

Kruskal-Factor Analysis Supplementary Material

1 Gibbs updates

KFA Likelihood and priors:

$$\begin{aligned}
 \mathcal{X}_n &\sim \mathcal{N}(\mathcal{D} \times_{T+1} (\mathbf{s}_n \circ \mathbf{z}_n), \gamma_\epsilon^{-1} \mathbf{I}), \quad d_{ik} = \sum_{r=1}^R \lambda_{rk} \prod_{t=1}^T u_{itr}^{(kt)} \\
 u_{itr}^{(kt)} &\sim \mathcal{N}(0, m_t^{-1}), \quad \lambda_{kr} \sim \mathcal{N}(0, \tau_{kr}^{-1}), \quad \tau_{kr} = \prod_{i=1}^r \delta_{ki} \\
 s_{kn} &\sim \mathcal{N}(0, \gamma_s^{-1}), \quad z_{kn} \sim \text{Bernoulli}(\pi_k) \\
 \pi_k &\sim \text{Beta}(a_\pi/K, b_\pi(K-1)/K), \quad \delta_{ki} \sim \text{Gam}(\alpha, 1)
 \end{aligned} \tag{1}$$

where $\mathcal{X}_n \in \mathbb{R}^{m_1 \times \dots \times m_T}$ is the n^{th} data item, and d_{ik} is the $i = [i_1, \dots, i_T]$ element of the k^{th} atom. The precisions $\gamma_\epsilon, \gamma_s$ have gamma priors.

$$\begin{aligned}
 d_{ik} &= \lambda_{rk} \prod_{t' \neq t} u_{itr'}^{(kt')} u_{itr}^{(kt)} + \sum_{r' \neq r} \lambda_{r'k} \prod_{t=1}^T u_{itr'}^{(kt)} = a_{it} u_{itr}^{(kt)} + b_{it} \\
 &= \prod_{t=1}^T u_{itr}^{(kt)} \lambda_{rk} + \sum_{r' \neq r} \lambda_{r'k} \prod_{t=1}^T u_{itr'}^{(kt)} = f_i \lambda_{rk} + g_i
 \end{aligned} \tag{2}$$

$$\begin{aligned}
 u_{itr}^{(kt)} &\sim \mathcal{N}(\hat{\mu}, \hat{\omega}^{-1}) \\
 \hat{\omega} &= m_t + \sum_{i: i_t=m} \gamma_{i_t k}^{\mathcal{D}_k} (a_{i_t})^2 \\
 \hat{\mu} &= \frac{1}{\hat{\omega}} \sum_{i: i_t=m} (\mu_{i_t k}^{\mathcal{D}_k} - b_{i_t}) a_{i_t} \\
 \gamma^{\mathcal{D}_k} &= \gamma_\epsilon \mathbf{I}_P \sum_{n=1}^N s_{kn}^2 z_{kn}^2 \\
 \mu^{\mathcal{D}_k} &= \gamma_\epsilon \sum_{n=1}^N s_{kn} z_{kn} \tilde{\mathbf{x}}_n^{\setminus k} \\
 \tilde{\mathbf{x}}_n^{\setminus k} &= \text{vec} \left[\mathcal{X}_n - \mathcal{D} \times_{T+1} (\mathbf{s}_n \circ \mathbf{z}_n) + s_{kn} z_{kn} \mathcal{D}_k \right]
 \end{aligned} \tag{3}$$

$$\begin{aligned}
\lambda_{kr} &\sim \mathcal{N}(\hat{\nu}, \hat{\tau}^{-1}) \\
\hat{\tau} &= \tau_{kr} + \gamma_\epsilon \sum_{n,i} (f_i s_{kn})^2 \\
\hat{\nu} &= \frac{\gamma_\epsilon}{\hat{\tau}} \sum_{n,i} (h_{in} - g_i s_{kn}) f_i s_{kn} \\
h_{in} &= \mathcal{X}_{in} - [\mathcal{D} \times_{T+1} (\mathbf{s}_n \circ \mathbf{z}_n)]_i
\end{aligned} \tag{4}$$

$$\begin{aligned}
p(s_{ki}|-) &\sim \mathcal{N}(\mu_{s_{ki}}, \Sigma_{s_{ki}}) \\
\Sigma_{s_{ki}} &= (\gamma_s + \gamma_\epsilon z_{kn}^2 \mathbf{d}_k^T \mathbf{d}_k)^{-1} \\
\mu_{s_{ki}} &= \gamma_\epsilon z_{kn} \mathbf{d}_k^T \tilde{\mathbf{x}}_n^{\setminus k} \\
\mathbf{d}_k &= \text{vec}[\mathcal{D}_k]
\end{aligned} \tag{5}$$

$$\begin{aligned}
z_{ki} &\sim \text{Bern}\left(\frac{p_1}{p_0 + p_1}\right) \\
p_1 &= \pi_k \exp\left[-\frac{\gamma_\epsilon}{2} (s_{kn}^2 \mathbf{d}_k^T \mathbf{d}_k - s_{kn} \mathbf{d}_k^T \tilde{\mathbf{x}}_n^{\setminus k})\right] \\
p_0 &= 1 - \pi_k
\end{aligned} \tag{6}$$

$$p(\gamma_\epsilon|-) \sim \Gamma\left(c + \frac{PN}{2}, d + \frac{1}{2} \sum_{n=1}^N \|\mathcal{X}_n - \mathcal{D} \times_{T+1} (\mathbf{s}_n \circ \mathbf{z}_n)\|_F^2\right) \tag{7}$$

$$p(\gamma_s|-) \sim \Gamma\left(e + \frac{KN}{2}, f + \frac{1}{2} \sum_{n=1}^N \mathbf{s}_n^T \mathbf{s}_n\right) \tag{8}$$

$$p(\pi_k|-) \sim \text{Beta}\left(\frac{a}{K} + \sum_{n=1}^N z_{kn}, b \frac{K-1}{K} + N - \sum_{n=1}^N z_{kn}\right) \tag{9}$$

2 Bayesian Conditional Density Filter for KFA

We record surrogate conditional sufficient statistics (SCSS) for the global parameters that directly interact with the data or local parameters $\mathcal{D}, \gamma_\epsilon, \gamma_s, \pi_k$. The SCSS are given by the posterior updates of the global parameters. The SCSS are accumulated for each subset of data, and then applied in the subsequent update as the prior parameters. We reset the SCSS after the number of data items seen is the same as the size of the dataset. In our implementation we update the model parameters 3 times using a single subset of data. The SCSS are held fixed until the 3rd iteration. This allows stale local parameters $\{\mathbf{S}, \mathbf{Z}\}$ to be refreshed before updating the SCSS. The algorithm below describes the process.

```

for  $epoch = 1$  to  $numEpoch$  do
  Randomly partition  $\mathcal{X}$  into minibatches  $\mathcal{X}^{(i)}$  of size  $batchSize$ 
  Set SCSS to prior parameters
  for all  $\mathcal{X}^{(i)}$  do
    Update  $\mathbf{S}^{(i)}, \mathbf{Z}^{(i)}$  and all hyperparameters with KFA
    Do a standard KFA update twice (including the dictionary)
    Update SCSS ( $\forall k$ ):
       $\boldsymbol{\mu}_{SCSS}^{\mathcal{D}_k} = \boldsymbol{\mu}_{SCSS}^{\mathcal{D}_k} + \boldsymbol{\mu}^{\mathcal{D}_k(i)}$ 
       $\gamma_{SCSS}^{\mathcal{D}_k} = \gamma_{SCSS}^{\mathcal{D}_k} + \gamma^{\mathcal{D}_k(i)}$ 
       $\pi_{kaSCSS} = \pi_{kaSCSS} + \sum_n z_{kn}^{(i)}$ 
       $\pi_{kbSCSS} = \pi_{kbSCSS} + N - \sum_n z_{kn}^{(i)}$ 
       $\gamma_{saSCSS} = \gamma_{saSCSS} + KN^{(i)}/2$ 
       $\gamma_{sbSCSS} = \gamma_{sbSCSS} + \frac{1}{2} \sum_n \mathbf{s}_n^{(i)\top} \mathbf{s}_n^{(i)}$ 
       $\gamma_{\epsilon aSCSS} = \gamma_{\epsilon aSCSS} + PN^{(i)}/2$ 
       $\gamma_{\epsilon bSCSS} = \gamma_{\epsilon bSCSS} + \frac{1}{2} \sum_n \|\mathcal{X}_n^{(i)} - \mathcal{D}_{T+1} \times (\mathbf{s}_n^{(i)} \circ \mathbf{z}_n^{(i)})\|_F^2$ 
    end for
  Reset SCSS to original prior parameters
  if  $epoch == 1$  then
    Update  $\mathbf{S}, \mathbf{Z}$  and all hyperparameters using the full data set (note that
    every column is independent, so this can be done in parallel).
  end if
end for

```

3 Barbara Images



Figure 1: Barbara reconstructions from a random 20% of the pixels: corrupted image (top), KFA (center), BPFA (bottom). KFA reconstructs the eyes and mouth better than BPFA. Also, the contrast in the stripe pattern is stronger in the KFA reconstruction.

4 Dictionary Recovery

For this example, we assume the dictionary is composed of the $5 \times 5 \times 5 \times 5$ DCT filters. Note that each filter is rank-1 (separable). We generate 500 data samples using normal random weights for \mathcal{S} , enforce that only 8 nonzero elements are present in each column of \mathcal{Z} and set the noise variance to 0.01. We use the same settings for BPFA and KFA, with $R = 1$. The dictionary size is $K = 625$, the same as the number of DCT filters. We ran both models for 100 iterations. Figure 2 shows the learned dictionaries from each algorithm. KFA learns the DCT filters. We find that even when we set R to the expected rank $R_E = 30$, the KFA dictionary atoms are low-rank. The BPFA dictionary looks noisy (high-rank structures) while both of the KFA dictionaries have low-rank atoms.

The RMSE for the three models is given in table 1. We can see from the dictionary and the RMSE that KFA is discovering the multi-way structure of the data. The number of samples is smaller than the dictionary, yet KFA is still able to find a well structured dictionary and achieve a lower error.

Table 1: RMSE on 4D-DCT example

	BPFA	KFA $R = 1$	KFA $R = 30$
RMSE	0.03795	0.01387	0.01095

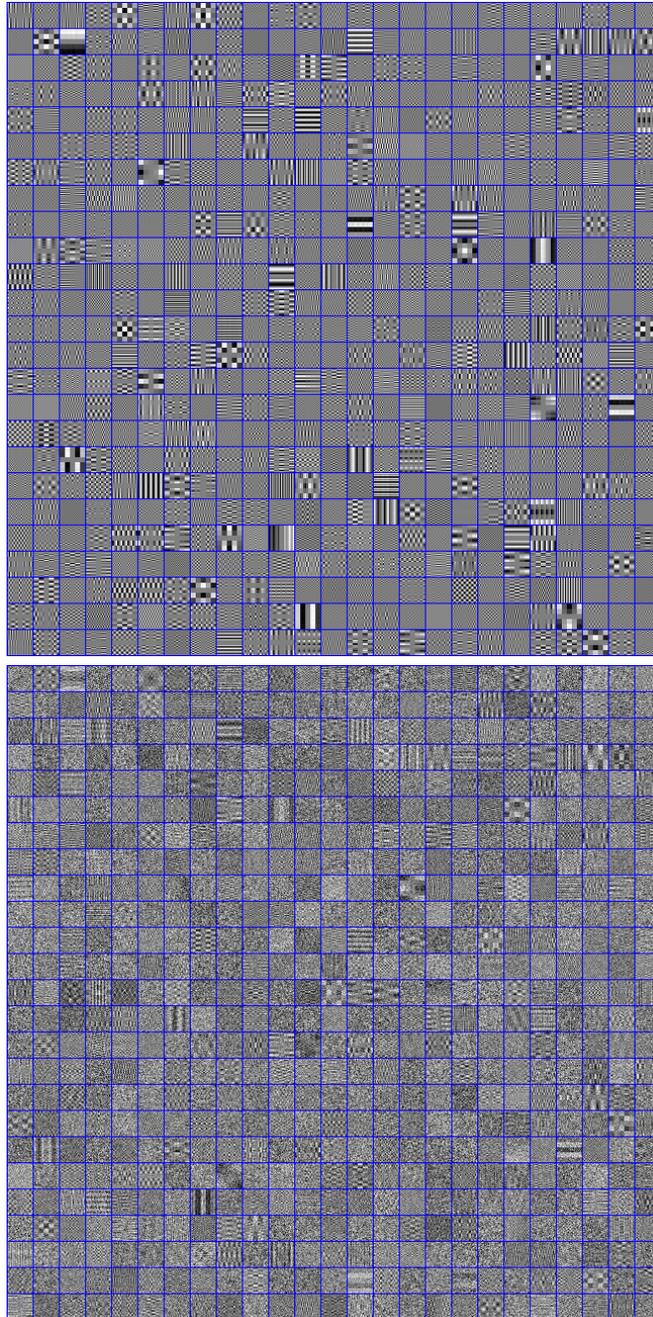


Figure 2: 4D-DCT example dictionaries collapsed to 2D: KFA (top), BPFA (bottom). Zoom-in for fine structure.