

## A Proof of Theorem 3

*Proof.* We prove by mapping (4) to an equivalent problem by performing a variable change.

Let  $G_S := \bigcup_{i \in S} G_i$ . Note that the inner optimization  $\min_{|G_S| \leq k, q \in \mathcal{F}_{G_S}} \text{KL}(q||p) = -\log p(\mathbf{x}_{G \setminus G_S})$  [Koyejo et al., 2014].

Define the function  $J : \mathfrak{p}(r) \rightarrow \mathbb{R}$  as  $J(S) := \log p(\mathbf{x}_{G \setminus G_S} = 0)$ , and the function  $\tilde{J} : \mathfrak{p}(r) \rightarrow \mathbb{R}$  as  $\tilde{J}(S) := J(S) - J(\emptyset)$ .

Define the costs associated with picking  $G_i$  as  $c_i = |G_i| \forall i \in [r]$ . The cost function of a set  $s \subset G$  can thus be written as  $c(s) := \sum_{\forall i \text{ s.t. } G_i \in s} c_i$ . The optimization problem 4 is then equivalent to  $\max_{\sum_{i \in S} c_i \leq k} \tilde{J}(S)$ .

The result follows from Theorem 1.  $\square$

## B Extension to multiple factors

The factor models defined in Sections 4.3,4.2 are readily extensible to multiple factors using deflation techniques, either by information projection or simple subtraction. See the work of [Khanna et al., 2017] for further discussion. An alternative is also to infer for all factors as is typically done in classic non-sparse settings. We present the detailed equations and explicitly derive EM algorithm and the relevant equations next for the case of the sparse PCA model (Section 4.2). Note that because of the partition constraint construction described in Section 4.3, the equations are also relevant as the EM algorithm for the probabilistic CCA. The only difference is in the Estep, where Algorithm 2 is to be used instead of 3.

We have  $n$  observations made in a  $d$  dimensional space, which are stacked in a matrix  $\mathbf{T} \in \mathbb{R}^{n \times d}$ . Drawing analogy from traditional PCA, we want to search for a few sparse basis vectors whose linear combination generates the observation matrix with small error. Additionally, the basis vectors themselves may have structure.

We model the observation matrix as a product of a parameter  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and a *sparse*  $\mathbf{W} \in \mathbb{R}^{p \times d}$ . The sparse basis vectors are stacked as rows of  $\mathbf{W}$ , and their linear combination is modelled by  $\mathbf{X}$ . We are looking at scenarios with  $n \gg d$ , so the above factorization is useful for small  $p$ , which is set according to the domain.  $\boldsymbol{\mu}$  is the matrix of column means generated as,  $\boldsymbol{\mu} = \text{columnMeans}(\mathbf{T})^\dagger \otimes \mathbf{1}$ , and gaussian noise is represented by  $\epsilon_{ij} \sim \mathcal{N}(0, \sigma^2), \forall i \in [n], \forall j \in [d]$ .

$$\mathbf{T} = \mathbf{X}\mathbf{W} + \boldsymbol{\mu} + \boldsymbol{\epsilon} \quad (6)$$

We have a normal prior on each row of  $\mathbf{W}$ .  $\forall i \in [p], \mathbf{W}_{i,\cdot} \sim \mathcal{N}(0, \mathbf{C})$ , for a known matrix  $\mathbf{C}$ , while the rows are in-

dependent. In the matrix-variate normal form,  $\mathbf{W} \sim \text{MVN}(0, \mathbf{C}, \mathbf{I})$ . For convenience that will be apparent in coming pages, the above equation can be vectorized and rewritten as

$$\vec{\mathbf{T}} = (\mathbf{I} \otimes \mathbf{X})\vec{\mathbf{W}} + \vec{\boldsymbol{\mu}} + \vec{\boldsymbol{\epsilon}} \quad (7)$$

Inference and learning can be performed for this model by an Expectation Maximization (EM) algorithm, We introduce sparsity into  $\mathbf{W}$  by constraining its support in the variational E-step as detailed next.

## C EM Algorithm

Neal and Hinton [1998] give a free energy function based interpretation of the EM algorithm wherein the E-step maximizes the energy function  $\mathcal{F}$  in the space of distributions of the missing data and the M-step maximizes it in the parameter space. In our model,  $\mathbf{X}$  and  $\sigma^2$  are the parameters, and  $\mathbf{W}$  can be treated as the missing data for the EM algorithm. With  $\text{KL}(\cdot, \cdot)$  as the Kullback-Liebler distance, from Neal and Hinton [1998], we have

$$\mathcal{F}(q(\mathbf{W}), \mathbf{X}, \sigma^2) = -\text{KL}(q(\mathbf{W})||p(\mathbf{W}|\mathbf{T}; \mathbf{X}, \sigma^2)) + \log P(\mathbf{T}; \mathbf{X}, \sigma^2). \quad (8)$$

E-step can then be viewed as the search for the best  $q(\mathbf{W})$ ,

$$\text{E-step: } \max_q \mathcal{F}(q(\mathbf{W}), \mathbf{X}, \sigma^2). \quad (9)$$

The M-step is the search for the best parameters,

$$\text{M-step: } \max_{\mathbf{X}, \sigma^2} \mathcal{F}(q(\mathbf{W}), \mathbf{X}, \sigma^2). \quad (10)$$

This view of the EM algorithm provides the flexibility to design algorithms with any E and M steps that monotonically increase  $\mathcal{F}$ .

### C.1 Variational E-step

Equation 9 finds the best  $q$  to maximize the free energy function  $\mathcal{F}$ . An unconstrained search returns the posterior  $p(\mathbf{W}|\mathbf{T}; \mathbf{X}, \sigma^2)$ , making the KL term 0.  $\mathbf{W}$  is size  $p \times d$ , so doing a full E-step is costly. To introduce sparsity, and to make the Estep more tractable, we present a variational E-step. Variational methods perform the search for best  $q$  over a constrained set. For sparsity, we introduce the constraint is that  $q(\vec{\mathbf{W}})$  is  $k$ -sparse i.e. it has support only on at most  $k$  out of total  $n \times d$  dimensions. From Equations 9 and 8, it follows that the variational E-step can be written using the vectorized  $\mathbf{W}$  as:

$$\min_{|q|=k} \text{KL}\left((q(\vec{\mathbf{W}})||p(\vec{\mathbf{W}}|\mathbf{T}; \mathbf{X}, \sigma^2))\right). \quad (11)$$

Let  $K$  be the enumerative set of constraints (group sparsity is a special case, but the equation is valid for other constraint sets as well). Expanding the KL equation, it is straightforward to see that the Equation 11 is equivalent to

$$\max_{\text{Supp}(P) \in K} \log(p(\vec{W}_{K^c} = \mathbf{0}_{K^c} | \mathbf{T}; \mathbf{X}, \sigma^2)) \quad (12)$$

Equation 12 is the resulting optimization problem to be solved for variational E-step. However, it requires a combinatorial search over the set  $k$  which is of size  $\frac{(p*d)!}{(p*d-k)!k!}$ , where  $p*d$  is dimension of the ambient space of  $\vec{W}$ . Searching over *atmost*  $k$ -sparse supports (instead of *exactly*  $k$ -sparse) is more costly as the size of the set to be searched over increases. Since this is to be done in every iteration of the EM, it is prohibitively expensive. As a workaround, we proposed selecting each of the  $k$  dimensions greedily.

### C.1.1 E-step Equations

Having defined the optimization function, the algorithm and its gaurantees, we now derive explicit equations for the variational E-step.

For deriving the E-step equations, we first require the posterior  $p(\vec{W} | \vec{T}; \mathbf{X}, \sigma^2)$ . In this section, we suppress  $\{\mathbf{X}, \sigma^2\}$  from the equations for brevity since dependence on them is obvious. From standard properties of the gaussians, we know that,

$$\begin{aligned} p(\vec{W} | \vec{T}) &\sim \mathcal{N}(\mathbf{m}, \Sigma) \\ \text{where,} \\ \Sigma^{-1} &= \frac{1}{\sigma^2} (\mathbf{I}_d \otimes \mathbf{X}^\dagger \mathbf{X}) + (\mathbf{C}^{-1} \otimes \mathbf{I}_p), \\ \mathbf{m} &= \Sigma \frac{1}{\sigma^2} (\mathbf{I}_d \otimes \mathbf{X}^\dagger) (\vec{T} - \vec{\mu}). \end{aligned} \quad (13)$$

We can now expand Equation 12 as

$$\max_K \mathbf{m}_{K^c}^\dagger [\Sigma_{K^c}]^{-1} \mathbf{m}_{K^c} - \log \det \Sigma_{K^c} + \text{const},$$

which is equivalent to:

$$\begin{aligned} \text{define } \mathbf{r} &= \Sigma^{-1} \mathbf{m}; \\ \max_K \mathbf{r}_K & [[\Sigma^{-1}]_K]^{-1} \mathbf{r}_K - \log \det [\Sigma^{-1}]_K. \end{aligned} \quad (14)$$

Recall that  $\Sigma_K$  is the submatrix of  $\Sigma$  supported on  $K$ , similarly for  $[\Sigma^{-1}]_K$ . Note that Equation 14 is in space of size  $p*d$ . Since  $[\Sigma^{-1}]_K$  and  $\mathbf{r}_K$  are easy to calculate by some smart indexing, this version of the equation is convenient to implement and use.

Say after solving the constrained optimization problem specified by Equations [13,14], we obtain the best support as  $K^*$ ,

then the resulting solution density, say  $q^*$ , which is known to be the conditional Koyejo et al. [2014], is

$$\begin{aligned} q^* &\sim \mathcal{N}(\mathbf{c}, \mathbf{D}) \\ \text{where,} \\ \mathbf{D}^{-1} &= [\Sigma^{-1}]_{K^*}, \\ \mathbf{c} &= \mathbf{D} \mathbf{r}_{K^*} \end{aligned} \quad (15)$$

Recall,  $q^*$  has support only on  $K^*$ , so in Equation 15,  $\mathbf{c} \in \mathbb{R}^{|K^*|}$ ,  $\mathbf{D} \in \mathbb{R}^{|K^*| \times |K^*|}$ .

## C.2 M-step

Equation 10 optimizes  $\mathcal{F}$  in the parameter space to get the argmax parameters that maximize  $\mathcal{F}$ . However, it turns out solving for  $\{X, \sigma^2\}$  over  $\mathcal{F}$  directly is costly. Since the free energy view of the EM shows that any M-step that increases  $\mathcal{F}$  suffices, we maximize the log likelihood portion of  $\vec{\mathcal{F}}$  instead for the M-step. Recall  $q^*$  is the distribution on  $\vec{W}$  obtained from the E-step, the effective M-step is:

$$\max_{\mathbf{X}, \sigma^2} \mathbb{E}_{q^*} [\log p(\vec{T} | \vec{W}; \mathbf{X}, \sigma^2)] \quad (16)$$

### C.2.1 M-step Equations

Since  $q^*$  has measure 0 outside  $K^*$ , let  $\hat{\mathbf{c}}$  represent the mean vector  $\mathbf{c}$  expanded from  $|K^*|$  to ambient dimension  $p*d$ , with zeroes padded as needed. Similarly,  $\hat{\mathbf{D}} \in \mathbb{R}^{p*d \times p*d}$  represents  $\mathbf{D}$  with zeroes padded as needed. For brevity, the following equations assume  $\vec{T}$  to be zero mean i.e.  $\vec{\mu} = 0$ , the derivation extends to non-zero mean  $\vec{T}$  trivially by replacing  $\vec{T}$  with  $\vec{T} - \vec{\mu}$ . Equation 16 can be written as:

$$\begin{aligned} \max_{\mathbf{X}, \sigma^2} \mathbb{E}_{q^*} & \left[ \frac{-1}{2\sigma^2} (\vec{T} - (\mathbf{I} \otimes \mathbf{X}) \vec{W})^\dagger (\vec{T} - (\mathbf{I} \otimes \mathbf{X}) \vec{W}) \right. \\ & \left. - nd \log \sigma^2 \right] \\ \equiv \max_{\mathbf{X}, \sigma^2} & \frac{-1}{2\sigma^2} \mathcal{V}(\mathbf{X}) - nd \log \sigma^2, \end{aligned}$$

where,

$$\mathcal{V}(\mathbf{X}) = \mathbf{r} (\mathbf{I} \otimes \mathbf{X}^\dagger \mathbf{X}) (\widehat{\mathbf{c}} \widehat{\mathbf{c}}^\dagger + \widehat{\mathbf{D}}) - 2 \widehat{\mathbf{c}}^\dagger (\mathbf{X}^\dagger \mathbf{T})$$

Clearly,  $\mathbf{X}$  and  $\sigma^2$  can be updated separately.

For  $\mathbf{X}$ , it is easy to take gradient of  $\mathcal{V}(\mathbf{X})$  w.r.t  $\mathbf{X}$ . We can matricize  $\widehat{\mathbf{c}}$  as  $\text{mat}(\widehat{\mathbf{c}}) \in \mathbb{R}^{p \times d}$  to rewrite the second term as a trace:  $\text{rmat}(\widehat{\mathbf{c}})^\dagger \mathbf{X}^\dagger \mathbf{T}$ . Note that  $\mathbf{I} \otimes \mathbf{X}^\dagger \mathbf{X}$  is a block diagonal matrix with the same block  $\mathbf{X}^\dagger \mathbf{X}$  repeating over and over. Thus, we can define  $\mathbf{M} = \sum_{\text{blocks}} \widehat{\mathbf{c}} \widehat{\mathbf{c}}^\dagger + \widehat{\mathbf{D}}$ , where the summation is over diagonal blocks of size  $p \times p$

(sp  $\mathbf{M} \in \mathbb{R}^{p \times p}$ ), to write  $\mathcal{V}(\mathbf{X})$  and its gradient as

$$\begin{aligned} \mathcal{V}(\mathbf{X}) &= \mathbf{r}\mathbf{X}^\dagger\mathbf{X}\mathbf{M} - \mathbf{r}\mathbf{mat}(\hat{\mathbf{c}})^\dagger\mathbf{X}^\dagger\mathbf{T} \\ \frac{\partial \mathcal{V}}{\partial \mathbf{X}} &= 2\mathbf{X}\mathbf{M} - \mathbf{T}\mathbf{mat}(\hat{\mathbf{c}})^\dagger \end{aligned} \quad (17)$$

Equation 17 gives closed form solution if  $\mathbf{M}$  is invertible. If that is not the case, gradient steps can be taken to update  $\mathbf{X}$ . Recall that even a single gradient step to update  $\mathbf{X}$  suffices as required by the free energy view of the EM.

Once  $\mathbf{X}$  is updated, updating  $\sigma^2$  is more straightforward by taking derivative and setting to 0 to get

$$\sigma^2 = \frac{\mathcal{V}(\mathbf{X})}{2nd} \quad (18)$$

---

**Algorithm 5:** EM Algorithm for sparse projections

---

- 1: **Input:**  $k, p, d, \mathbf{C}, \mathbf{T}$
  - 2: Initialize  $\mathbf{X}$  randomly
  - 3: **while** not converged **do**
  - 4:   **E-Step**
  - 5:   Init:  $\mathbf{K}^* = \{\}$
  - 6:   **for**  $i = 1 \dots k$  **do**
  - 7:     Solve Equation 14 with  
 $\mathbf{K} = \mathbf{H} \cup \{j\}, \forall j \in [p * d], j \notin \mathbf{K}^*$
  - 8:     Set  $\mathbf{K}^* \cup \{j^*\}$ , where  $j^*$  is argmax from  
previous step
  - 9:   **end for**
  - 10:   Use Equation 15 to get  $q^*$
  - 11:   **M-Step**
  - 12:   Solve for  $\mathbf{X}$  using Equation 17
  - 13:   Solve for  $\sigma^2$  using Equation 18
  - 14: **end while**
  - 15: return( $q^*, \mathbf{X}, \sigma^2$ )
-