} - \mathcal{P}_{\mathcal{M}}\| = \sqrt{1 - \min_i \gamma_i^2}$ . To obtain (10), one can simply apply the trigonometric identity of sines yielding

$$\|\mathcal{P}_{\mathcal{G}} - \mathcal{P}_{\mathcal{M}}\| = \gamma_{\min} \sqrt{1 - \sigma_{\min}^2} - \sigma_{\min} \sqrt{1 - \gamma_{\min}^2} = \max \left\{ 0, \epsilon \sqrt{1 - \sigma_{\min}^2} - \sigma_{\min} \sqrt{1 - \epsilon^2} \right\},$$

where we denote by  $\gamma_{\min} := \min_i \gamma_i$ .  $\square$

## A.2 Implementation

Algorithm 1 gives the Matlab-style pseudo-code for our approach which can handle multiple protected attributes. This algorithm use the SDR procedure described in Algorithm 2 to compute the desired model representation with a specified trade-off  $\epsilon$ .

---

**Algorithm 1:**  $E = \text{MBasis}(\mathbf{K}, \mathbf{y}, \mathbf{S}, m, d, \epsilon)$  — Compute the basis  $\phi E$  for  $\mathcal{M}$

---

- ```

[1] Initialize  $\mathbf{W} = []$ ,  $n$  with the number of rows of  $\mathbf{K}$ , as well as indices  $\text{pos} = (\mathbf{y} == 1)$  and
       $\text{neg} = (\mathbf{y} \neq 1)$ .
[2] foreach column  $s$  of  $\mathbf{S}$  do
[3]   if EqualizedOdds or EqualityOfOpportunity then
[4]     Set  $\mathbf{B} = \mathbf{0}_{n \times m}$  and update  $\mathbf{B}(\text{pos} :) = \text{SDR}(\mathbf{K}(\text{pos}, \text{pos}), s, m)$ .
[5]     Append basis  $\mathbf{B}$  to  $\mathbf{W}$ :  $\mathbf{W} = [\mathbf{W} \ \mathbf{B}]$ .
[6]     if EqualizedOdds then
[7]       Set  $\mathbf{B} = \mathbf{0}_{n \times m}$ , then  $\mathbf{B}(\text{neg}, :) = \text{SDR}(\mathbf{K}(\text{neg}, \text{neg}), s, m)$ .
[8]       Let  $\mathbf{W} = [\mathbf{W} \ \mathbf{B}]$ .
      end
[9]   else
      Compute  $\mathbf{B} = \text{SDR}(\mathbf{K}, s, m)$ , and update  $\mathbf{W} = [\mathbf{W} \ \mathbf{B}]$ .
    end
  end
[10] Predictive Subspace: Compute the predictive subspace as the SDR subspace  $\mathbf{A} = \text{SDR}(\mathbf{K}, \mathbf{y}, d)$ .
[11] Fair Subspace: Obtain  $\mathbf{K}'$  by subtracting the mean of each column of  $\mathbf{K}$ . Let  $\tilde{\mathbf{K}} = \mathbf{K}'^\top \mathbf{K}'$ , and
      use QR decomposition to compute  $\mathbf{Q}$  as the nullspace basis of the column space of  $\tilde{\mathbf{K}}\mathbf{W}$ .
[12] Perform the eigenvalue decompositions to obtain Equation (7), and then use Theorem 1 to compute
       $E$ .

```
- 

---

**Algorithm 2:**  $\mathbf{W} = \text{SDR}(\mathbf{K}, \mathbf{s}, m)$  — Compute the SDR subspace  $\phi\mathbf{W}$

---

- ```

[1] Sort  $\mathbf{s}$  such that  $s(\text{idx})$  is non-decreasing. Let  $\text{invIdx}$  be the inverse of  $\text{idx}$  satisfying
       $\text{idx}(\text{invIdx}) = 1 : n$ , where  $n$  is the number of rows of  $\mathbf{K}$ .
[2] Slice  $\mathbf{s}$  approximately evenly as described in § 3.1 such that entries with the same value are in the
      same partition. Denote by  $n_i$  the size of partition  $i$ .
[3] Initialize  $\eta = 10^{-4}$ , i.e., a small constant. Let  $\mathbf{K}' := \mathbf{K}(\text{idx}, \text{idx})$ , and solve
       $\mathbf{\Gamma}_n \mathbf{K}' \mathbf{A}_i = \tau_i [\text{diag}(\mathbf{\Gamma}_{n_i}) \mathbf{K}' + n\eta \mathbf{I}_n] \mathbf{A}_i$  for  $\mathbf{A}$ .
[4] Return  $\mathbf{W} = \mathbf{A}(\text{invIdx}, :)$ .

```
- 

## References

- Stamatis Cambanis, Steel Huang, and Gordon Simons. On the theory of elliptically contoured distributions. *Journal of Multivariate Analysis*, 11(3):368–385, 1981.
- G.H. Golub and C.F. Van Loan. *Matrix Computations*. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, 2013. ISBN 9781421407944.
- G.W. Stewart. *Matrix Algorithms Volume 2: Eigensystems*. Other Titles in Applied Mathematics. Society for Industrial and Applied Mathematics, 2001. ISBN 9780898715033.