# Probabilistic Contrastive Principal Component Analysis

Didong Li[1,2*], Andrew Jones[1*], and Barbara E. Engelhardt[1,3]

*Department of Computer Science, Princeton University*[1]
*Department of Biostatistics, University of California, Los Angeles*[2]
*Center for Statistics and Machine Learning, Princeton University*[3]

May 4, 2021

## Abstract

Dimension reduction is useful for exploratory data analysis. In many applications, it is of interest to discover variation that is enriched in a "foreground" dataset relative to a "background" dataset. Recently, contrastive principal component analysis (CPCA) was proposed for this setting. However, the lack of a formal probabilistic model makes it difficult to reason about CPCA and to tune its hyperparameter. In this work, we propose probabilistic contrastive principal component analysis (PCPCA), a model-based alternative to CPCA. We discuss how to set the hyperparameter in theory and in practice, and we show several of PCPCA's advantages over CPCA, including greater interpretability, uncertainty quantification and principled inference, robustness to noise and missing data, and the ability to generate data from the model. We demonstrate PCPCA's performance through a series of simulations and case-control experiments with datasets of gene expression, protein expression, and images.

## 1 Introduction

Principal component analysis (PCA) is a popular technique for dimension reduction and data visualization (Hotelling, 1933). PCA has been widely used to understand the low-dimensional structure of datasets in a variety of scientific applications (Jirsa et al., 1994; Brenner et al., 2000; Novembre and Stephens, 2008; Darbyshire and Hamish, 2016; Pasini, 2017). In addition to its practical utility in data exploration tasks, estimation in PCA is computationally feasible using, for example, singular value decomposition (SVD). Moreover, PCA offers a satisfying geometric interpretation, namely, that the PCs capture orthogonal directions of maximum variation in the data. There is an immense literature on non-linear

---

[*]Equal contribution

Code for the model and experiments is available at `https://github.com/andrewcharlesjones/pcpca`.

generalizations of PCA including kernel PCA (Schölkopf et al., 1998), generalized PCA (Vidal et al., 2005), and principal curves (Hastie and Stuetzle, 1989), as well as modifications to PCA that incorporate sparsity (Tibshirani, 1996; Zou and Hastie, 2005; Zou et al., 2006), robustness (Candès et al., 2011), and more. In addition, probabilistic PCA (PPCA, Roweis 1998; Tipping and Bishop 1999) was developed to provide a model-based alternative to PCA, where the traditional objective function is re-interpreted as the likelihood estimate of a latent variable model that is a special homoskedastic version of Gaussian factor analysis (Fruchter, 1954). A non-linear version of probabilistic PCA was described soon afterwards in a Gaussian process latent variable model (GPLVM, Lawrence 2003).

However vast, these collective PCA methods are still not suitable for some applications. In this work, we consider settings in which the dataset consists of two groups — a *foreground group* and a *background group* — and we are interested in identifying structure, variation, and information unique to the foreground group. This situation arises naturally in many scientific experiments with two or more subpopulations, such as case-control studies. For example, in a genomics context, the foreground data could be gene expression measurements from patients with a disease, and the background data could be measurements from healthy patients (Twine et al., 2011; Zheng et al., 2017; Young et al., 2018). In this case, the goal is to identify transcriptional structure that is enriched in patients with the disease relative to healthy patients. Clearly, PCA is not suitable in this contrastive setting because PCA only identifies structure that exists across the union of the two groups or structure in each group in isolation.

Contrastive modeling approaches have recently been proposed for this purpose. As a first push in this direction, a general contrastive learning framework was developed for mixture models (Zou et al., 2013). More recently, contrastive PCA (CPCA) was developed (Abid et al., 2017, 2018) to find contrastive principal components (CPCs) that maximize variance in the foreground and minimize the variance in the background. However, in its original formulation, CPCA lacks a formal probabilistic model, so it is difficult to perform statistical inference within this framework. Moreover, the current CPCA framework does not allow a geometric interpretation.

In this paper, we develop probabilistic contrastive principal component analysis (PCPCA), a model-based alternative to CPCA for contrastive variation estimation. We recast the CPCA objective in a way that is amenable to a geometric interpretation, and we extend this analysis to the probabilistic setting. We then present a novel contrastive objective function which takes the form of a relative likelihood, and we provide a simple maximum relative likelihood estimate (MRLE) for the model. Furthermore, we develop a gradient descent algorithm that optimizes the objective in the presence of missing data.

We show that PCPCA is a more general model than PCA, PPCA, or CPCA, and that these three methods can be recovered as special cases of PCPCA, thus providing a unifying framework to understand these methods. Unlike CPCA, our model is both generative, providing a model-based approach that allows for uncertainty quantification and principled inference. Unlike PPCA, our proposed method extracts variation that is unique to the

2

foreground data while excluding variation shared between the foreground and background data, which is a critical goal in many experimental settings.

PCPCA may be applied to a variety of statistical and machine learning problem domains including dimension reduction, synthetic data generation, missing data imputation, and clustering. We demonstrate the model's behavior and capabilities through an extensive series of simulations and experiments with datasets of case/control gene and protein expression, and biological image data.

The specific contributions of our work to this field of PCA-based methods are the following. First, we present probabilistic contrastive component analysis (PCPCA), a model-based alternative to CPCA. Next, we show that three existing dimension reduction methods — PCA, PPCA, and CPCA — are special cases of our model. Then, we demonstrate several advantages of PCPCA, including a more principled probability model, a geometric interpretation analogous to that of PCA, a generalized inference procedure, robustness to missing data, and the ability to generate data from the model. Finally, we provide theoretical insight into the tuning parameter $\gamma$ in both CPCA and PCPCA, which controls the degree to which the model focuses on variation in the background or foreground data.

This paper is organized as follows. First, we review related dimension reduction methods, including PCA, PPCA, and CPCA. Second, we provide a novel geometric interpretation of CPCA, along with conditions under which CPCA is well-defined. Third, we present PCPCA, derive its maximum likelihood estimators, and show that PCA, PPCA, and CPCA are special cases of this model. Fourth, we present a generalized Bayes approach for performing posterior inference. Fifth, we present a gradient descent algorithm for fitting our model in the presence of latent variables or missing data. Finally, we demonstrate our model's performance through a series of experiments with simulated, biomedical, and image data. Proofs are in the Appendix.

## 2 Background

### 2.1 Principal Component Analysis (PCA)

Let $x_1, \cdots, x_n \in \mathbb{R}^D$ be i.i.d. observations and $X \in \mathbb{R}^{n \times D}$ with the $i$th row $x_i^\top$. PCA is designed to find the best $d$-dimensional affine subspace to represent the data, where $d < D$. There are several equivalent definitions of PCA. We review two of them below.

The first definition is derived from a geometric perspective, where PCA finds a hyperplane $V$ that minimizes the distance between the samples and this hyperplane:

$$\underset{V^\top V = \mathrm{I}_d}{\mathrm{argmin}} \sum_{i=1}^n d^2(x_i, V) = \underset{V^\top V = \mathrm{I}_d}{\mathrm{argmin}} \sum_{i=1}^n \|x_i - VV^\top x_i\|^2 = \underset{V^\top V = \mathrm{I}_d}{\mathrm{argmin}} \frac{1}{n} \sum_{i=1}^n \|x_i - VV^\top x_i\|^2, \quad (1)$$

where $V \in \mathbb{R}^{D \times d}$ has orthonormal column(s) , representing a $d$-dimensional subspace of $\mathbb{R}^D$.

The solution is given by

$$V = [v_1, \cdots, v_d], \ v_j = \text{eig}_j \left( \sum_{i=1}^n x_i x_i^\top \right) = \text{eig}_j(C),$$

where $\text{eig}_j$ is the $j$th eigenvalue of $C = \sum_{i=1}^n x_i x_i^\top$ in the descending order.

The second definition, which leads to an equivalent solution as Equation (1), is motivated from a statistical perspective. In particular, PCA maximizes the variance of the projected data onto each principal component, subject to the components being orthogonal to one another (assume $d = 2$ for simplicity):

$$\max_{v_1^\top v_1 = 1} \text{var}(v_1^\top x_i) = \max_{v_1^\top v_1 = 1} \sum_{i=1}^n v_1^\top x_i x_i^\top v_1 = \max_{v_1^\top v_1 = 1} v_1^\top C v_1, \tag{2}$$

$$\max_{v_2^\top v_2 = 1, v_1^\top v_2 = 0} \text{var}(v_2^\top (x_i - v_1 v_1^\top x_i)) = \max_{v_2^\top v_2 = 1, v_1^\top v_2 = 0} v_2^\top C v_2. \tag{3}$$

The geometric and statistical frameworks for PCA yield equivalent solutions, but having multiple perspectives gives greater insight into the method. Our work is motivated by these complementary perspectives (Theorem 1).

Note that we drop the mean parameter since, in practice, the data can easily be centered to have zero mean. Thus, for simplicity, throughout this paper we assume all data include features that are centered at zero.

## 2.2 Probabilistic PCA (PPCA)

PCA may be generalized in the form of a probabilistic model. Assume $z \sim N(0, \text{I}_d)$, $x = Wz + \epsilon$ with $W \in \mathbb{R}^{D \times d}$, $\epsilon \sim N(0, \sigma^2 \text{I}_D)$. Then

$$x \sim N(0, WW^\top + \sigma^2 \text{I}_D).$$

The objective is to maximize the likelihood with respect to the parameters $W$ and $\sigma^2$:

$$\underset{W, \sigma^2}{\text{argmax}} \ p(X|W, \sigma^2). \tag{4}$$

The MLE of $W$ and $\sigma^2$ are given by (Roweis, 1998; Tipping and Bishop, 1999):

$$\widehat{W}_{ML} = U(\Lambda - \widehat{\sigma}_{ML}^2 I_d)^{1/2} R, \ \ \widehat{\sigma}_{ML}^2 = \frac{1}{D-d} \sum_{i=d+1}^D \lambda_i,$$

where $U$ consists of the first $d$ eigenvectors of $C = \sum_{i=1}^n x_i x_i^\top$ with eigenvalues $\lambda_1 \geq \cdots \geq \lambda_D > 0$, $\Lambda = \text{diag}\{\lambda_1, \cdots, \lambda_d\}$ and $R \in \text{O}(d)$ is any rotation matrix.

That is, the hyperplane obtained by PPCA only differs by a re-scaling of the basis from the PCA hyperplane. In other words, PPCA is "equivalent" to PCA, and this becomes exact when $\sigma^2 \to 0$.

## 2.3   Contrastive PCA (CPCA)

PCA can also be generalized for contrastive modeling of two datasets. For foreground observations $x_1, \cdots, x_n \in \mathbb{R}^D$ and background observations $y_1, \cdots, y_m \in \mathbb{R}^D$, contrastive PCA (CPCA, Abid et al. 2018) is designed to discover low-dimensional structure that is unique to or enriched in the foreground dataset $X$ relative to the background dataset $Y$. Let $C_X = \frac{1}{n} \sum_{i=1}^{n} x_i x_i^\top$ be the empirical covariance matrix for $X$ and $C_Y = \frac{1}{m} \sum_{j=1}^{m} y_j y_j^\top$ for $Y$.

Recall the statistical perspective of PCA given by Equation 2. In the contrastive setting, for any unit vector $v$, we have two variances — the foreground variance and background variance — given by $v^\top C_X v$ and $v^\top C_Y v$, respectively. The objective of CPCA is to identify directions $v$ that account for a large amount of variance in the foreground and a small amount of variance in the background. Specifically, CPCA solves the following optimization problem:

$$\underset{v^\top v = 1}{\operatorname{argmax}} \; v^\top C_X v - \gamma v^\top C_Y v = \underset{v^\top v = 1}{\operatorname{argmax}} \; v^\top C v, \tag{5}$$

where $\gamma \in [0, \infty]$ is a tuning parameter, and $C = C_X - \gamma C_Y$.

The solution of CPCA is the same as PCA if we replace $C_X$ by $C$, namely, the optimal $v$ is the top eigenvector of $C$. From the definition of CPCA, it is clear that CPCA reduces to PCA when $\gamma = 0$.

# 3   A deeper look at CPCA

In this section, we analyze some important aspects of CPCA that were not discussed in previous studies. These analyses provide insight into when CPCA is well-defined, and in turn provide motivation for our proposed model, PCPCA, which is described in the next section.

## 3.1   Geometric interpretation of CPCA

CPCA was originally defined from a statistical perspective (Equation 5, Abid et al. 2018). Recalling the geometric definition of PCA, Equation 1, it is natural to consider whether there also exists a geometric interpretation for CPCA.

**Theorem 1.** *The statistical objective function of CPCA in Equation* (5) *is equivalent to the following geometric objective function*

$$\underset{v^\top v}{\operatorname{argmin}} \; \frac{1}{n} \sum_{i=1}^{n} \|x_i - vv^\top x_i\|^2 - \gamma \frac{1}{m} \sum_{j=1}^{m} \|y_j - vv^\top y_j\|^2. \tag{6}$$

The proof can be found in Appendix 9.2. From a geometric perspective, the objective of CPCA is to find a hyperplane that is close to the foreground data but far from the

background data. This coincides with the intuition of CPCA's overall goal, which is to identify the unique information in the foreground data.

In addition to the geometric intuition provided by Theorem 1, this theorem also allows us to adapt distance-based algorithms to the contrastive setting. Specifically, in any of these algorithms, one can consider replacing the distance by the "constrastive distance." For example, sparse CPCA has been developed following this philosophy (Boileau et al., 2020), although without this justification. However, it is important to note that the contrastive distance is not a well-defined distance, which may violate the assumptions of traditional distance-based algorithms, and so cannot be used to replace distance metrics in existing algorithms without some luck.

## 3.2   Positive definiteness of $C$

Another crucial consideration of CPCA is the positive definiteness of $C$, which may be treated as a "covariance matrix." However, $C$ is not necessarily positive definite unless $\gamma = 0$, in which case $C = C_X$.

Here, we derive a sufficient condition on $\gamma$ such that $C$ is positive definite. Let the eigenvalues of $C$, $C_X$, and $C_Y$ be $\lambda_1 \geq \cdots \lambda_D \geq 0$, $\mu_1 \geq \cdots \mu_D \geq 0$ and $\rho_1 \geq \cdots \geq \rho_D \geq 0$, respectively.

**Lemma 1.** *$C$ is positive definite if*

$$\gamma < \frac{\min\{\mu_1, \cdots, \mu_D\}}{\max\{\rho_1, \cdots, \rho_D\}}.$$

The proof can be found in Appendix 9.3. However, in CPCA, the positive definiteness of $C$ is not strictly required since the target is a $d \ll D$ dimensional subspace, and $D$ is large for high-dimensional data, such as biomedical data. Instead, CPCA only requires that the first $d$ eigenvalues of $C$ must be positive. In many applications for visualization and clustering, $d = 2$ (Abid et al., 2018), which allows $\gamma$ to be defined over a wide range.

The following theorem provides a necessary and sufficient condition for the first $d$ eigenvalues of $C$ being positive, with Lemma 1 as a special case when $d = D$.

**Theorem 2.** *The first $d$ eigenvalues of $C$ are positive if*

$$\gamma < \max\left\{\frac{\mu_d}{\rho_1}, \frac{\mu_{d+1}}{\rho_2}, \cdots, \frac{\mu_D}{\rho_{D-d+1}}\right\}.$$

*Otherwise, there exists a $C$ such that the dth eigenvalue is negative. That is, the upper bound is tight.*

The proof can be found in Appendix 9.4.

**Corollary 1.** *For a fixed $d$, a larger $\gamma$ corresponds to a smaller loss. For a fixed $\gamma$, the loss will decrease when $d$ is increased to $d + 1$ if $\gamma < \max\left\{\frac{\mu_{d+1}}{\rho_1}, \frac{\mu_{d+2}}{\rho_2}, \cdots, \frac{\mu_D}{\rho_{D-d}}\right\}$.*

6

The proof can be found in Appendix 9.5. The above corollary explains why, when $\gamma$ is large, a smaller $d$ is preferable: when $\gamma$ is large enough (such that $\lambda_3 < 0$), CPCA with $d = 2$ is better than higher dimensional CPCA in terms of mean squared error (MSE). See Section 6 for more details.

## 3.3   The tuning parameter $\gamma$

In CPCA, the tuning parameter $\gamma$ can be any non-negative real number, making it difficult to tune. Although a tuning method was suggested in the original CPCA proposal (Abid et al., 2018), the procedure depends on an almost exhaustive search, making it inefficient. We first analyze the role of $\gamma$ and propose a new parameterization such that the new tuning parameter $\gamma$ always lies in a small range, typically close to $[0, 1]$, making it easier to tune.

Recall that, for PCA, minimizing the sum of squared error and minimizing the mean squared error are equivalent, since the scale $\frac{1}{n}$ only changes the eigenvalues of the sample covariance, not its eigenvectors:

$$\underset{v^\top v=1}{\operatorname{argmin}} \sum_{i=1}^{n} \|x_i - vv^\top x_i\|^2 = \underset{v^\top v=1}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^{n} \|x_i - vv^\top x_i\|^2.$$

However, in the contrastive setting, the scale matters. Specifically, the following two optimization problems are not equivalent unless $m = n$, which rarely happens in practice:

$$\underset{v^\top v=1}{\operatorname{argmin}} \sum_{i=1}^{n} \|x_i - vv^\top x_i\|^2 - \gamma' \sum_{i=1}^{m} \|y_i - vv^\top y_i\|^2 \tag{7}$$

$$\neq \underset{v^\top v=1}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^{n} \|x_i - vv^\top x_i\|^2 - \frac{1}{m}\gamma' \sum_{i=1}^{m} \|y_i - vv^\top y_i\|^2.$$

Comparing Equations (6) and (7), we conclude that they are equivalent when $\gamma' = \gamma\frac{n}{m}$. If the sample sizes of the two groups are not the same, then the choice of $\gamma$ is different from the choice of $\gamma'$. We will show that this reparameterization, which is adjusted by the relative sample size, makes $\gamma$ more interpretable and easier to tune.

# 4   Probabilistic CPCA (PCPCA)

In this section, we present a probabilistic approach to contrastive learning. First, we present adjacent work on contrastive learning performed through probabilistic modeling. Then, we present our model, PCPCA, and analyze it through the lens of CPCA, PPCA, and PCA.
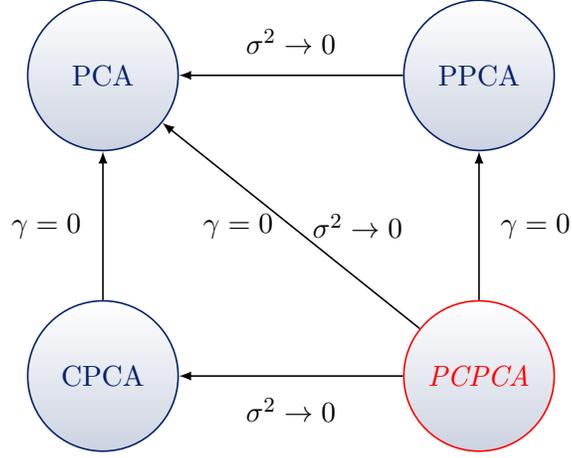
Figure 1: Target commutative diagram for PCA family.

## 4.1 Contrastive latent variable model

As a slightly different model than CPCA, the contrastive latent variable model (CLVM) has been proposed (Severson et al., 2019):

$$z \sim N(0, \mathrm{I}_d), \ t \sim N(0, \mathrm{I}_{d'}), \ x = Sz + Wt + \epsilon_x, \ y = Sz + \epsilon_y, \ \epsilon_x, \epsilon_y \sim N(0, \sigma^2 \mathrm{I}_D). \tag{8}$$

The marginals are given by

$$x \sim N(0, SS^\top + WW^\top + \sigma^2 \mathrm{I}_D), \ \ y \sim N(0, SS^\top + \sigma^2 \mathrm{I}_D).$$

The objective function in inference is the likelihood or log likelihood. When the background dimension is zero, that is, when $d' = 0$, the above model becomes PPCA for $X \cup Y$.

However, we are interested in characterizing the foreground data $X$ while controlling for variation in the background $Y$. For this reason, it is more desirable to recover PPCA for $X$ as a special case of the model rather than PPCA for $X \cup Y$ or $Y$. Recall that similar a relation holds for PCA and CPCA: when $\gamma = 0$, CPCA becomes PCA on $X$ only. In fact, there does not exist any $\gamma$ such that CPCA is equivalent to PCA on $X \cup Y$.

In addition, there is no clear link between the CLVM (Equation 8) and CPCA, even if $\sigma^2 \to 0$. As a result, it is of interest to develop a general model that is simultaneously a probabilistic version of CPCA and a contrastive version of PPCA.

## 4.2 PCPCA model

Consider the following model

$$z_x, z_y \sim N(0, \mathrm{I}_d), \ x = Wz_x + \epsilon_x, \ y = Wz_y + \epsilon_y, \tag{9}$$

8

where $\epsilon_x$, $\epsilon_y \sim N(0, \sigma^2 \, \mathrm{I}_D)$ are i.i.d. Gaussian noise vectors. Recall the equivalent statistical and geometric interpretations of PPCA: maximizing the likelihood of $X$ is equivalent to minimizing the distance between $X$ and the hyperplane $W$. For CPCA, we expect such $W$ to be far away from the background $Y$, which is exactly the (geometric) objective of CPCA (see Theorem 1). In the probabilistic setting, maximizing the distance from $Y$ is equivalent to minimizing the likelihood of $Y$. This is counter-intuitive, but it coincides with our model's motivation to account for variation in the foreground data, not the background data. Thus, we have the following objective function

$$\underset{W, \sigma^2}{\mathrm{argmax}} \; \frac{p(X|W, \sigma^2)}{p(Y|W, \sigma^2)^\gamma}. \tag{10}$$

The above objective function becomes the PPCA objective function when $\gamma = 0$, and a relative likelihood when $\gamma = 1$. For general $\gamma \in [0, \infty)$, we refer to Equation 10 as the relative likelihood, as it captures the likelihood of the foreground data with respect to the likelihood of the background data, scaled by $\gamma$.

The non-traditional nature of this objective requires further comment. Notice that this objective is not a traditional likelihood ratio, which is typically defined as a ratio of the likelihood of one dataset under two different parameter settings. Rather, ours is a ratio of likelihoods of two different datasets under a shared parameter setting. Furthermore, Equation (10) is not a well-defined likelihood unless $\gamma = 0$. These caveats preclude the use of traditional estimation and inference procedures based on likelihoods and relative likelihoods. For this reason, we present alternative procedures: one based on a direct maximization of Equation (10) and another based on generalized posterior inference.

First, we investigate the closed-form solution for this objective.

**Theorem 3.** *The $W$ and $\sigma^2$ that maximize Equation* (10), *denoted by* $\widehat{W}_{ML}, \widehat{\sigma}^2_{ML}$, *are given by*

$$\widehat{\sigma}^2_{ML} = \frac{1}{(n - \gamma m)(D - d)} \sum_{i=d+1}^{D} \lambda_i$$

$$\widehat{W}_{ML} = U_d \left( \frac{\Lambda_d}{n - \gamma m} - \widehat{\sigma}^2_{ML} \, \mathrm{I}_d \right)^{1/2} R,$$

*where $U_d$ consists of the first $d$ eigenvectors of $C = \sum_{i=1}^{n} x_i x_i^\top - \gamma \sum_{j=1}^{m} y_j y_j^\top$, $\Lambda_d = \mathrm{diag}\{\lambda_1, \cdots, \lambda_d\}$ contains the corresponding eigenvalues, and $R$ is any $d$ by $d$ rotation matrix. Moreover, $\widehat{W}_{ML}$ is equivalent to the $\widehat{W}_{PPCA}$ that maximizes the PPCA objective when $\gamma = 0$, and is equivalent to the $\widehat{W}_{CPCA}$ that maximizes the CPCA objective as $\sigma^2 \to 0$.*

The proof can be found in Appendix 9.6. As a result, we find last missing piece in the commutative diagram 1, namely, we have a complete algorithm for PCPCA that allows PCA, PPCA, and CPCA to be recovered as subcases of this general framework.

9

**Remark 1.** *The above PCPCA solution highlights two hidden assumptions for PCPCA:*

1. *$n - \gamma m > 0$ so that $\widehat{W}_{ML}$ is well defined, that is, $\gamma < \frac{n}{m}$.*

2. *$\sum_{i=d+1}^{D} \lambda_i > 0$ so that $\widehat{\sigma}_{ML}^2 > 0$. By Theorem 2, a sufficient condition is $\gamma < \frac{\sum_{i=d+1}^{D} \mu_i}{(D-d)\rho_1}$.*

*These seemingly strong constraints restrict the range of $\gamma$ from $[0, \infty)$ to a small interval, often a subinterval of $[0, 1]$. This more restricted interval makes PCPCA easier to tune than CPCA. In addition, the performance of PCPCA is robust to the choice of $\gamma$ within this interval, which is not observed for CPCA.*

## 5 Generalized Bayesian approach

Next, we present a generalized Bayesian framework for performing posterior inference in the PCPCA model. Recall that our objective (Equation 10) is not a likelihood, so we cannot simply place a prior on $\theta = (W, \sigma^2)$ and perform Bayesian inference in the traditional fashion. For this reason, we leverage more general loss-based inference methods based on Gibbs posteriors.

### 5.1 Gibbs Posterior

Let $\Theta$ be the space of parameters, which can be a finite or infinite dimensional space, and $U$ be the feature space, then we denote the loss function $l : U \times \Theta \to \mathbb{R}$ and $l_\theta : U \to \mathbb{R}$. For a given measure $P$ on $U$ (often the true measure), define the risk function $R : \Theta \to \mathbb{R}$ to be

$$R : \Theta \to \mathbb{R}, \ \theta \mapsto \mathbb{E}_P[l_\theta].$$

Then the goal is to minimize the risk: $\min_{\theta \in \Theta} R(\theta)$. It is common to assume the minimizer is unique, denoted by $\theta^* = \arg\min_{\theta \in \Theta} R(\theta)$.

However, $P$ is often unknown. Instead, we have observations $u_1, \cdots, u_n$ and we have the corresponding empirical measure $P_n = \frac{1}{n} \sum_{i=1}^{n} \delta_{u_i}$. Then the empirical risk function is

$$R_n : \Theta \to \mathbb{R}, \ \theta \mapsto \mathbb{E}_{P_n}[l_\theta].$$

The goal in this more practical setting is to minimize the empirical risk:

$$\min_{\theta \in \Theta} R_n(\theta).$$

Let $\widehat{\theta}_n = \arg\min_{\theta \in \Theta} R_n(\theta)$ be the unique minimizer of the empirical risk.

If a statistical model exists with density function $p_\theta$ and the loss function is $l_\theta(u) = -\log p_\theta(u)$, then this minimization reduces to maximum likelihood estimation. In particular, in this setting, $R_n(\theta)$ is the negative log likelihood and $\widehat{\theta}_n$ is the MLE.

**Definition 1.** *Given a prior $\Pi$ on $\Theta$, the Gibbs posterior is defined as*

$$\Pi_n(d\theta) \propto e^{-wnR_n(\theta)}\Pi(d\theta), \ \theta \in \Theta,$$

*where $\Pi$ is the prior and $w > 0$ is the learning rate.*

The Gibbs posterior becomes the true posterior if $R_n$ is the negative log-likelihood.

**Definition 2.** *The Gibbs posterior $\Pi_n$ asymptotically concentrates around $\theta^*$ at the rate $\varepsilon_n \to 0$ w.r.t. divergence measure $d$ on $\Theta$ if*

$$\mathbb{E}_{P_n}\Pi_n\left(\{\theta : d(\theta, \theta^*) > M\varepsilon_n\}\right) \to 0,$$

*where $M$ is a constant.*

Note that $d$ is only required to be positive semi-definite, that is, $d(\theta; \theta') \geq 0$ with equality iff $\theta = \theta'$.

## 5.2   Gibbs posterior for CPCA

We next propose a loss function and Gibbs posterior for CPCA. Let $u = (x, \alpha)$ where $x \in \mathbb{R}^D$ is the observation, $\alpha \in \{0, 1\}$ with $\alpha = 0$ represents foreground data while $\alpha = 1$ represents background data and $\theta = V \in \text{Gr}(D, d)$, the Grassmannian manifold. Let the loss function be

$$l_\theta(u) = (-\gamma)^\alpha d^2(x, V) = (-\gamma)^\alpha \|x - VV^\top x\|^2 = (-\gamma)^\alpha \left(x^\top x - xVV^\top x\right).$$

For simplicity, assume $u_1, \cdots, u_n$ are foreground data while $u_{n+1}, \cdots, u_{n+m}$ are background data. As a result, the empirical risk function is

$$
\begin{aligned}
R_n(\theta) &= \frac{1}{n+m} \sum_{i=1}^{n+m} l_\theta(u_i) \\
&= \frac{1}{n+m} \left( \sum_{i=1}^{n} \left(x_i^\top x_i - x_i^\top VV^\top x_i\right) - \gamma \sum_{j=1}^{m} \left(x_{n+j}^\top x_{n+j} - x_{n+j}^\top VV^\top x_{n+j}\right) \right) \\
&= -\sum_{i=1}^{n} x_i^\top VV^\top x_i + \gamma \sum_{j=1}^{m} x_{n+j}^\top VV^\top x_{n+j} + M \\
&= -\text{tr}(VV^\top C) + M,
\end{aligned}
$$

where $C = \sum_{i=1}^{n} x_i x_i^\top - \gamma \sum_{j=1}^{m} x_{n+j} x_{n+j}^\top$ and $M$ is independent of $\theta$. We conclude that the empirical risk function coincides with the objective function of CPCA. Furthermore, if the prior $\Pi$ is chosen to be the uniform prior, then the maximum a posteriori estimation (MAP) matches the solution of CPCA, which is the subspace spanned by the first $d$ eigenvectors of $C$.

For the population version of the risk, assume $P = \beta P_F + (1 - \beta) P_B$ where $\beta \in (0, 1)$, $P_F$ is the foreground measure with zero mean and covariance $C_F$, and $P_B$ is the background measure with zero mean and covariance $C_B$. Then the risk function is

$$
\begin{aligned}
R(\theta) &= \mathbb{E}_P l_\theta(u) \\
&= \beta \mathbb{E}_{x \sim P_F} \left( x^\top x - x^\top V V^\top x \right) - (1 - \beta) \gamma \mathbb{E}_{x \sim P_B} \left( x^\top x - x^\top V V^\top x \right) \\
&= -\beta \operatorname{tr}(V V^\top C_F) - (1 - \beta) \gamma \operatorname{tr}(V V^\top C_B) + M \\
&= -\operatorname{tr}(V V^\top C) + M,
\end{aligned}
$$

where $C = \beta C_F - (1 - \beta) \gamma C_B$ and $M$ is independent of $\theta$. So the minimizer is given by

$$
\theta^* = V^* = [\operatorname{eig}_1(C), \cdots, \operatorname{eig}_d(C)].
$$

We consider the risk divergence $d(\theta; \theta^*) = (R(\theta) - R(\theta^*))^{1/2} = \operatorname{tr}((V^* V^{*\top} - V V^\top) C)^{1/2}$, which measures the difference between risks. We now consider the contraction rate of this Gibbs posterior.

**Theorem 4.** *Assume $P = \beta N(0, C_F) + (1 - \beta) N(0, C_B)$ and let the prior $\Pi$ be uniform on $\operatorname{Gr}(D, d)$ w.r.t. the Haar measure, then the Gibbs posterior $\Pi_n$ asymptotically contracts to $\theta^*$ w.r.t. $d$ at rate $n^{-1/2}$.*

The proof can be found in Appendix 9.7. As a result, the Gibbs posterior will contract to the optimal parameter as the sample size increases, which provides theoretical support for the generalized Bayesian version of CPCA.

## 5.3   Gibbs Posterior for PCPCA

We now consider the Gibbs posterior for PCPCA. Let $u = (x, \alpha)$ where $x \in \mathbb{R}^D$ is the observation and $\alpha \in \{0, 1\}$ indicates the sample's condition, with $\alpha = 0$ representing foreground data while $\alpha = 1$ representing background data. As before, let $\theta = (W, \sigma^2)$ be the parameter. Let the loss function be $l_\theta(u) = -(-\gamma)^\alpha \log N(v; 0, W W^\top + \sigma^2 \mathrm{I}_D)$. For simplicity, assume $u_1, \cdots, u_n$ are foreground data while $u_{n+1}, \cdots, u_{n+m}$ are background data. As a result, the empirical risk function is

$$
\begin{aligned}
R_n(\theta) &= \frac{1}{n+m} \sum_{i=1}^{n+m} l_\theta(u_i) \\
&= \frac{1}{n+m} \left( \sum_{i=1}^{n} -\log N(x_i; 0, W W^\top + \sigma^2 \mathrm{I}_d) + \gamma \sum_{j=1}^{m} \log N(x_{n+j}; 0, W W^\top + \sigma^2 \mathrm{I}_d) \right) \\
&= \sum_{i=1}^{n} \left( \frac{1}{2} \log |A| + \frac{1}{2} x_i^\top A^{-1} x_i \right) + \gamma \sum_{j=1}^{m} \left( -\frac{1}{2} \log |A| - \frac{1}{2} x_{n+j}^\top A^{-1} x_{n+j} \right) + M \\
&= \frac{n - \gamma m}{2} \log |A| + \frac{1}{2} \operatorname{tr}(A^{-1} C) + M,
\end{aligned}
$$

12

where $A = WW^\top + \sigma^2 I_D$, $C = \sum_{i=1}^n x_i x_i^\top - \gamma \sum_{j=1}^m x_{n+j} x_{n+j}^\top$ and $M$ is independent of $\theta$. We conclude that the empirical risk function coincides with the (negative log) objective function of PCPCA. Furthermore, if the prior $\Pi$ is chosen to be the uniform prior, then the maximum a posteriori (MAP) estimate matches the solution in Theorem 3.

For the population version of the risk, assume the same model as in the previous section. Specifically, we assume $P = \beta P_F + (1-\beta)P_B$ where $\beta \in (0,1)$, $P_F$ is the foreground measure with zero mean and covariance $C_F$, and $P_B$ is the background measure with zero mean and covariance $C_B$. Then the risk function is

$$
\begin{aligned}
R(\theta) &= \mathbb{E}_P l_\theta(u) \\
&= -\beta \mathbb{E}_{x \sim P_F} \log p(x|0, A) + (1-\beta)\gamma \mathbb{E}_{x \sim P_B} \log p(x|0, A) \\
&= \frac{\beta}{2} \left( \log |A| + \mathrm{tr}(A^{-1} C_F) \right) - \frac{(1-\beta)}{2} \gamma \left( \log |A| + \mathrm{tr}(A^{-1} C_B) \right) + M \\
&= \frac{\beta - (1-\beta)\gamma}{2} \log |A| + \frac{1}{2} \mathrm{tr}(A^{-1} C) + M,
\end{aligned}
\tag{11}
$$

where $C = \beta C_F - (1-\beta)\gamma C_B$. So the minimizer is given by

$$
\sigma^{*2} = \frac{1}{D-d} \sum_{i=d+1}^D \lambda_i, \quad W^* = U_d \left( \frac{\Lambda_d}{\beta - (1-\beta)\gamma} - \sigma^{*2} I \right)^{1/2} R
\tag{12}
$$

where $\Lambda_d = \mathrm{diag}\{\lambda_i\}$ consists of the largest $d$ eigenvalues of $C$, and $U_d$ consists of the corresponding $d$ eigenvectors. We consider the same risk divergence as in the previous section:

$$
d(\theta; \theta^*) = (R(\theta) - R(\theta^*))^{1/2} = \left( \frac{\beta - (1-\beta)\gamma}{2} (\log |A| - \log |A^*|) + \frac{\mathrm{tr}((A^{-1} - A^{*-1})C)}{2} \right)^{1/2}.
$$

We now consider the contraction rate of the PCPCA Gibbs posterior under this divergence.

**Theorem 5.** *Assume $P = \beta N(0, C_F) + (1-\beta)N(0, C_B)$ and $\sigma^2 \geq \sigma_0^2 > 0$. Let the prior $\Pi$ be uniform on $\mathbb{R}^{D \times d} \times [\sigma_0^2, \infty)$, then the Gibbs posterior $\Pi_n$ asymptotically contracts to $\theta^*$ w.r.t. $d$ at rate $n^{-1/2}$.*

The proof can be found in Appendix 9.8. As a result, the Gibbs posterior contracts to the optimal parameter as the sample size increases, which supports the generalized Bayesian PCPCA.

# 6 Experiments

To demonstrate the behavior and usefulness of PCPCA, we fit the model on a series of simulated, gene and protein expression, and image datasets. Note that for most plots, we refer to the sample size-adjusted hyperparameter $\gamma' = \frac{m}{n}\gamma$.

## 6.1 Visualizing the role of the hyperparameter $\gamma$

First, to demonstrate the role of the hyperparameter $\gamma$ in the PCPCA model, we fit the PCPCA model on a two-dimensional simulated dataset. In this simple dataset, the foreground data contain two subgroups, each of which shares an axis of variation with the background data. In particular, we generated the foreground and background by sampling $x_i \sim \mathcal{N}(\mu_x, \Sigma)$ and $y_j \sim \mathcal{N}(0, \Sigma)$ where $\mu_x = \left(\begin{smallmatrix} 1 \\ -1 \end{smallmatrix}\right)$ for half of the foreground samples, and $\mu_x = \left(\begin{smallmatrix} -1 \\ 1 \end{smallmatrix}\right)$ for the other half. For all samples, $\Sigma = \left(\begin{smallmatrix} 2.7 & 2.6 \\ 2.6 & 2.7 \end{smallmatrix}\right)$. We set the foreground and background sample sizes to be equal, $n = m = 200$. We fit the PCPCA model for $\gamma' \in \{0, 0.2, 0.6, 0.9\}$, and we visualize the 1-dimensional line defined by $\widehat{W}_{ML}$, where $d = 1$ and $D = 2$ (Figure 2).

Recall that when $\gamma' = 0$, PCPCA reduces to PPCA. In this case, $W$ captures the variation that is shared between the background and foreground data (Figure 2a). At higher values of $\gamma'$, we observed that PCPCA captures the variation that is unique to the foreground dataset, which divides the two foreground subgroups (Figure 2d). Note that $\widehat{W}_{ML}$ rotates nearly 90 degrees to capture the direction of maximal variation unique to the foreground data when $\gamma' = 0.9$ relative to the PPCA solution when $\gamma' = 0$. At intermediate values of $\gamma'$, $\widehat{W}_{ML}$ balances between capturing the shared and foreground-specific variation (Figure 2c).
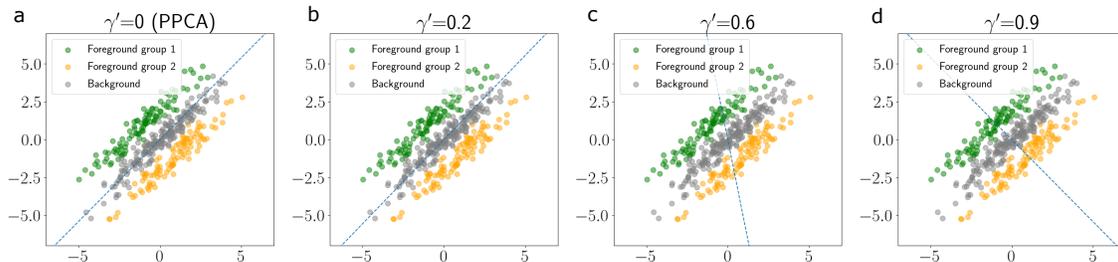


Figure 2: **PCPCA MLE on simulated data.** PCPCA estimates for toy data at varying values of $\gamma'$. The dotted line represents $\widehat{W}_{ML}$ ($d = 1$ in this case). PCPCA recovers PPCA when $\gamma' = 0$ (a) and captures the axis of variation between the two foreground subgroups when $\gamma' = 0.9$ (d).

## 6.2 Tuning $\gamma$ in experimental settings

### 6.2.1 Mouse protein expression

We next tested PCPCA using a dataset of mouse protein expression (Higuera et al., 2015). In this experiment, the foreground data are protein expression samples from the cortex of mice with and without Down Syndrome who were subjected to shock therapy. The background dataset consists of a set of protein expression measurements from mice without

Down Syndrome who did not receive shock therapy. In total, there are $n = 270$ foreground samples and $m = 135$ background samples, each measuring the expression of 77 proteins. The foreground samples contain 135 mice with Down Syndrome and 135 mice without Down Syndrome, and the intervention we model in this experiment is how shock therapy affects protein expression levels differently for mice with Down Syndrome and those without.

We fit PCPCA using a range of values for the tuning parameter $\gamma'$, setting $d = 2$ in each case. We found that, at higher values of $\gamma'$, PCPCA was able to separate the mice with and without Down Syndrome that received shock therapy (Figure 3b, c). Furthermore, PCPCA separated the background samples from the foreground samples (Figure 3d). When $\gamma' = 0$, the model is equivalent to PPCA, and visually there is minimal separation of the two groups of foreground mice (Figure 3a).
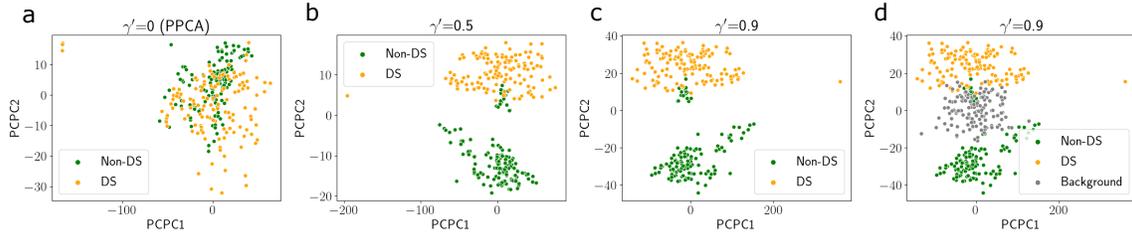


Figure 3: **PCPCA on mouse protein expression data.** PCPCA applied to a dataset of mouse protein expression measurements from shock therapy-treated mice with and without Down Syndrome. Plotted in each panel are the projections of the foreground samples onto the first two components at varying values for $\gamma'$. There are two subgroups in the foreground: *DS* (yellow) represents mice with Down Syndrome who are exposed to shock therapy, and *Non-DS* (green) represents mice without Down Syndrome who are exposed to shock therapy. (a)-(c) show the two foreground groups at each value of $\gamma'$. (d) also includes the *Background* dataset (gray), which is made up of mice without Down Syndrome who were not exposed to shock therapy.

We measured the degree of separation using the silhouette score (SS) of the two foreground groups of mice (Down syndrome and control) when projected into PCPCA's latent space. SS is a measure of cluster tightness (Rousseeuw, 1987); higher scores represent better clustering of sample labels in the space. We found that the maximum silhouette score achieved by CPCA and PCPCA were comparable (CPCA: 0.425, PCPCA: 0.404).

However, we observed different behavior between the methods in the tuning process for $\gamma'$. For PCPCA, we found that SS increased monotonically with $\gamma'$ (Figure 4b). In contrast, CPCA showed better clustering performance at lower values of $\gamma'$, and the SS decreased with a higher $\gamma'$ (Figure 4a). Additionally, the range of allowable values for $\gamma'$ differed substantially between the two methods. The looser constraint on $\gamma$ in CPCA allowed for high values of $\gamma'$ — going as high as $\gamma' = 241$ in the mouse dataset. The reason for the large

allowable values of $\gamma$ in CPCA can be understood in the context of Corollary 1. Furthermore, at these large values of $\gamma'$, the CPCA projection of the background dataset reduces to a single point (Figure 4c). Together, these results suggest that PCPCA's parameterization allows for an easier interpretation of the tuning parameter $\gamma'$, and $\gamma'$ is restricted to a reasonable range in PCPCA compared with the parameter's range in CPCA.
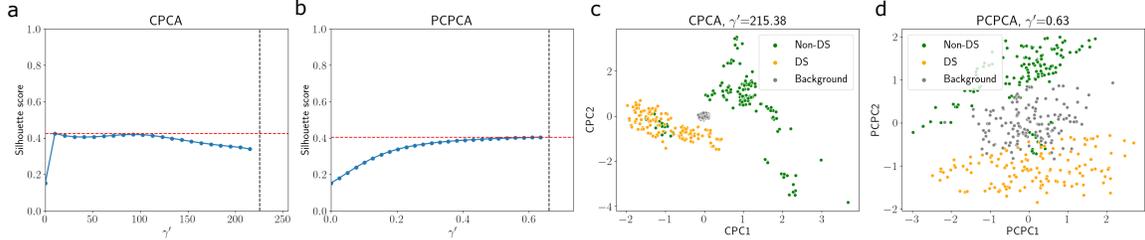


Figure 4: **Comparison of CPCA and PCPCA on the mouse protein expression dataset.** (a) and (b) show the silhouette score of the foreground latent variables at a range of values for $\gamma'$ for each method. The dotted vertical lines shows the first value of $\gamma'$ at which each method "failed," and the red horizontal lines indicate the maximium silhouette score achieved by each method. (c) and (d) show the latent variables for each method at the largest value of $\gamma'$ at which each method "succeeded." Color labels are the same as Figure 3.

### 6.2.2   Single-cell RNA sequencing data

To test our model in a high-dimensional setting, we fit PCPCA with $d = 2$ to a single-cell RNA sequencing (scRNAseq) dataset (Zheng et al., 2017). Here, the foreground dataset contains gene expression measurements from bone marrow mononuclear cells (BMMCs) derived from a patient with acute myeloid leukemia (AML) before and after they received a stem-cell transplant ($n = 4501$). The background dataset contains gene expression measurements of BMMCs from a healthy patient ($m = 1985$). We preprocessed the data by log-transforming and subsetting to the 500 most variable genes, in accordance with previous analyses on these data (Zheng et al., 2017; Abid et al., 2018).

Visualizing the two-dimensional latent variables from PCPCA, we found that the model separates the pre- and post-transplant cells effectively at higher values of $\gamma'$, while PPCA ($\gamma' = 0$) fails to do so (Figure 5). Furthermore, we measured the silhouette score for these two foreground subgroups in the CPCA and PCPCA reduced-dimension spaces. Similar to our observation with the mouse protein expression dataset, we found that, for PCPCA, the silhouette score monotonically increased with $\gamma'$, while CPCA's performance peaked at lower allowable values of $\gamma'$ (Figure 6). Additionally, the allowable range for $\gamma'$ in CPCA was again much larger than that for PCPCA. The maximum silhouette scores achieved by each method were roughly equivalent (CPCA: 0.20, PCPCA: 0.23). These results imply that

PCPCA is effective with high-dimensional data and further demonstrate the advantage of PCPCA's parameterization over CPCA.
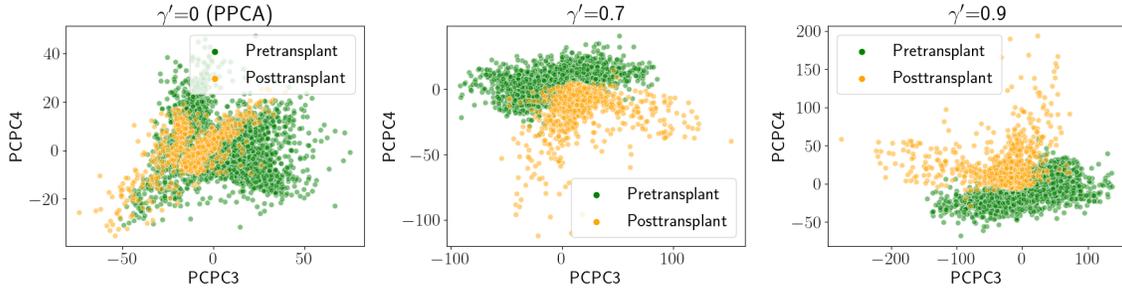


Figure 5: **PCPCA applied to single-cell RNA-seq data.** This dataset contains $n = 4501$ foreground samples (cells from AML patient, plotted here) and $m = 1985$ background samples (cells from a healthy patient, not plotted here). The foreground cells contain two subgroups: *Pre-transplant* (green) and *Post-transplant* (orange). Plotted here are the two foreground groups projected onto PCPCA's third and fourth components for varying values of $\gamma'$.
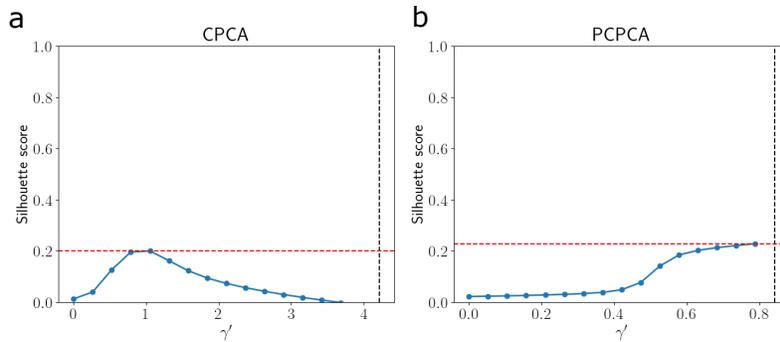


Figure 6: **Silhouette scores for PCPCA and CPCA clusters in scRNA-seq data.** Shown here are the silhouette score across a range of $\gamma'$ for (a) CPCA and (b) PCPCA. The dotted vertical lines shows the first value of $\gamma'$ at which each method "failed," and the red horizontal lines indicate the maximum silhouette score achieved by each method.

## 6.3    Robustness to noise

An advantage of PCPCA's model-based approach to contrastive learning is its ability to explicitly account for noise in the data. To test this directly, we again fit PCPCA and CPCA

17

to the mouse protein expression dataset, but this time we injected additive, independent Gaussian noise across the features. In particular, we transformed every foreground and background sample $x_i$ and $y_j$ as

$$\widetilde{x}_i = x_i + \epsilon_i$$
$$\widetilde{y}_j = y_j + \epsilon_j,$$

where $\epsilon_i, \epsilon_j \sim \mathcal{N}(0, \sigma^2 I_D)$. We generated ten datasets for $\sigma^2 \in \{0.5, 1, \ldots, 5\}$. We also included the case when $\sigma^2 = 0$, which is the original dataset with no additional noise.

We fit PCPCA and CPCA on each of these noisy datasets and measured the silhouette score of PCPCA and CPCA with $d = 2$. We repeated this experiment 100 times for each value of $\sigma^2$. We tuned $\gamma$ independently for PCPCA and CPCA for each value of $\sigma^2$ and took the $\gamma$ with the highest silhouette score. We found that, while the performance of both methods declined with more noise, PCPCA showed better performance than CPCA at higher noise levels (Figure 7). This suggests that PCPCA is more robust to noise than CPCA, demonstrating another advantage of our model-based approach.
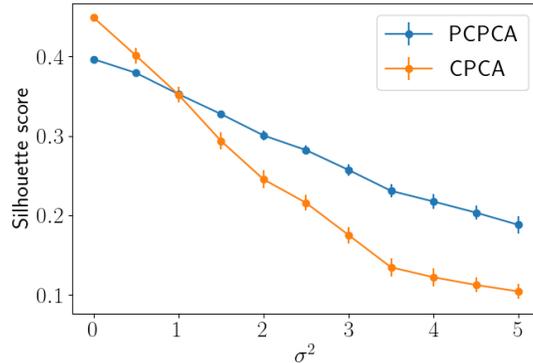


Figure 7: **Model performance with increasing noise.** We injected independent additive Gaussian noise with variance $\sigma^2$ to the mouse protein expression dataset. We then measured the silhouette score of the foreground latent variables. CPCA is shown in orange, and PCPCA is show in blue, both with 95% confidence interval whiskers.

## 6.4 Generating data from the foreground distribution

Another advantage of PCPCA's model-based approach is the ability to generate data from the foreground data distribution. In CPCA, this is not possible because there is no associated generative model. Note that, in PCPCA, we cannot reasonably generate data from the background distribution because the objective function is a relative likelihood with the goal

of minimizing the relative likelihood of the background model. This is not a problem in most settings, as we are typically interested in exploring the variance unique to the foreground data.

To demonstrate PCPCA's ability to generate realistic foreground data, we used the corrupted MNIST dataset (Abid et al., 2018). In this dataset, the foreground samples are MNIST digits (0s and 1s) superimposed onto natural images of grass from ImageNet (Russakovsky et al., 2015). The background samples are unaltered natural images of grass (Figure 8).
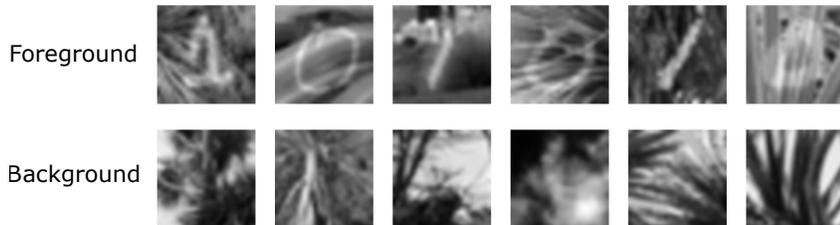


Figure 8: **Examples of the corrupted MNIST dataset.** The top row contains examples from the foreground data, and the bottom row contains examples from the background data.

We fit PCPCA with $d = 2$ and $\gamma' = 0.8$ to obtain $\widehat{W}_{ML}$. For comparison, we also fit PPCA ($\gamma' = 0$). Examining the latent variables, we found that PCPCA showed substantially better clustering of the two MNIST digits than PPCA — the silhouette score for PCPCA was 0.33, while the score for PPCA was just 0.007 (Figure 9a, c).

To generate new data, we sampled $S = 300$ i.i.d. latent variables $z_s \sim \mathcal{N}(0, I_d)$ for $s = 1, \ldots, S$, and projected these to the data space to obtain synthetic images. Specifically, each generated image is computed as $\widehat{x}_s = \widehat{W}_{ML} z_s + \mu_x$ where $\mu_x$ is the mean of the foreground data. We found that these samples recovered the variation in the MNIST digits in the foreground data (Figure 9d). In contrast, samples generated from PPCA did not show as much of the digit structure (Figure 9b). These results suggest that PCPCA can generate realistic data from the foreground distribution, which is useful for exploratory data analysis.

Furthermore, using $\widehat{W}_{ML}$ estimated for PPCA and PCPCA fit to the corrupted MNIST data, we computed the log likelihood of a set of held-out samples of MNIST digits without any corruption. We found that PCPCA has a higher test likelihood than PPCA on these uncorrupted digits (Figure 10). This suggests that the foreground model for PCPCA more accurately captures the uncorrupted MNIST digits relative to PPCA.
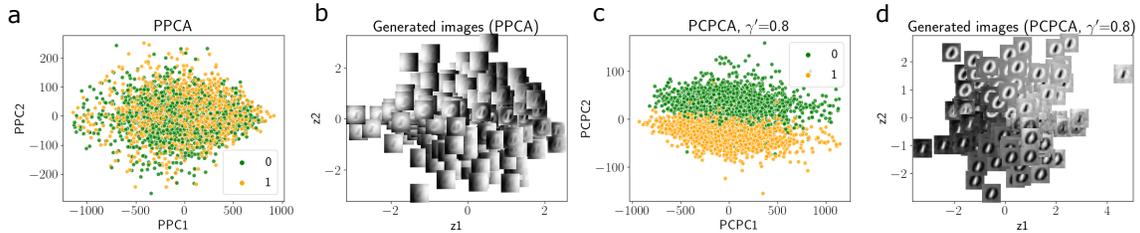
Figure 9: **MNIST digits generated from the PCPCA model.** (a) and (c) show the projected foreground samples for PPCA ($\gamma' = 0$) and PCPCA ($\gamma' = 0.8$), respectively. (b) and (d) show new foreground samples generated from the foreground distribution of PPCA and PCPCA, respectively.
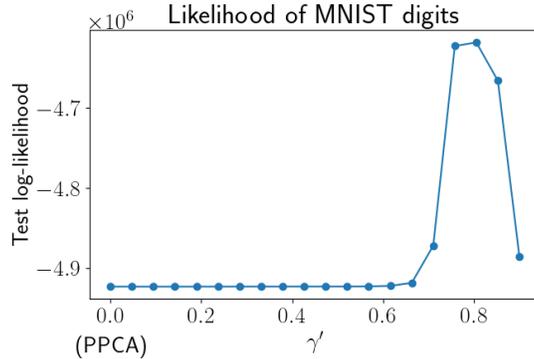


Figure 10: **Log-likelihood of held-out MNIST digits.** Using a PCPCA model fit on the corrupted MNIST dataset, plotted here is the log-likelihood of a held-out set of uncorrupted MNIST digit images. Note that $\gamma' = 0$ corresponds to PPCA.

## 6.5   Gibbs posterior sampling

We next sought to numerically evaluate the Gibbs posterior for PCPCA. To estimate the posterior, any sampling-based inference methods can be applied. We use the No U-Turn Sampler (Hoffman and Gelman, 2014) — which is an extension of Hamiltonian Monte Carlo — as implemented in the Stan programming language (Carpenter et al., 2017). We place uniform priors on $W$ and $\sigma^2$, as required by our theoretical results.

### 6.5.1   Visualizing Gibbs posterior samples

First, we sought to visualize the posterior for $W$. To do so, we used the same toy dataset as our initial experiments in Figure 2. Recall that these samples are generated from a

mixture model in which the background distribution is a two-dimensional Gaussian, and the foreground distribution is a mixture of two Gaussians. Specifically,

$$p(x) = \beta \left\{ \pi \mathcal{N}(x|\mu_{f1}, \Sigma_{f1}) + (1 - \pi)\mathcal{N}(x|\mu_{f2}, \Sigma_{f2}) \right\} + (1 - \beta)\mathcal{N}(x|\mu_b, \Sigma_b)$$

where $\beta$ controls the mixture proportion between the background and foreground, and $\pi$ controls the mixture proportion between the foreground subgroups. In this case, we set $\beta = \pi = 0.5$, $\Sigma_b = \Sigma_f = \left( \begin{smallmatrix} 4.0 & 2.6 \\ 2.6 & 4.0 \end{smallmatrix} \right)$, $\mu_b = \left( \begin{smallmatrix} 0 \\ 0 \end{smallmatrix} \right)$, $\mu_{f1} = \left( \begin{smallmatrix} -1.5 \\ 1.5 \end{smallmatrix} \right)$, $\mu_{f1} = \left( \begin{smallmatrix} 1.5 \\ -1.5 \end{smallmatrix} \right)$, and $\gamma = 0.85$. Using this model, we generated three datasets with increasing sample sizes, respectively containing 30, 60, and 300 samples in each condition. After estimating the Gibbs posterior using each dataset, we drew 100 samples from the posterior for $W$ (Figure 11).

As expected, we found that the posterior increasingly concentrated around the axis separating the foreground subgroups as the sample size increased. With $n = 30$, the posterior draws for $W$ were close to uniformly distributed, but with $n = 300$, the posterior became tightly concentrated around the desired axis. This suggests that the PCPCA Gibbs posterior is a viable tool for accounting for uncertainty in the context of our loss-based modeling framework. Furthermore, it confirms that the posterior can be estimated using well-known MCMC methods, not requiring any model-specific algorithms.
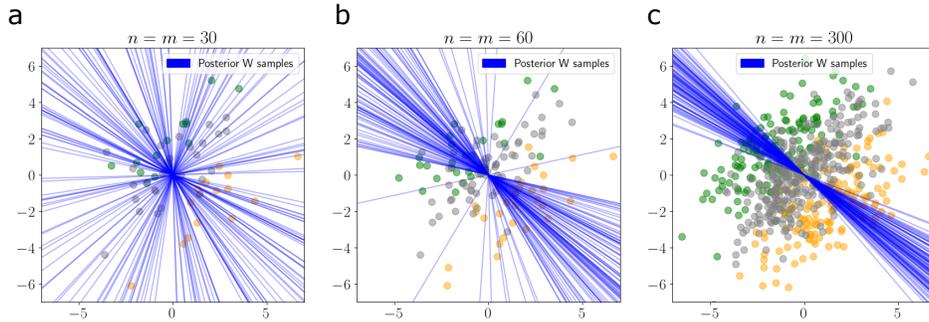


Figure 11: **Sampling from the Gibbs posterior.** Shown here are samples from the Gibbs posterior for $W$ using the toy data. We drew 100 samples from the posterior for varying sample sizes: (a) 30, (b) 60, (c) 300.

### 6.5.2 Posterior convergence rate

Next, we sought to validate the posterior convergence rate of $n^{-1/2}$ for the PCPCA Gibbs posterior. To do so, we again simulated data from a mixture of two-dimensional Gaussians. However, in this experiment, we set the foreground to be a single multivariate Gaussian, rather than a mixture of two Gaussians.

To estimate the convergence rate, we first fit the Gibbs posterior $\Pi_n$ to the simulated data, setting $d = 2$ in this case. We then sampled $T = 1,000$ parameter values from the

posterior, $(W_1, \sigma_1^2), \ldots, (W_T, \sigma_T^2) \sim \Pi_n$. We estimated the divergence using each of these samples and the true risk-minimizing parameter values $(W^\star, \sigma^{2\star})$ in Equation 12. Recall that, in this case, the divergence is $d(\theta, \theta^\star) = (R(\theta) - R(\theta^\star))^{1/2}$, where $\theta = (W, \sigma^2)$ and $R(\cdot)$ is the risk (Equation 11). Finally, we computed the fraction of these divergences that exceeded $n^{-1/2}$. Specifically, we computed

$$\widehat{D} = \frac{1}{T} \sum_{t=1}^{T} I\left(d(\theta_t, \theta^\star) > Cn^{-1/2}\right),$$

where $I$ is the indicator function. We estimated this quantity for $n \in \{50, 200, 300, 400, 500\}$, where $n$ is the total number of samples across conditions. We repeated this five times for each value of $n$.

We found that $\widehat{D}$ fell to zero as $n$ increased, which matches our theoretical result (Figure 12). This result numerically validates PCPCA's posterior convergence rate of $n^{-1/2}$, which is optimal. It also provides further evidence that the PCPCA Gibbs posterior is a principled tool for performing inference in our framework.
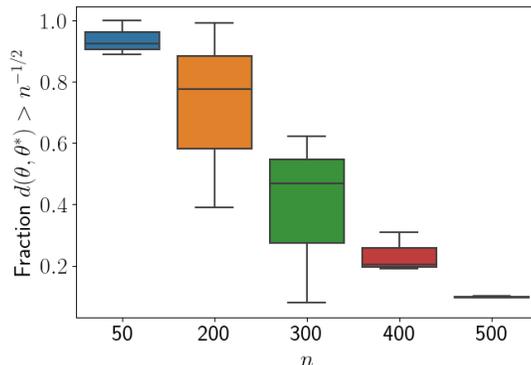


Figure 12: **Empirical validation of the posterior contraction rate for PCPCA Gibbs posterior.** For increasing values of $n$, we randomly drew 1000 samples for $W$; computed the divergence for each; and computed the fraction of these divergences that exceeded $n^{-1/2}$. The boxes show the results for ten repetitions for each value of $n$.

# 7 Contrastive PCA for missing data

Another advantage of our probabilistic modeling approach is the ability to handle missing data in a principled way. Missing or incomplete data is extremely common in real-world datasets. Due to its lack of a probabilistic model, CPCA is unable to deal with missing data.

In this section, we show how to find maximum likelihood estimates for PCPCA in settings with missing data, and we demonstrate this method through simple experiments.

PPCA can handle missing data, where the MLE relies on an EM algorithm that iteratively reconstructs the missing matrix elements from the PCs, and re-estimates the PCs from the expected complete matrix (Tipping and Bishop, 1999; Roweis, 1998). However, in PCPCA, the target function (10) to be maximized is not a likelihood, so we cannot apply the EM algorithm. Instead, we propose a data augmentation method by introducing a indicator matrix representing the location of missing elements to obtain closed-form gradients of the objective so that we can make use of existing gradient-based optimization algorithms, such as gradient descent. We first present the details of our approach and then demonstrate its performance through experiments.

## 7.1   Gradient descent with missing data

Assume some elements of both the background and foreground matrices are missing. Let $x = [x^o, x^u]^\top$, where $x^o$ is the sub-vector of observed features and $x^u$ unobserved, and $y = [y^o, y^u]$ with the same partition. Consider the missing-at-random setting, where the MLE of the complete data is the same as the MLE of the non-missing data only. Let $x_i^o$ be the observed subvector of $x_i$ with length $D_i$, where the locations observed are $i_1, \cdots, i_{D_i}$. Then, we introduce a indicator matrix $L_i$ with dimension $D_i \times D$ such that $x_i^o = L_i x_i$:

$$(L_i)_{kl} = \begin{cases} 1 & l = i_k \\ 0 & \text{otherwise.} \end{cases}$$

Similarly, let $y_j^o$ be the observed subvector of $y_j$ with length $E_j$, and define $M_j \in \mathbb{R}^{E_j \times D}$ as before such that $y_j^o = M_j y_j$. Observe that

$$x_i^o \sim N(0, A_i), \ \ y_j^o \sim N(0, B_j), \ \ A_i := L_i(WW^\top + \sigma^2 \mathrm{I}_D)L_i^\top, \ \ B_j := M_j(WW^\top + \sigma^2 \mathrm{I}_D)M_j^\top.$$

As a result, the objective function of the observed data is

$$l(W, \sigma^2) = l(X^o | W, \sigma^2) - \gamma l(Y^o | W, \sigma^2)$$
$$= -\frac{1}{2}\sum_{i=1}^{n}\left(D_i \log(2\pi) + \log\det(A_i) + \mathrm{tr}(A_i^{-1}x_i^o x_i^{o\top})\right) + \frac{\gamma}{2}\sum_{j=1}^{m}\left(E_i \log(2\pi) + \log\det(B_j) + \mathrm{tr}(B_j^{-1}y_j^o y_j^{o\top})\right).$$

Then we take the derivative w.r.t. to $W$:

$$\frac{\partial l}{\partial W} = -\left\{\sum_{i=1}^{n}L_i^\top A_i^{-1}\left(\mathrm{I}_{D_i} - x_i^o x_i^{o\top}A_i^{-1}\right)L_i - \gamma\sum_{j=1}^{m}M_j^\top B_j^{-1}\left(\mathrm{I}_{E_j} - y_j^o y_j^{o\top}B_j^{-1}\right)M_j\right\}W.$$

Similarly, the derivative w.r.t. $\sigma^2$ is

$$
\begin{aligned}
\frac{\partial l}{\partial \sigma^2} = & -\frac{1}{2} \sum_{i=1}^{n} \left( \mathrm{tr}(A_i^{-1} L_i L_i^\top) - \mathrm{tr}(A_i^{-1} x_i^o x_i^{o\top} A_i^{-1} L_i L_i^\top) \right) \\
& + \frac{\gamma}{2} \sum_{j=1}^{m} \left( \mathrm{tr}(B_j^{-1} M_j M_j^\top) - \mathrm{tr}(B_j^{-1} y_j^o y_j^{o\top} B_j^{-1} M_j M_j^\top) \right).
\end{aligned}
$$

We can then use iterative optimization algorithms, such as gradient descent, to find $\widehat{W}_{ML}$ and $\widehat{\sigma}^2_{ML}$.

## 7.2   Imputing missing data

After finding $\widehat{W}_{ML}$ and $\widehat{\sigma}^2_{ML}$ as above, the unobserved foreground values can be imputed. Let $P_i$ be an indicator matrix with dimension $U_i \times D$, where $U_i = D - D_i$, such that $x_i^u = P_i x_i$. Further, define

$$
C_i := P_i(\widehat{W}_{ML}\widehat{W}_{ML}^\top + \widehat{\sigma}^2_{ML} \,\mathrm{I}_D)P_i^\top, \quad F_i := P_i(\widehat{W}_{ML}\widehat{W}_{ML}^\top + \widehat{\sigma}^2_{ML} \,\mathrm{I}_D)L_i^\top.
$$

Continuing to assume mean-centered data, observe that

$$
x_i^u | x_i^o \sim N(F_i A_i^{-1} x_i^o, C_i - F_i A_i^{-1} F_i^\top).
$$

The unobserved values can then be imputed using the conditional mean $\widehat{x}_i^u = F_i A_i^{-1} x_i^o$.

## 7.3   Experiments with missing data

### 7.3.1   Simulated data

To test PCPCA in the presence of missing data, we first fit the model to a synthetic dataset.

To construct the dataset, we generated foreground and background data from separate PPCA models in order to give them separate covariance structures. In particular, for $i \in [n]$ and $j \in [m]$, we sampled $x_i \sim \mathcal{N}(W^f z_i, \sigma^2 \,\mathrm{I}_D)$ and $y_j \sim \mathcal{N}(W^b z_j, \sigma^2 \,\mathrm{I}_D)$ where $z_i, z_j \sim \mathcal{N}(0, \mathrm{I}_d)$. In our experiments, we set $n = m = 100$, $D = 10$, $d = 2$, and $\sigma^2 = 1$. We sampled the elements of $W^f$ and $W^b$ independently at random from a standard Gaussian.

To test the performance of PCPCA, we randomly removed elements of these two matrices with probability $p$, simulating a missing-at-random scenario. In our experiments, we used $p \in \{0, 0.1, 0.2, \ldots, 0.7\}$. After removing the randomly chosen values, we fit PCPCA on the partially observed dataset using the gradients derived in the previous section, along with the Adam optimizer for additional stability (Kingma and Ba, 2014). Using the fitted model, we computed the log likelihood of a held-out dataset of foreground data. For comparison, we also fit PPCA on the pooled data $X \cup Y$ and reported the log likelihood of the held-out foreground dataset.

We found that PCPCA showed relatively steady test log likelihood for $p \leq 0.3$ (Figure 13a). For higher levels of missing data, PCPCA showed a steady decline in test log likelihood. In contrast, PPCA showed a substantially lower test log-likelihood than PCPCA at all values of $p$. Meanwhile, CPCA and PCA do not allow for settings where $p > 0$.

For each of these partially-observed datasets, we also imputed the missing values in the foreground data and computed the reconstruction error. As before, let $x_i^u \in \mathbb{R}^{U_i}$ be the true values for the unobserved portion of $x_i$, and let $\widehat{x}_i^u$ be the PCPCA reconstruction of these values. We computed the mean-squared error of these reconstructions:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{U_i} \|x_i^u - \widehat{x}_i^u\|_2^2.$$

We found that PCPCA achieved low reconstruction error when few values were missing, and the error increased steadily for higher fractions of unobserved values (Figure 13b). In all cases, PCPCA performed better than PPCA.

These results suggest that PCPCA is robust in the presence of missing data, even when a relatively large fraction of the data is missing at random.
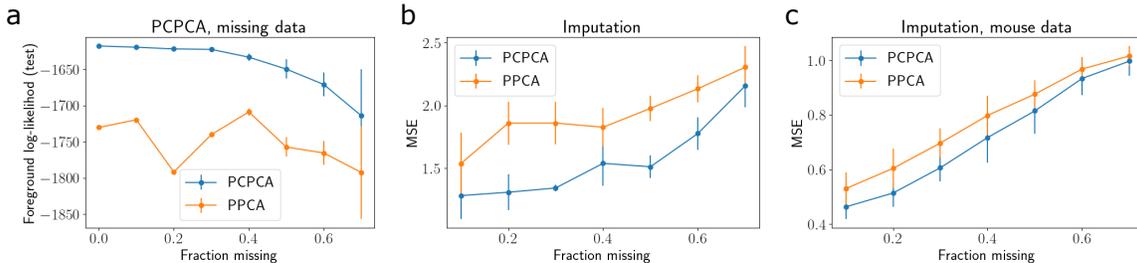


Figure 13: **PCPCA and PPCA with missing data.** In (a) and (b), we used a simulated dataset in which $n = m = 100$, $D = 10$, and $X$ and $Y$ were generated with separate PPCA models. We fit PCPCA on $X$ and $Y$ and PPCA on the concatenation of $X$ and $Y$. Values were dropped with increasing probability. (a) shows the test log-likelihood was computed for a held-out foreground dataset. (b) shows the mean-squared error of the imputed values for the training data from the PCPCA model. (c) shows the imputation error for the mouse protein expression data.

### 7.3.2   Mouse protein expression data

To further validate PCPCA's ability to handle missing data, we applied the model to the mouse protein expression dataset. Here, we randomly removed elements from the foreground and background data with probability $p$, where $p \in \{0, 0.2, 0.6, 0.9\}$. We fit the model to each partially-observed dataset using gradient descent and the Adam optimizer to obtain

$\widehat{W}_{ML}$. We set $d = 2$ and $\gamma = 0.4$ for all runs based on previous experiments. Finally, we projected the fully-observed dataset to the latent space and computed the silhouette score.

We observed that the two subgroups of mice were preserved even with a large fraction of the data masked (Figure 14). The silhouette scores confirmed this, staying steady for $p < 0.6$. Furthermore, we imputed the missing foreground values using the PCPCA model, and we found that the model's reconstructions consistently showed lower error than PPCA (Figure 13c).

These results imply that PCPCA could be used in many real-world settings in which datasets are only partially observed.
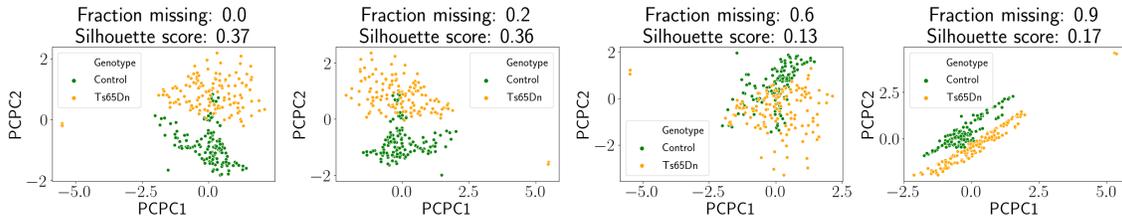


Figure 14: **PCPCA on mouse protein expression data with missing data.** We randomly removed values from the data and fit PCPCA. Shown here are the latent projections of the fully-observed foreground dataset for increasing missing probabilities.

# 8 Discussion

In this paper, we presented a probabilistic model, PCPCA, for learning the contrastive dimensions between a foreground dataset and a background dataset. We derived conditions for the tuning parameter $\gamma$ under which our model and a previous method, CPCA, are well-defined. In contrast to CPCA, our model-based approach allows for uncertainty quantification, is robust to noise, includes the ability to sample from the fitted model, and is able to accept and impute missing data. We developed a generalized Bayesian framework that allows for principled inference in our model, despite it not having a well-defined likelihood. To find the loss-minimizer in the presence of missing data, we derived a gradient descent algorithm that only relies on the observed data. We demonstrated the utility of PCPCA in several applications using protein expression, gene expression, and image data. We found that PCPCA outperformed PPCA and CPCA in capturing subgroup structure and in its robustness to noise and missing data.

Several future directions remain to be explored. First, more general inference procedures could be developed for the PCPCA relative likelihood objective function. Likelihood ratios have been well-studied for ratios comparing two sets of parameters under a single shared dataset (Anderson, 1962). However, there has been little work studying relative likelihoods

that compare one set of shared parameters under two datasets. Typical inference procedures — such as expectation maximization (EM) – cannot be used in this contrastive setting because the objective function is not a well-defined likelihood. Future work will benefit from adapting well-known methods, such as EM, for estimation, inference, and optimization to the novel relative likelihood objective presented in this work.

Second, more sophisticated optimization procedures could be used for performing gradient descent in the presence of missing data. In both PPCA and PCPCA, the gradient of the likelihood w.r.t. the noise variance $\sigma^2$ shows a sharp increase as $\sigma^2 \to 0$. Modern optimization techniques could be used to further stabilize the gradient in this regime.

Finally, the PCPCA model itself could be extended in several ways. Future versions could incorporate various data likelihoods in order to capture non-Gaussian data. Other extensions might allow multiple foreground datasets, possibly allowing structured relationships between those foreground matrices. Furthermore, a non-linear version of PCPCA could be considered.

# References

Abid, A., Zhang, M. J., Bagaria, V. K., and Zou, J. (2017). Contrastive principal component analysis. *arXiv preprint arXiv:1709.06716*.

Abid, A., Zhang, M. J., Bagaria, V. K., and Zou, J. (2018). Exploring patterns enriched in a dataset with contrastive principal component analysis. *Nature Communications*, 9(1):1–7.

Anderson, T. W. (1962). An introduction to multivariate statistical analysis. Technical report, Wiley New York.

Boileau, P., Hejazi, N. S., and Dudoit, S. (2020). Exploring high-dimensional biological data with sparse contrastive principal component analysis. *Bioinformatics*, 36(11):3422–3430.

Brenner, N., Bialek, W., and Van Steveninck, R. d. R. (2000). Adaptive rescaling maximizes information transmission. *Neuron*, 26(3):695–702.

Candès, E. J., Li, X., Ma, Y., and Wright, J. (2011). Robust principal component analysis? *Journal of the ACM*, 58(3):1–37.

Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M. A., Guo, J., Li, P., and Riddell, A. (2017). Stan: a probabilistic programming language. *Grantee Submission*, 76(1):1–32.

Darbyshire, J. and Hamish, J. (2016). The pricing and hedging of interest rate derivatives: A practical guide to swaps.

Fruchter, B. (1954). *Introduction to factor analysis.* Van Nostrand.

Hastie, T. and Stuetzle, W. (1989). Principal curves. *Journal of the American Statistical Association*, 84(406):502–516.

Higuera, C., Gardiner, K. J., and Cios, K. J. (2015). Self-organizing feature maps identify proteins critical to learning in a mouse model of down syndrome. *PloS One*, 10(6):e0129126.

Hoffman, M. D. and Gelman, A. (2014). The no-u-turn sampler: adaptively setting path lengths in hamiltonian monte carlo. *J. Mach. Learn. Res.*, 15(1):1593–1623.

Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology*, 24(6):417.

Jirsa, V. K., Friedrich, R., Haken, H., and Kelso, J. S. (1994). A theoretical model of phase transitions in the human brain. *Biological Cybernetics*, 71(1):27–35.

Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Lawrence, N. (2003). Gaussian process latent variable models for visualisation of high dimensional data. *Advances in Neural Information Processing Systems*, 16:329–336.

Novembre, J. and Stephens, M. (2008). Interpreting principal component analyses of spatial population genetic variation. *Nature Genetics*, 40(5):646–649.

Pasini, G. (2017). Principal component analysis for stock portfolio management. *International Journal of Pure and Applied Mathematics*, 115(1):153–167.

Rousseeuw, P. J. (1987). Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20:53–65.

Roweis, S. T. (1998). EM algorithms for PCA and SPCA. In *Advances in Neural Information Processing Systems*, pages 626–632.

Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. (2015). Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252.

Schölkopf, B., Smola, A., and Müller, K.-R. (1998). Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10(5):1299–1319.

Severson, K. A., Ghosh, S., and Ng, K. (2019). Unsupervised learning with contrastive latent variable models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4862–4869.

Syring, N. and Martin, R. (2020). Gibbs posterior concentration rates under sub-exponential type losses. *arXiv preprint arXiv:2012.04505*.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B*, 58(1):267–288.

Tipping, M. E. and Bishop, C. M. (1999). Probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B*, 61(3):611–622.

Twine, N. A., Janitz, K., Wilkins, M. R., and Janitz, M. (2011). Whole transcriptome sequencing reveals gene expression and splicing differences in brain regions affected by alzheimer's disease. *PloS One*, 6(1):e16266.

Vidal, R., Ma, Y., and Sastry, S. (2005). Generalized principal component analysis (GPCA). *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(12):1945–1959.

Young, M. D., Mitchell, T. J., Braga, F. A. V., Tran, M. G., Stewart, B. J., Ferdinand, J. R., Collord, G., Botting, R. A., Popescu, D.-M., Loudon, K. W., et al. (2018). Single-cell transcriptomes from human kidneys reveal the cellular identity of renal tumors. *Science*, 361(6402):594–599.

Zheng, G. X., Terry, J. M., Belgrader, P., Ryvkin, P., Bent, Z. W., Wilson, R., Ziraldo, S. B., Wheeler, T. D., McDermott, G. P., Zhu, J., et al. (2017). Massively parallel digital transcriptional profiling of single cells. *Nature Communications*, 8(1):1–12.

Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B*, 67(2):301–320.

Zou, H., Hastie, T., and Tibshirani, R. (2006). Sparse principal component analysis. *Journal of Computational and Graphical tatistics*, 15(2):265–286.

Zou, J. Y., Hsu, D. J., Parkes, D. C., and Adams, R. P. (2013). Contrastive learning using spectral methods. *Advances in Neural Information Processing Systems*, 26:2238–2246.

# 9 Appendix

## 9.1 Data availability

All data are available via their respective papers. Preprocessing scripts are included in the code repository: `https://github.com/andrewcharlesjones/pcpca`.

## 9.2 Proof of Theorem 1

*Proof.* We simplify the geometric objective function in (6) until it matches (5), the statistical objective function. Let $v_1$ be the eigenvector of $C$ corresponding to the largest eigenvalue,

then

$$v_1 = \operatorname*{argmin}_{v^\top v=1} \frac{1}{n} \sum_{i=1}^{n} \|x_i - vv^\top x_i\|^2 - \gamma \frac{1}{m} \sum_{j=1}^{m} \|y_j - vv^\top y_j\|^2$$

$$= \operatorname*{argmin}_{v^\top v=1} \frac{1}{n} \sum_{i=1}^{n} \left( x_i^\top x_i - 2x_i^\top vv^\top x_i + x_i vv^\top vv^\top x \right) - \gamma \frac{1}{m} \sum_{j=1}^{m} \left( y_j^\top y_j - y_j^\top vv^\top y_j + y_j^\top vv^\top vv^\top y_j \right)$$

$$= \operatorname*{argmin}_{v^\top v=1} \frac{1}{n} \sum_{i=1}^{n} \left( -2x_i^\top vv^\top x_i + x_i vv^\top x \right) - \gamma \frac{1}{m} \sum_{j=1}^{m} \left( -y_j^\top vv^\top y_j + y_j^\top vv^\top y_j \right)$$

$$= \operatorname*{argmax}_{v^\top v=1} \frac{1}{n} \sum_{i=1}^{n} x_i^\top vv^\top x_i - \gamma \frac{1}{m} \sum_{j=1}^{m} y_j^\top vv^\top y_j$$

$$= \operatorname*{argmax}_{v^\top v=1} v^\top \frac{1}{n} \sum_{i=1}^{n} x_i x_i^\top v - \gamma v^\top \frac{1}{m} \sum_{j=1}^{m} y_j y_j^\top v$$

$$= \operatorname*{argmax}_{v^\top v=1} v^\top C_x v - \gamma v^\top C_Y v.$$

$\square$

## 9.3  Proof of Lemma 1

*Proof.* Let $\lambda_1 \geq \cdots \geq \lambda_D$ and $\rho_1 \geq \cdots \geq \rho_D$ be the eigenvalues of $C_X$ and $C_Y$ in descending order. Assume $\gamma < \lambda_D/\rho_1$; then for any unit vector $v \in \mathbb{R}^D$ with $\|v\| = 1$,

$$v^\top C v = v^\top C_X v - \gamma v^\top C_Y v$$
$$\geq \lambda_D v^\top v - \gamma \rho_1 v^\top v$$
$$= \lambda_D - \gamma \rho_1 > 0,$$

so $C$ is positive definite.

Assume $\gamma \geq \lambda_D/\rho_1$ with corresponding eigenvectors $u_D$ and $v_1$ where $u_D = v_1$. Let $u = u_D = v_1$ with $\|u\| = 1$, then

$$u^\top C u = u^\top C_X v - \gamma u^\top C_Y u$$
$$= \lambda_D - \gamma \rho_1 \leq 0.$$

So $C$ is not positive definite.

$\square$

Note that the second half of the proof is the worst case, where the last eigenvector of $C_X$ matches the first eigenvector of $C_Y$. In order to make $C$ positive definite (PD), the strong condition is necessary. However, in practice, much larger values of $\gamma$ are sometimes allowed such that $C$ is still PD.

## 9.4 Proof of Theorem 2

*Proof.* Recall that the eigenvalues of $C$, $C_X$ and $C_Y$ are $\lambda_1 \geq \cdots \lambda_D \geq 0$, $\mu_1 \geq \cdots \mu_D \geq 0$ and $\rho_1 \geq \cdots \geq \rho_D \geq 0$. Then the eigenvalues of $-\gamma C_Y$ are $-\gamma\rho_D \geq \cdots \geq -\gamma\rho_1$. Since $C = C_X - \gamma C_Y = C_X + (-\gamma C_Y)$, by Weyl's inequalities, for any $j + k \geq D + d$,

$$\lambda_d \geq \mu_j - \gamma\rho_{D-k+1} > 0.$$

Similar to the proof of Lemma 1, if the condition is violated, there exists a $C$ such that the first $d$ eigenvalues are negative. $\qquad\square$

## 9.5 Proof of Corollary 1

*Proof.* By the same proof as the proof for the PCA loss, the CPCA loss is $\sum_{i=d+1}^{D} \lambda_i$. For a fixed $d$, we first show that increasing $\gamma$ will result in a smaller $\lambda_i$. Let $\gamma_1 < \gamma_2$, $C_1 = C_X - \gamma_1 C_Y$ and $C_2 = C_Y - \gamma_2 C_Y$. Then $C_2 - C_1 = (\gamma_1 - \gamma_2)C_Y$ is positive definite, hence the eigenvalues of $C_2$ are greater than those of $C_1$. This implies that increasing $\gamma$ will decrease the loss.

Then, for a fixed $\gamma < \max\left\{\frac{\mathrm{eig}_{d+1}(C_X)}{\mathrm{eig}_1(C_Y)}, \frac{\mathrm{eig}_{d+1}(C_X)}{\mathrm{eig}_2(C_Y)}, \cdots, \frac{\mathrm{eig}_D(C_X)}{\mathrm{eig}_{D-d}(C_Y)}\right\}$, Theorem 1 implies $\lambda_{d+1} > 0$, so raising $d$ to $d+1$ results in a smaller tail sum of the eigenvalues, hence a smaller loss. $\qquad\square$

## 9.6 Proof of Theorem 3

First, we find the maximizer of $W$ given $\sigma^2$. Recall the marginals: $x, y \sim N(0, WW^\top + \sigma^2 I_D)$, and denote $A = WW^\top + \sigma^2 I_D$. Then, taking the log of the objective function, we have

$$l(W, \sigma^2 | X, Y, \sigma^2) = -\frac{n}{2}\left(D\ln(2\pi) + \ln|A| + \mathrm{tr}(A^{-1}C_X)\right) + \frac{\gamma m}{2}\left(D\ln(2\pi) + \ln|A| + \mathrm{tr}(A^{-1}C_Y)\right).$$

We drop all constants and the log likelihood becomes

$$l(W, \sigma^2) = -\frac{n - \gamma m}{2}\ln|A| - \frac{1}{2}\mathrm{tr}(A^{-1}(nC_X - \gamma m C_Y)) = -\frac{n - \gamma m}{2}\ln|A| - \frac{1}{2}\mathrm{tr}(A^{-1}C).$$

Denote $C = nC_X - \gamma m C_Y$ where $C_X = \frac{1}{n}\sum_{i=1}^{n} x_i x_i^\top$ and $C_Y = \frac{1}{m}\sum_{j=1}^{m} y_j y_j^\top$. Then, we take the derivative of $l$:

$$\frac{\partial l}{\partial W} = -\frac{n - \gamma m}{2} 2A^{-1}W + \frac{1}{2}2A^{-1}CA^{-1}W.$$

Letting $\frac{\partial l}{\partial W} = 0$, we have

$$A^{-1}\frac{\sum_{i=1}^{n} x_i x_i^\top - \gamma \sum_{j=1}^{m} y_j y_j^\top}{n - \gamma m}A^{-1}W = A^{-1}W,$$

32

that is, $\frac{1}{n-\gamma m}CA^{-1}W = W$. $C = (n-\gamma m)A$ solves this equation, that is,

$$WW^\top = \frac{1}{n-\gamma m}C - \sigma^2\,\mathrm{I}_D\,.$$

Assume $C = U\Lambda U^\top$, then

$$\widehat{W}_{ML} = U_d\left(\frac{\Lambda_d}{n-\gamma m} - \sigma^2\,\mathrm{I}_d\right)^{1/2}R,$$

where $U_d$ consists of the first $d$ eigenvectors of $C$, $\Lambda_d$ contains the corresponding eigenvalues, and $R$ is any rotation matrix.

Next, we consider $\widehat{\sigma}^2_{ML}$. Plugging $\widehat{W}_{ML}$ into the objective, we have

$$l(\sigma^2|X, Y, \widehat{W}_{ML}) = -\frac{n-\gamma m}{2}\ln|A| - \frac{1}{2}\mathrm{tr}(A^{-1}C)$$

$$= -\frac{n-\gamma m}{2}\left(\sum_{i=1}^{d}\ln\frac{\lambda_i}{n-\gamma m} + (D-d)\ln\sigma^2\right) - \frac{1}{2}\left(d(n-\gamma m) + \frac{1}{\sigma^2}\sum_{j=d+1}^{D}\lambda_j\right).$$

$$= -\frac{n-\gamma m}{2}\left(\sum_{i=1}^{d}\ln\widetilde{\lambda}_i + (D-d)\ln\sigma^2 + d + \frac{1}{\sigma^2}\sum_{j=d+1}^{D}\widetilde{\lambda}_j\right),$$

where $\widetilde{\lambda}_i = \frac{\lambda_j}{n-\gamma m}$. The derivative of $l$ is:

$$\frac{\partial l}{\partial\sigma^2} = -\frac{n-\gamma m}{2}\left(\frac{D-d}{\sigma^2} - \frac{1}{\sigma^4}\sum_{j=d+1}^{D}\widetilde{\lambda}_j\right).$$

Letting $\frac{\partial l}{\partial\sigma^2} = 0$, we have $\frac{D-d}{\sigma^2} = \frac{1}{\sigma^4}\sum_{j=d+1}^{D}\widetilde{\lambda}_j$, so the MLE of $\sigma^2$ is given by

$$\widehat{\sigma}^2_{ML} = \frac{1}{D-d}\sum_{i=d+1}^{D}\widetilde{\lambda}_i = \frac{1}{(D-d)(n-\gamma m)}\sum_{i=d+1}^{D}\lambda_i.$$

We next connect PCPCA to CPCA and PPCA. When $\gamma = 0$, the objective function (10) is the same as the objective of PPCA, so the MLEs are also the same. Alternatively, when $\sigma^2 \to 0$,

$$\widehat{W}_{ML} = U_d\Lambda^{1/2}R$$

is exactly the solution of CPCA with the new parameterization, which corresponds to eigenvectors of $C = \sum_i x_i x_i^\top - \gamma\sum_{j=1}^{m} y_j y_j^\top$, and is equivalent to the CPCA proposed by Abid et al. (2018) with $\gamma' = \gamma m/n$. In this sense, our parameterization is more natural since it corresponds to the likelihood function.

## 9.7 Proof of Theorem 4

First recall the following Lemma.

**Definition 3.** *The the loss function $l$ is said to be of sub-exponential type if there exists $\overline{w}, K, r > 0$ such that for any $w \in (0, \overline{w})$,*

$$d(\theta; \theta^*) > \delta \implies \mathbb{E}_P e^{-w(l_\theta - l_{\theta^*})} \le e^{-Kw\delta^r}. \tag{13}$$

Let $m(\theta, \theta^*) = \mathbb{E}_P(l_\theta - l_{\theta^*}) = R(\theta) - R(\theta^*)$ and $v(\theta, \theta^*) = \mathbb{E}_P \left( l_\theta - l_{\theta^*} - m(\theta, \theta^*) \right)^2$.

**Lemma 2** (Syring and Martin (2020)). *Assume $\varepsilon_n \to 0$ and $n\varepsilon_n^r \to \infty$ for $r > 0$, the prior satisfies*

$$\log \Pi(\{\theta : m(\theta, \theta^*) \vee v(\theta, \theta^*) \le \varepsilon_n^r\}) \gtrsim -Mn\varepsilon_n^r$$

*and the loss function is of sub-exponential type, then the Gibbs posterior distribution has asymptotic concentration rate $\varepsilon_n$ for all large enough constants $M > 0$.*

By the definition of the divergence: $d(\theta; \theta^*) = (R(\theta) - R(\theta^*))^{1/2}$, we know that

$$d(\theta; \theta^*) > \delta \implies \mathbb{E}_P e^{-w(l_\theta - l_{\theta^*})} \le e^{-w(R(\theta) - R(\theta^*))} \le e^{-w\delta^2},$$

so the loss is of sub-exponential type. Then it suffices to check that the prior satisfies the conditions in Lemma 2. First we calculate $m$,

$$m(\theta, \theta^*) = R(\theta) - R(\theta^*) = \text{tr}\left( (V^* V^{*\top} - VV^\top)C) \right). \tag{14}$$

Recall that $v(\theta, \theta^*) = \mathbb{E}_P \left( l_\theta - l_{\theta^*} \right)^2 - m(\theta, \theta^*)^2$. We show the following two lemmas to calculate $v$.

**Lemma 3.** $X \sim N(0, C)$, *then* $\mathbb{E}[X^\top X X^\top X] = \text{tr}(C)^2 + 2\text{tr}(C^2)$.

*Proof.* Since $\mathbb{E}[X^\top X X^\top X] = \mathbb{E}\left( \sum_{i,j} X_i^2 X_j^2 \right) = \sum_{i,j} \mathbb{E}(X_i^2 X_j^2)$, we start with $\mathbb{E}(X_i^2 X_j^2)$. Observe that $(X_i, X_j) \sim N\left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} C_{ii} & C_{ij} \\ C_{ij} & C_{jj} \end{bmatrix} \right)$, so $\mathbb{E}(X_i^2 X_j^2) = C_{ii} C_{jj} + 2C_{ij}^2$. Then we have

$$\mathbb{E}[X^\top X X^\top X] = \mathbb{E}\left( \sum_{i,j} X_i^2 X_j^2 \right) = \sum_{i,j} \mathbb{E}(X_i^2 X_j^2)$$

$$= \sum_{i,j} \left( C_{ii} C_{jj} + 2C_{ij}^2 \right) = \text{tr}(C)^2 + 2\text{tr}(CC^\top).$$

$\square$

**Lemma 4.** $X \sim N(0, C)$, *and* $A$ *is symmetric, then* $\mathbb{E}[X^\top X X^\top A X] = \text{tr}(C)\text{tr}(AC) + 2\text{tr}(AC^2)$.

*Proof.* Let $Y = A^{1/2}X$, so $Y \sim N(0, A^{1/2}CA^{1/2})$. Since $\mathbb{E}[X^\top X X^\top A X] = \mathbb{E}\left(\sum_{i,j} X_i^2 Y_j^2\right) = \sum_{i,j} \mathbb{E}(X_i^2 Y_j^2)$, we start with $\mathbb{E}(X_i^2 Y_j^2)$. Observe that

$$(X_i, Y_j) \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} C_{ii} & (CA^{1/2})_{ij} \\ (CA^{1/2})_{ij} & (A^{1/2}CA^{1/2})_{jj} \end{bmatrix}\right),$$

so $\mathbb{E}(X_i^2 Y_j^2) = C_{ii}(A^{1/2}CA^{1/2})_{jj} + 2C_{ij}^2$. Then we have

$$\mathbb{E}[X^\top X X^\top A X] = \mathbb{E}\left(\sum_{i,j} X_i^2 Y_j^2\right) = \sum_{i,j} \mathbb{E}(X_i^2 Y_j^2)$$

$$= \sum_{i,j} \left(C_{ii}(A^{1/2}CA^{1/2})_{jj} + 2(CA^{1/2})_{ij}^2\right) = \operatorname{tr}(C)\operatorname{tr}(A^{1/2}CA^{1/2}) + 2\operatorname{tr}(A^{1/2}CC^\top A^{1/2})$$

$$= \operatorname{tr}(C)\operatorname{tr}(AC) + 2\operatorname{tr}(AC^2).$$

$\square$

Let $\Delta = V^*V^{*\top} - VV^\top$, then observe that

$$\mathbb{E}_p\left[l_\theta(u) - l_{\theta^*}(u)\right]^2 = \mathbb{E}_P\left[(-\gamma)^{2\alpha}(x^\top V^*V^{*\top}x - x^\top VV^\top x)^2\right]$$

$$= \beta \mathbb{E}_{x \sim P_F}\left[x^\top \Delta x x^\top \Delta x\right] + (1-\beta)\gamma^2 \mathbb{E}_{x \sim P_B}\left[x^\top \Delta x x^\top \Delta x\right]$$

$$= \beta \left[\operatorname{tr}(\Delta C_F)^2 + 2\operatorname{tr}(\Delta C_F \Delta C_F)\right] + (1-\beta)\gamma^2 \left[\operatorname{tr}(\Delta C_B)^2 + 2\operatorname{tr}(\Delta C_B \Delta C_B)\right].$$

Rewrite $m$ in a similar form:

$$m(\theta, \theta^*) = R(\theta) - R(\theta^*) = \mathbb{E}_P(l_\theta - l_{\theta^*})$$

$$= \mathbb{E}_P\left[(-\gamma)^\alpha x^\top \Delta x\right]$$

$$= \beta \operatorname{tr}(\Delta C_F) + (1-\beta)\gamma \operatorname{tr}(\Delta C_B).$$

Now we can calculate $v(\theta, \theta^*)$:

$$v(\theta, \theta^*) = \mathbb{E}_p\left[l_\theta(u) - l_{\theta^*}(u)\right]^2 - m(\theta, \theta^*)^2$$

$$= \beta \left[\operatorname{tr}(\Delta C_F)^2 + 2\operatorname{tr}(\Delta C_F \Delta C_F)\right] + (1-\beta)\gamma^2 \left[\operatorname{tr}(\Delta C_B)^2 + 2\operatorname{tr}(\Delta C_B \Delta C_B)\right]$$

$$- \left[\beta \operatorname{tr}(\Delta C_F) + (1-\beta)\gamma \operatorname{tr}(\Delta C_B)\right]^2$$

$$= (\beta - \beta^2)\operatorname{tr}(\Delta C_F)^2 + \gamma^2(\beta - \beta^2)\operatorname{tr}(\Delta C_B)^2 + 2\beta \operatorname{tr}(\Delta C_F \Delta C_F)$$

$$+ 2\gamma^2(1-\beta)\operatorname{tr}(\Delta C_B \Delta C_B) - 2\gamma\beta(1-\beta)\operatorname{tr}(\Delta C_F)\operatorname{tr}(\Delta C_B).$$

*Proof.* Observe that $\|\Delta\|\lambda_D^F \leq |\operatorname{tr}(\Delta C_F)| \leq \|\Delta\|\lambda_1^F$, where $\lambda_1^F$ and $\lambda_D^F$ are the largest and smallest eigenvalue of $C_F$. Then by the above calculation, we know that

$$m(\theta, \theta^*) \sim \|VV^\top - V^*V^{*\top}\|, \quad v(\theta, \theta^*) \sim \|VV^\top - V^*V^{*\top}\|^2.$$

When $n$ is sufficiently large, there exists constant $c$ such that

$$\left\{\theta : \|VV^\top - V^*V^{*\top}\| \leq cn^{-1/2}\right\} \subset \left\{\theta : m(\theta, \theta^*) \vee v(\theta, \theta^*) \leq n^{-1/2}\right\}.$$

So it suffices to check the prior $\Pi$ assigns enough mass around $V^*$ w.r.t. the operator norm. Recall that the Riemannian volume measure on $\operatorname{Gr}(D, d)$, denoted by $\mathcal{H}$, is $O(n)$ invariant, also known as the Haar measure, while the distance on $\operatorname{Gr}(D, d)$ is given by

$$d(V_1, V_2) = \|V_1 V_1^\top - V_2 V_2^\top\|.$$

Denote the ball centered at $V$ with radius $r$ w.r.t. this distance by $B(V, r)$, then there exists constant $c$ such that

$$\frac{1}{c} r^{d(D-d)} \leq \mathcal{H}(B(V, r)) \leq cr^{d(D-d)}, \quad \forall V \in \operatorname{Gr}(D, d), \ r > 0.$$

When the prior $\Pi$ is uniform w.r.t. the Haar measure,

$$\Pi\left\{\theta : m(\theta, \theta^*) \vee v(\theta, \theta^*) \leq n^{-1/2}\right\} \gtrsim \Pi\left\{\theta : \|\Delta\|_o \leq cn^{-1/2}\right\} \gtrsim \left(n^{-1/2}\right)^{d(D-d)},$$

where $d(D-d) = \dim(\operatorname{Gr}(D, d))$, so by (Syring and Martin, 2020, Theorem 3.3), the posterior contraction rate of the Gibbs posterior is $n^{-1/2}$, which is optimal. $\qquad\square$

## 9.8 Proof of Theorem 5

As in the proof of Theorem 4, the loss is of sub-exponential type. Then it suffices to check that the prior $\Pi$ satisfies Lemma 2. To do so, we first calculate $m$ and $v$. Let $\delta = \log|A| - \log|A^*|$ and $\Delta = A^{-1} - A^{*-1}$, then

$$m(\theta, \theta^*) = R(\theta) - R(\theta^*) = \frac{\beta - (1-\beta)\gamma}{2}(\log|A| - \log|A^*|) + \frac{\operatorname{tr}((A^{-1} - A^{*-1})C)}{2}$$
$$= \frac{\beta - (1-\beta)\gamma}{2}\delta + \frac{1}{2}\operatorname{tr}(\Delta C) = O(\sqrt{\delta^2 + \|\Delta\|^2}).$$

Now, to calculate $v$, we have

$$\mathbb{E}_p \left[l_\theta(u) - l_{\theta^*}(u)\right]^2 = \mathbb{E}_P \left[(-\gamma)^{2\alpha} \left(\frac{1}{2}(\log|A| - \log|A^*|) + \frac{1}{2}x^\top(A^{-1} - A^{*-1})x\right)^2\right]$$

$$= \frac{\beta}{4}\mathbb{E}_{x\sim P_F}\left[\delta^2 + 2\delta x^\top \Delta x + x^\top \Delta x x^\top \Delta x\right] + \frac{(1-\beta)\gamma^2}{4}\mathbb{E}_{x\sim P_B}\left[\delta^2 + 2\delta x^\top \Delta x + x^\top \Delta x x^\top \Delta x\right]$$

$$= \frac{\beta}{4}\left[\delta^2 + 2\delta\operatorname{tr}(\Delta C_F) + \operatorname{tr}(\Delta C_F)^2 + 2\operatorname{tr}(\Delta C_F \Delta C_F)\right]$$

$$\quad + \frac{(1-\beta)\gamma^2}{4}\left[\delta^2 + 2\delta\operatorname{tr}(\Delta C_B) + \operatorname{tr}(\Delta C_B)^2 + 2\operatorname{tr}(\Delta C_B \Delta C_B)\right]$$

$$= O(\delta^2 + \|\Delta\|^2).$$

As a result,

$$v(\theta, \theta^*) = \mathbb{E}_p\left[l_\theta(u) - l_{\theta^*}(u)\right]^2 - m(\theta, \theta^*)^2 = O(\delta^2 + \|\Delta\|^2).$$

*Proof.* By the assumption that $\sigma^2 \geq \sigma_0^2$, all eigenvalues of $A = WW^\top + \sigma^2 I_D$ and $A^* = W^*W^{*\top} + \sigma^{*2}I_D$ are lower-bounded by $\sigma_0^2$. As a result, both $\log|\cdot|$ and $\operatorname{tr}(\cdot^{-1}C)$ are Lipschitz, and both $m$ and $v$ can be bounded by the distance between parameters. Then by the above calculation, there exists a $c$ such that

$$m(\theta, \theta^*) \leq c\sqrt{\|WW^\top - W^*W^{*\top}\|^2 + (\sigma^2 - \sigma^{*2})^2},$$

$$v(\theta, \theta^*) \leq c\left(\|WW^\top - W^*W^{*\top}\|^2 + (\sigma^2 - \sigma^{*2})^2\right).$$

When $n$ is sufficiently large, there exists a constant $c$ such that

$$\left\{\theta : \sqrt{\|WW^\top - W^*W^{*\top}\|^2 + (\sigma^2 - \sigma^{*2})^2} \leq cn^{-1/2}\right\} \subset \left\{\theta : m(\theta, \theta^*) \vee v(\theta, \theta^*) \leq n^{-1/2}\right\}.$$

So it suffices to check that the prior $\Pi$ assigns enough mass around $(V^*, \sigma^{*2})$ w.r.t. the product measure. Denote the ball centered at $(V, \sigma^2)$ with radius $r$ w.r.t. this distance by $B(V, \sigma^2, r)$, then there exists a constant $c$ such that

$$\frac{1}{c}r^{Dd+1} \leq \operatorname{Vol}(B(V, \sigma^2, r)) \leq cr^{Dd+1}, \quad \forall (W, \sigma^2) \in \mathbb{R}^{D\times d} \times [\sigma_0^2, \infty), \ r > 0.$$

When the prior $\Pi$ is uniform w.r.t. the Haar measure,

$$\Pi\left\{\theta : m(\theta, \theta^*) \vee v(\theta, \theta^*) \leq n^{-1/2}\right\}$$

$$\gtrsim \Pi\left\{\theta : \sqrt{\|WW^\top - W^*W^{*\top}\|^2 + (\sigma^2 - \sigma^{*2})^2} \leq cn^{-1/2}\right\} \gtrsim \left(n^{-1/2}\right)^{dD+1}.$$

We conclude that, by (Syring and Martin, 2020, Theorem 3.3), the posterior contraction rate of the Gibbs posterior is $n^{-1/2}$, which is optimal. $\qquad\square$

37