

## PCG for the RKHS CP-ALS mode- $k$ subproblem (missing data)

Fix all CP factors except the (possibly infinite-dimensional) RKHS mode  $k$ . Write the mode- $k$  unfolding as  $T \in \mathbb{R}^{n \times M}$  with missing entries set to zero, and let  $S \in \mathbb{R}^{N \times q}$  (with  $N = nM$ ) be the selection matrix so that  $S^\top \text{vec}(T)$  extracts the  $q$  observed entries. Let  $Z \in \mathbb{R}^{M \times r}$  be the Khatri–Rao product of the other CP factors and  $B = TZ$ . Assume the RKHS representer form  $A_k = KW$ , where  $K \in \mathbb{R}^{n \times n}$  is symmetric psd. The ALS subproblem in  $W \in \mathbb{R}^{n \times r}$  is the linear system

$$[(Z \otimes K)^\top SS^\top(Z \otimes K) + \lambda(I_r \otimes K)] \text{vec}(W) = (I_r \otimes K) \text{vec}(B) = \text{vec}(KB), \quad (1)$$

of size  $nr \times nr$ . The goal is to solve (??) without forming the dense matrix and without any  $O(N)$  work, assuming  $n, r < q \ll N$ .

**1. SPD and why PCG applies.** Let  $P \equiv SS^\top$  (a diagonal mask,  $P = P^\top = P^2$ ) and define

$$A \equiv (Z \otimes K)^\top P(Z \otimes K) + \lambda(I_r \otimes K), \quad b \equiv \text{vec}(KB).$$

Then  $A$  is symmetric. If  $K \succ 0$  and  $\lambda > 0$ , for any  $x \neq 0$ ,

$$x^\top Ax = \|P^{1/2}(Z \otimes K)x\|_2^2 + \lambda x^\top(I_r \otimes K)x > 0,$$

so  $A \succ 0$  and (preconditioned) conjugate gradients (PCG) is applicable. If  $K$  is only psd, add a small nugget  $\varepsilon I$  to  $K$  (standard in kernel ridge regression) or reduce to the rank- $m$  eigenspace of  $K$  to obtain an SPD system of size  $mr$ .

## 2. Matvecs in $O(n^2r + qr)$ using gather/scatter

Write  $x \in \mathbb{R}^{nr}$  as  $x = \text{vec}(X)$  with  $X \in \mathbb{R}^{n \times r}$  (column-stacked). Use the identity

$$(Z \otimes K) \text{vec}(X) = \text{vec}(KXZ^\top), \quad (2)$$

so that the action of  $(Z \otimes K)$  is “form the prediction matrix”  $U \equiv KXZ^\top \in \mathbb{R}^{n \times M}$ .

**Observed index list.** Store the  $q$  observed indices in unfolding coordinates as pairs  $(i_t, j_t)$ ,  $t = 1, \dots, q$ . Then  $S^\top \text{vec}(U) = (U_{i_t, j_t})_{t=1}^q$  (gather). Conversely, for  $u \in \mathbb{R}^q$ , the vector  $Su \in \mathbb{R}^N$  has entries  $(Su)_{(i_t, j_t)} = u_t$  and zeros elsewhere; reshaping  $Su$  into  $n \times M$  via  $\text{vec}^{-1}$  gives the corresponding sparse matrix (scatter). These operations cost  $O(q)$  given the index arrays.

**Matvec**  $y = Ax$ . Given  $X$ :

$$1. \Gamma \leftarrow KX \quad (O(n^2r)).$$

2. For each observation  $t$  compute a row vector  $z_t \equiv Z_{j_t,:} \in \mathbb{R}^r$  and the scalar

$$u_t \leftarrow \langle \Gamma_{i_t,:}, z_t \rangle. \quad (3)$$

3. Accumulate  $H \in \mathbb{R}^{n \times r}$  via

$$H_{i_t,:} += u_t z_t, \quad t = 1, \dots, q. \quad (4)$$

4. Output  $\text{vec}(KH + \lambda\Gamma)$ .

To see correctness:  $U = KXZ^\top$  implies  $P\text{vec}(U) = S(S^\top\text{vec}(U)) = Su$  where  $u_t = U_{i_t,j_t}$ . Let  $\tilde{U} \in \mathbb{R}^{n \times M}$  be the reshape of  $Su$ , so  $\text{vec}(\tilde{U}) = P\text{vec}(U)$  and  $\tilde{U}_{i_t,j_t} = u_t$ . The adjoint Kronecker identity gives

$$(Z \otimes K)^\top \text{vec}(\tilde{U}) = \text{vec}(K\tilde{U}Z),$$

and the sparse accumulation (??) computes  $H = \tilde{U}Z$  without ever materializing  $\tilde{U}$ . Thus the first term equals  $\text{vec}(KH)$ , and adding  $\lambda(I_r \otimes K)x = \text{vec}(\lambda KX) = \text{vec}(\lambda\Gamma)$  yields  $y = Ax$ .

**Avoiding explicit  $Z$  (avoiding  $M$  and  $N$ ).** Although  $Z$  is of size  $M \times r$ , we never form it. Given an observed tensor multi-index  $(i_1^{(t)}, \dots, i_d^{(t)})$ , the required row is

$$z_t = A_d(i_d^{(t)}, :) \odot \cdots \odot A_{k+1}(i_{k+1}^{(t)}, :) \odot A_{k-1}(i_{k-1}^{(t)}, :) \odot \cdots \odot A_1(i_1^{(t)}, :),$$

computable on the fly in  $O((d-1)r)$  time. If memory allows, cache all  $z_t$  once in a  $q \times r$  array to make each PCG iteration cost  $O(qr)$  for the sparse part.

**RHS.** Compute  $B = TZ$  without forming  $T$ : for each observed value  $y_t$  at  $(i_t, j_t)$ , do  $B_{i_t,:} += y_t z_t$  (same sparse accumulation as above), then set  $b = \text{vec}(KB)$ . Cost:  $O(qr + n^2r)$  (or  $O(qdr + n^2r)$  if computing  $z_t$  on the fly).

### 3. A Kronecker preconditioner and fast application

A convenient SPD preconditioner replaces  $P$  by a scaled identity  $\alpha I$ . Under uniform sampling,  $\mathbb{E}[P] = (q/N)I$  so a natural default is  $\alpha = q/N$ . More generally one can choose  $\alpha$  by matching traces,

$$\alpha := \frac{\text{Tr}((Z \otimes K)^\top P(Z \otimes K))}{\text{Tr}((Z^\top Z) \otimes K^2)} = \frac{\sum_{t=1}^q \|K_{:,i_t}\|_2^2 \|z_t\|_2^2}{\text{Tr}(K^2) \text{Tr}(Z^\top Z)}, \quad (5)$$

which is computable from the observed indices in  $O(qr + n^2)$  time (or  $O(qr)$  if  $\|K_{:,i_t}\|_2^2$  are precomputed).

With this scalar approximation,

$$A_0 \equiv \alpha(Z^\top Z) \otimes (K^2) + \lambda(I_r \otimes K). \quad (6)$$

When there are no missing entries ( $P = I$ ) and  $\alpha = 1$ ,  $A_0$  equals  $A$  exactly.

**Computing  $\Phi = Z^\top Z$  without forming  $Z$ .** With  $Z = A_d \odot \cdots \odot A_{k+1} \odot A_{k-1} \odot \cdots \odot A_1$ , the Gram matrix satisfies the standard Khatri–Rao identity

$$\Phi \equiv Z^\top Z = \underset{i \neq k}{\circ} (A_i^\top A_i),$$

where  $\circ$  denotes the entrywise (Hadamard) product. This costs  $O(\sum_{i \neq k} n_i r^2)$  and is already computed in many CP-ALS implementations.

**Applying  $A_0^{-1}$ .** Let  $K = U\Lambda U^\top$  and  $\Phi = V\Sigma V^\top$ . Then  $A_0$  diagonalizes in the Kronecker basis  $(V \otimes U)$ , so for  $x = \text{vec}(X)$ ,

$$\hat{X} \leftarrow U^\top X V, \quad \hat{X}_{b,a} \leftarrow \hat{X}_{b,a}/(\alpha\sigma_a\lambda_b^2 + \lambda\lambda_b), \quad X \leftarrow U\hat{X}V^\top.$$

Each application costs  $O(n^2r + nr^2)$  after one-time eigendecompositions ( $O(n^3 + r^3)$ ).

**Cheaper block-diagonal alternative.** Replacing  $\Phi$  by  $\text{diag}(\Phi)$  yields  $r$  independent column-wise preconditioners: if  $X = [x_1, \dots, x_r]$  then  $(A_0^{\text{bd}})^{-1}$  applies

$$x_a \mapsto (\alpha \Phi_{aa} K^2 + \lambda K)^{-1} x_a = (K(\alpha \Phi_{aa} K + \lambda I))^{-1} x_a,$$

which can be done using only the eigendecomposition of  $K$  in  $O(n^2)$  per column.

## 4. Complexity and iteration count (no $O(N)$ terms)

Let  $m$  be the PCG iteration count. Per iteration:

$$\text{matvec } Ax : O(n^2r + qr) \text{ (or } O(n^2r + qdr)), \quad \text{preconditioner } A_0^{-1} : O(n^2r + nr^2).$$

Thus the PCG solve costs  $O(m(n^2r + qr + nr^2))$  time and  $O(nr + q)$  memory (plus optional  $O(qr)$  cache). One-time setup (per outer ALS sweep) includes forming  $\Phi$  in  $O(\sum_{i \neq k} n_i r^2)$  and eigendecompositions in  $O(n^3 + r^3)$ , which are amortized over the  $m$  iterations.

**How  $m$  depends on the preconditioner.** Standard PCG theory gives for the  $m$ th iterate  $x_m$ :

$$\frac{\|x_m - x_\star\|_A}{\|x_0 - x_\star\|_A} \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^m, \quad \kappa \equiv \kappa(A_0^{-1} A),$$

so  $m = O(\sqrt{\kappa} \log(1/\varepsilon))$  iterations suffice for relative  $A$ -norm error  $\varepsilon$ . Moreover, if  $\|A_0^{-1/2}(A - A_0)A_0^{-1/2}\|_2 \leq \eta < 1$  then the eigenvalues of  $A_0^{-1} A$  lie in  $[1 - \eta, 1 + \eta]$  and

$$\kappa(A_0^{-1} A) \leq \frac{1 + \eta}{1 - \eta}.$$

Under near-uniform sampling and bounded ‘‘leverage’’ of the observed rank-one terms, one can bound  $\eta$  with high probability using matrix Bernstein/Chernoff inequalities, implying  $\kappa$  (and hence  $m$ ) is  $O(1)$  once  $q$  is moderately larger than  $nr$  (up to log factors).