

## Supplemental Material

### A Proofs

#### A.1 On the connection of ONMF with NNPCA

**Lemma 2.** Let  $\mathcal{E}_* \triangleq \min_{\mathbf{W} \geq \mathbf{0}, \mathbf{H} \geq \mathbf{0}, \mathbf{W}^\top \mathbf{W} = \mathbf{I}_k} \|\mathbf{M} - \mathbf{W}\mathbf{H}^\top\|_F^2$  be the optimal ONMF approximation error for a given  $m \times n$  real nonnegative matrix  $\mathbf{M}$  and target dimension  $k$ , as defined in (1). Then,

$$\mathcal{E}_* = \|\mathbf{M}\|_F^2 - \max_{\substack{\mathbf{W} \geq \mathbf{0}_{m \times k} \\ \mathbf{W}^\top \mathbf{W} = \mathbf{I}_k}} \|\mathbf{M}^\top \mathbf{W}\|_F^2, \quad (6)$$

If  $\mathbf{W}_*$  is a solution of the maximization in (6), then the pair  $\mathbf{W}_*, \mathbf{H}_* \triangleq \mathbf{M}^\top \mathbf{W}_*$  is a feasible solution to the ONMF problem in (1), achieving the minimum error  $\mathcal{E}_*$ , i.e.,  $\|\mathbf{M} - \mathbf{W}_* \mathbf{W}_*^\top \mathbf{M}\|_F^2 = \mathcal{E}_*$ .

*Proof.* Recall that by assumption,  $\mathbf{W}$  is a  $m \times k$  nonnegative matrix with orthonormal columns. The subsequent analysis holds even in the case where  $\mathbf{W}$  is allowed to contain all-zero columns, as such columns do not contribute to the objective function and can be ignored effectively reducing the dimension  $k$  of the factorization.

Given a real nonnegative  $m \times k$  matrix  $\mathbf{W}$ , the real  $n \times k$  matrix  $\mathbf{H}$  that minimizes the Frobenius error  $\|\mathbf{M} - \mathbf{W}\mathbf{H}^\top\|_F^2$  over all real  $n \times k$  matrices (ignoring temporarily the fact what we seek a nonnegative  $\mathbf{H}$ ), is given by  $\mathbf{H}^\top = \mathbf{W}^\dagger \mathbf{M}$ , where  $\mathbf{W}^\dagger$  denotes the pseudo-inverse of  $\mathbf{W}$ . Here, however, the columns of  $\mathbf{W}$  are orthonormal and hence  $\mathbf{W}^\dagger = \mathbf{W}^\top$ . Moreover, since  $\mathbf{M}$  is nonnegative,  $\mathbf{H}^\top = \mathbf{W}_*^\dagger \mathbf{M} = \mathbf{W}_*^\top \mathbf{M}$  automatically satisfies the additional nonnegativity constraint. Therefore, the ONMF problem (defined in (1)) reduces to a minimization in a single variable:

$$\mathcal{E}_* = \min_{\substack{\mathbf{W} \geq \mathbf{0}_{m \times k} \\ \mathbf{W}^\top \mathbf{W} = \mathbf{I}_k}} \|\mathbf{M} - \mathbf{W}\mathbf{W}^\top \mathbf{M}\|_F^2. \quad (7)$$

Expanding the objective in (7),

$$\begin{aligned} \|\mathbf{M} - \mathbf{W}\mathbf{W}^\top \mathbf{M}\|_F^2 &= \|\mathbf{M}\|_F^2 - 2 \cdot \text{Tr}(\mathbf{W}^\top \mathbf{M} \mathbf{M}^\top \mathbf{W}) + \text{Tr}(\mathbf{M}^\top \mathbf{W} \mathbf{W}^\top \mathbf{M}) \\ &= \|\mathbf{M}\|_F^2 - \text{Tr}(\mathbf{W}^\top \mathbf{M} \mathbf{M}^\top \mathbf{W}) \\ &= \|\mathbf{M}\|_F^2 - \|\mathbf{M}^\top \mathbf{W}\|_F^2, \end{aligned} \quad (8)$$

where the first step follows from the cyclic property of the trace and the fact that  $\mathbf{W}^\top \mathbf{W} = \mathbf{I}_k$ . This concludes the proof.  $\square$

#### A.2 Proof of Lemma 1

**Lemma 1.** For any real  $m \times n$  matrix  $\bar{\mathbf{M}}$  with rank  $r$ , desired number of components  $k$ , and accuracy parameter  $\epsilon \in (0, 1)$ , Algorithm 1 outputs  $\bar{\mathbf{W}} \in \mathcal{W}_k$  such that

$$\|\bar{\mathbf{M}}^\top \bar{\mathbf{W}}\|_F^2 \geq (1 - \epsilon) \cdot \|\bar{\mathbf{M}}^\top \bar{\mathbf{W}}_*\|_F^2,$$

where  $\bar{\mathbf{W}}_*$  is the optimal solution defined in (3), in time  $T_{\text{SVD}} + O\left(\left(\frac{2}{\epsilon}\right)^{r \cdot k} \cdot k \cdot m\right)$ .

*Proof.* Let  $\bar{\mathbf{M}} = \bar{\mathbf{U}} \bar{\Sigma} \bar{\mathbf{V}}^\top$  denote the truncated eigenvalue decomposition of  $\bar{\mathbf{M}}$ ;  $\bar{\Sigma}$  is a diagonal  $r \times r$  matrix with  $\Sigma_{ii}$  being equal to the  $i$ th largest singular value of  $\bar{\mathbf{M}}$ . For any  $\mathbf{w} \in \mathbb{R}^m$ ,

$$\|\bar{\mathbf{M}}^\top \mathbf{w}\|_2^2 = \|\bar{\mathbf{V}} \bar{\Sigma} \bar{\mathbf{U}}^\top \mathbf{w}\|_2^2 = \|\bar{\Sigma} \bar{\mathbf{U}}^\top \mathbf{w}\|_2^2 \geq \langle \bar{\Sigma} \bar{\mathbf{U}}^\top \mathbf{w}, \mathbf{c} \rangle^2, \quad \forall \mathbf{c} \in \mathbb{R}^r : \|\mathbf{c}\|_2 = 1, \quad (9)$$

where the first equality follows from the fact that the columns of  $\bar{\mathbf{V}}$  are orthonormal and span the entire row space of  $\bar{\mathbf{M}}$ , and the inequality is due to Cauchy-Schwartz. In fact, equality is achieved for  $\mathbf{c}$  colinear to  $\bar{\Sigma} \bar{\mathbf{U}}^\top \mathbf{w}$ , appropriately scaled to unit-length, and hence,

$$\|\bar{\mathbf{M}}^\top \mathbf{w}\|_2^2 = \max_{\mathbf{c} \in \mathbb{S}^{r-1}} \langle \bar{\Sigma} \bar{\mathbf{U}}^\top \mathbf{w}, \mathbf{c} \rangle^2. \quad (10)$$

In turn,

$$\|\overline{\mathbf{M}}^\top \mathbf{W}\|_{\mathbb{F}}^2 = \sum_{j=1}^k \|\overline{\mathbf{M}}^\top \mathbf{w}_j\|_2^2 = \sum_{j=1}^k \max_{\mathbf{c}_j \in \mathbb{S}^{r-1}} \langle \overline{\Sigma} \overline{\mathbf{U}}^\top \mathbf{w}_j, \mathbf{c}_j \rangle^2. \quad (11)$$

Recall that  $\overline{\mathbf{W}}_*$  by definition maximizes the left hand side of (11) over all  $\mathbf{W} \in \mathcal{W}_k$ . Let  $\tilde{\mathbf{c}}_{*1}, \dots, \tilde{\mathbf{c}}_{*k} \in \mathbb{S}^{r-1}$  be the set of  $k$  vectors achieving equality in (11) for  $\mathbf{W} = \overline{\mathbf{W}}_*$ , and let  $\tilde{\mathbf{C}}_* \in \mathbb{R}^{r \times k}$  be the matrix formed by stacking the  $k$  vectors. Algorithm 1 iterates over a set  $\mathcal{N}_{\epsilon/2}^{\otimes k}(\mathbb{S}^{r-1})$  of points ( $r \times k$  matrices)  $\mathbf{C}$ . Recall that  $\mathcal{N}_{\epsilon/2}^{\otimes k}(\mathbb{S}^{r-1})$  is the  $k$ th cartesian power of an  $\epsilon/2$ -net of the  $r$ -dimensional  $\ell_2$ -unit sphere. By construction, the set contains a matrix  $\mathbf{C}_\#$  such that

$$\|\mathbf{C}_\# - \tilde{\mathbf{C}}_{*j}\|_{\infty, 2} \leq \epsilon/2. \quad (12)$$

Then, for all  $j \in \{1, \dots, k\}$ ,

$$\begin{aligned} \|\overline{\mathbf{M}}^\top \overline{\mathbf{W}}_{*j}\|_2 &= |\langle \overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{*j}, \tilde{\mathbf{c}}_{*j} \rangle| \\ &= |\langle \overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{*j}, \mathbf{c}_{\#j} \rangle + \langle \overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{*j}, (\tilde{\mathbf{c}}_{*j} - \mathbf{c}_{\#j}) \rangle| \\ &\leq |\langle \overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{*j}, \mathbf{c}_{\#j} \rangle| + |\langle \overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{*j}, (\tilde{\mathbf{c}}_{*j} - \mathbf{c}_{\#j}) \rangle| \\ &\leq |\langle \overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{*j}, \mathbf{c}_{\#j} \rangle| + \|\overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{*j}\|_2 \cdot \|\tilde{\mathbf{c}}_{*j} - \mathbf{c}_{\#j}\|_2 \\ &\leq |\langle \overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{*j}, \mathbf{c}_{\#j} \rangle| + (\epsilon/2) \cdot \|\overline{\mathbf{M}}^\top \overline{\mathbf{W}}_{*j}\|_2. \end{aligned} \quad (13)$$

The first step follows by the definition of  $\tilde{\mathbf{C}}_*$ , the second by the linearity of the inner product, the third by the triangle inequality, the fourth by Cauchy-Schwarz inequality and the last by the fact that  $\|\tilde{\mathbf{c}}_{*j} - \mathbf{c}_{\#j}\| \leq \epsilon/2, \forall i \in \{1, \dots, k\}$  (by (12)). Rearranging the terms in (13),

$$|\langle \overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{*j}, \mathbf{c}_{\#j} \rangle| \geq (1 - \frac{\epsilon}{2}) \cdot \|\overline{\mathbf{M}}^\top \overline{\mathbf{W}}_{*j}\|_2 \geq 0,$$

which in turn implies (by taking the square on both sides) that

$$\langle \overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{*j}, \mathbf{c}_{\#j} \rangle^2 \geq \|\overline{\mathbf{M}}^\top \overline{\mathbf{W}}_{*j}\|_2^2 \geq (1 - \epsilon) \cdot \|\overline{\mathbf{M}}^\top \overline{\mathbf{W}}_{*j}\|_2^2 \quad (14)$$

Summing the terms in (14) over all  $j \in \{1, \dots, k\}$ ,

$$\sum_{j=1}^k \langle \overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{*j}, \mathbf{c}_{\#j} \rangle^2 \geq (1 - \epsilon) \cdot \|\overline{\mathbf{M}}^\top \overline{\mathbf{W}}_*\|_{\mathbb{F}}^2. \quad (15)$$

Let  $\mathbf{W}_\# \in \mathcal{W}_k$  be the candidate solution produced by the algorithm at  $\mathbf{C}_\#, i.e.,$

$$\mathbf{W}_\# \triangleq \arg \max_{\mathbf{W} \in \mathcal{W}_k} \sum_{j=1}^k \langle \mathbf{w}_j, \overline{\mathbf{U}} \overline{\Sigma} \mathbf{c}_{\#j} \rangle^2 \quad (16)$$

Then,

$$\begin{aligned} \|\overline{\mathbf{M}}^\top \mathbf{W}_\#\| &\stackrel{(\alpha)}{=} \sum_{j=1}^k \max_{\mathbf{c}_j \in \mathbb{S}^{r-1}} \langle \overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{\#j}, \mathbf{c}_j \rangle^2 \\ &\stackrel{(\beta)}{\geq} \sum_{j=1}^k \langle \overline{\Sigma} \overline{\mathbf{U}}^\top \overline{\mathbf{W}}_{\#j}, \mathbf{c}_{\#j} \rangle^2 \\ &\stackrel{(\gamma)}{\geq} \sum_{j=1}^k \langle \overline{\mathbf{W}}_{*j}, \overline{\mathbf{U}} \overline{\Sigma} \mathbf{c}_{\#j} \rangle^2 \\ &\stackrel{(\delta)}{\geq} (1 - \epsilon) \cdot \|\overline{\mathbf{M}} \overline{\mathbf{W}}_*\|_{\mathbb{F}}^2, \end{aligned} \quad (17)$$

where  $(\alpha)$  follows from the observation in (11),  $(\beta)$  from the suboptimality of  $\mathbf{C}_\#$ ,  $(\gamma)$  from the fact that  $\mathbf{W}_\#$  maximizes the sum by its definition in (16), while  $(\delta)$  follows from (15). According to (17), at least one of the candidate solutions produced by Algorithm 1, namely  $\mathbf{W}_\#$ , achieves an objective value within a multiplicative factor  $(1 - \epsilon)$  from the optimal, implying the guarantees of the lemma.

Finally, the running time of Algorithm 1 follows immediately from the cost per iteration and the cardinality of the  $\epsilon/2$ -net on the unit-sphere. Note that matrix multiplications can exploit the available singular value decomposition which is performed once.  $\square$

### A.3 Proof of Theorem 1

We first prove some auxiliary lemmata. The proof of the Theorem is given in the end of this section.

**Lemma 3.** *For any real  $m \times n$  matrices  $\mathbf{M}$  and  $\bar{\mathbf{M}}$ , let*

$$\mathbf{W}_\star \triangleq \arg \max_{\mathbf{W} \in \mathcal{W}_k} \|\mathbf{M}^\top \mathbf{W}\|_F^2 \quad \text{and} \quad \bar{\mathbf{W}}_\star \triangleq \arg \max_{\mathbf{W} \in \mathcal{W}_k} \|\bar{\mathbf{M}}^\top \mathbf{W}\|_F^2, \quad (18)$$

respectively. Then, for any  $\bar{\mathbf{W}} \in \mathcal{W}_k$  such that  $\|\bar{\mathbf{M}}^\top \bar{\mathbf{W}}\|_F^2 \geq \gamma \cdot \|\bar{\mathbf{M}}^\top \bar{\mathbf{W}}_\star\|_F^2$  for some  $0 < \gamma < 1$ ,

$$\|\mathbf{M}^\top \bar{\mathbf{W}}\|_F^2 \geq \gamma \cdot \|\mathbf{M}^\top \mathbf{W}_\star\|_F^2 - 2 \cdot k \cdot \|\mathbf{M} - \bar{\mathbf{M}}\|_2^2.$$

*Proof.* By the optimality of  $\bar{\mathbf{W}}_\star$  for  $\bar{\mathbf{M}}$ ,

$$\|\bar{\mathbf{M}}^\top \bar{\mathbf{W}}_\star\|_F^2 \geq \|\bar{\mathbf{M}}^\top \mathbf{W}_\star\|_F^2.$$

In turn, for any  $\bar{\mathbf{W}} \in \mathcal{W}_k$  satisfying the assumptions of the lemma,

$$\|\bar{\mathbf{M}}^\top \bar{\mathbf{W}}\|_F^2 \geq \gamma \cdot \|\bar{\mathbf{M}}^\top \mathbf{W}_\star\|_F^2. \quad (19)$$

Let  $\mathbf{A} \triangleq \mathbf{M}\mathbf{M}^\top$ ,  $\tilde{\mathbf{A}} \triangleq \bar{\mathbf{M}}\bar{\mathbf{M}}^\top$ , and  $\mathbf{E} \triangleq \mathbf{A} - \tilde{\mathbf{A}}$ . By the linearity of the trace,

$$\begin{aligned} \|\bar{\mathbf{M}}^\top \bar{\mathbf{W}}\|_F^2 &= \text{Tr}(\bar{\mathbf{W}}^\top \mathbf{A} \bar{\mathbf{W}}) - \text{Tr}(\bar{\mathbf{W}}^\top \mathbf{E} \bar{\mathbf{W}}) \\ &\leq \text{Tr}(\bar{\mathbf{W}}^\top \mathbf{A} \bar{\mathbf{W}}) + |\text{Tr}(\bar{\mathbf{W}}^\top \mathbf{E} \bar{\mathbf{W}})|. \end{aligned} \quad (20)$$

By Lemma 10,

$$|\text{Tr}(\bar{\mathbf{W}}^\top \mathbf{E} \bar{\mathbf{W}})| \leq \|\bar{\mathbf{W}}\|_F^2 \cdot \|\mathbf{E}\|_2 \leq k \cdot \|\mathbf{E}\|_2 \triangleq R, \quad (21)$$

where the last inequality follows from the fact that  $\|\mathbf{W}\|_F^2 \leq k$  for any  $\mathbf{W} \in \mathcal{W}_k$ . Continuing from (20),

$$\|\bar{\mathbf{M}}^\top \bar{\mathbf{W}}\|_F^2 \leq \text{Tr}(\bar{\mathbf{W}}^\top \mathbf{A} \bar{\mathbf{W}}) + R. \quad (22)$$

Similarly,

$$\begin{aligned} \|\bar{\mathbf{M}}^\top \mathbf{W}_\star\|_F^2 &= \text{Tr}(\mathbf{W}_\star^\top \mathbf{A} \mathbf{W}_\star) - \text{Tr}(\mathbf{W}_\star^\top \mathbf{E} \mathbf{W}_\star) \\ &\geq \text{Tr}(\mathbf{W}_\star^\top \mathbf{A} \mathbf{W}_\star) - |\text{Tr}(\mathbf{W}_\star^\top \mathbf{E} \mathbf{W}_\star)| \\ &\geq \text{Tr}(\mathbf{W}_\star^\top \mathbf{A} \mathbf{W}_\star) - R. \end{aligned} \quad (23)$$

Combining the above, we have

$$\begin{aligned} \text{Tr}(\bar{\mathbf{W}}^\top \mathbf{A} \bar{\mathbf{W}}) &\geq \|\bar{\mathbf{M}}^\top \bar{\mathbf{W}}\|_F^2 - R \\ &\geq \gamma \cdot \|\bar{\mathbf{M}}^\top \mathbf{W}_\star\|_F^2 - R \\ &\geq \gamma \cdot (\text{Tr}(\mathbf{W}_\star^\top \mathbf{A} \mathbf{W}_\star) - R) - R \\ &= \gamma \cdot \text{Tr}(\mathbf{W}_\star^\top \mathbf{A} \mathbf{W}_\star) - (1 + \gamma) \cdot R \\ &\geq \gamma \cdot \text{Tr}(\mathbf{W}_\star^\top \mathbf{A} \mathbf{W}_\star) - 2 \cdot R, \end{aligned}$$

where the first inequality follows from (22) the second from (19), the third from (23), and the last from the fact that  $R \geq 0$  and  $0 < \gamma \leq 1$ . This concludes the proof.  $\square$

**Remark 1.** If in Lemma 3  $\bar{\mathbf{M}}$  is such that  $\mathbf{M}\mathbf{M}^\top - \bar{\mathbf{M}}\bar{\mathbf{M}}^\top$  is PSD, then

$$\|\mathbf{M}^\top \bar{\mathbf{W}}\|_{\text{F}}^2 \geq \gamma \cdot \|\mathbf{M}^\top \mathbf{W}_*\|_{\text{F}}^2 - \sum_{i=1}^k \lambda_i(\mathbf{M}\mathbf{M}^\top - \bar{\mathbf{M}}\bar{\mathbf{M}}^\top).$$

*Proof.* This follows from the fact that if  $\mathbf{E} \triangleq \mathbf{A} - \tilde{\mathbf{A}}$  is PSD, then

$$\text{Tr}(\tilde{\mathbf{X}}^\top \mathbf{E} \tilde{\mathbf{X}}) = \sum_j \mathbf{x}_j^\top \mathbf{E} \mathbf{x}_j \geq 0,$$

and the bound in (20) can be improved to

$$\begin{aligned} \|\bar{\mathbf{M}}^\top \bar{\mathbf{W}}\|_{\text{F}}^2 &= \text{Tr}(\tilde{\mathbf{X}}^\top \mathbf{A} \tilde{\mathbf{X}}) - \text{Tr}(\tilde{\mathbf{X}}^\top \mathbf{E} \tilde{\mathbf{X}}) \\ &\leq \text{Tr}(\tilde{\mathbf{X}}^\top \mathbf{A} \tilde{\mathbf{X}}). \end{aligned}$$

Further, by Lemma 10 (Corollary 3) the bound in (21) becomes

$$\text{Tr}(\bar{\mathbf{W}}^\top \mathbf{E} \bar{\mathbf{W}}) \leq \|\bar{\mathbf{W}}\|_{\text{F}}^2 \cdot \|\mathbf{E}\|_2 \leq \sum_{i=1}^k \lambda_i(\mathbf{E}).$$

The rest of the proof follows.  $\square$

**Theorem 1.** For any real  $m \times n$  (not necessarily nonnegative) matrix  $\mathbf{M}$  and desired number of components  $k$ , let  $\mathbf{W}_* \triangleq \arg \max_{\mathbf{W} \in \mathcal{W}_k} \|\mathbf{M}^\top \mathbf{W}\|_{\text{F}}^2$ . Let  $\bar{\mathbf{M}}$  be the best rank- $r$  approximation of  $\mathbf{M}$ . Algorithm 1 with input  $\bar{\mathbf{M}}$  and accuracy parameters  $\epsilon$  and  $r$ , outputs  $\bar{\mathbf{W}} \in \mathcal{W}_k$  such that

$$\|\mathbf{M}^\top \bar{\mathbf{W}}\|_{\text{F}}^2 \geq (1 - \epsilon) \cdot \|\mathbf{M}^\top \mathbf{W}_*\|_{\text{F}}^2 - k \cdot \|\mathbf{M} - \bar{\mathbf{M}}\|_2^2$$

in time  $T_{\text{SVD}} + O\left(\left(\frac{1}{\epsilon}\right)^{r-k} \cdot k \cdot m\right)$ .

*Proof.* Let  $\bar{\mathbf{W}}$  be the output of Algorithm 1 with input the best rank- $r$  approximation of  $\mathbf{M}$ ,  $\bar{\mathbf{M}}$ . By the guarantees of Algorithm 1, (Lemma 1), the output  $\bar{\mathbf{W}} \in \mathcal{W}_k$  of Algorithm 1 is such that

$$\|\bar{\mathbf{M}}^\top \bar{\mathbf{W}}\|_{\text{F}}^2 \geq (1 - \epsilon) \cdot \|\bar{\mathbf{M}}^\top \mathbf{W}_*\|_{\text{F}}^2,$$

where  $\bar{\mathbf{W}}_* \triangleq \arg \max_{\mathbf{W} \in \mathcal{W}_k} \|\bar{\mathbf{M}}^\top \mathbf{W}\|_{\text{F}}^2$ . In turn, by Lemma 3 (and in particular taking into account the remark 1 whose conditions are satisfied since  $\mathbf{M}\mathbf{M} - \bar{\mathbf{M}}\bar{\mathbf{M}}^\top$  is PSD), we have

$$\begin{aligned} \|\mathbf{M}^\top \bar{\mathbf{W}}\|_{\text{F}}^2 &\geq (1 - \epsilon) \cdot \|\mathbf{M}^\top \mathbf{W}_*\|_{\text{F}}^2 - \sum_{i=1}^k \lambda_i(\mathbf{M}\mathbf{M}^\top - \bar{\mathbf{M}}\bar{\mathbf{M}}^\top) \\ &= (1 - \epsilon) \cdot \|\mathbf{M}^\top \mathbf{W}_*\|_{\text{F}}^2 - \sum_{i=r+1}^{r+k} \sigma_i^2(\mathbf{M}). \end{aligned} \quad (24)$$

The desired result readily follows.  $\square$

#### A.4 Proof of Theorem 2

**Lemma 4.** For any  $m \times n$  real nonnegative matrix  $\mathbf{M}$ , target dimension  $k$ , and accuracy parameters  $r \in [n]$  and  $\epsilon > 0$ , Algorithm 2 outputs an ONMF pair  $\bar{\mathbf{W}}, \bar{\mathbf{H}}$ , such that

$$\|\mathbf{M} - \bar{\mathbf{W}}\bar{\mathbf{H}}^\top\|_{\text{F}}^2 \leq \mathcal{E}_* + \epsilon \cdot \sum_{i=1}^k \sigma_i^2(\mathbf{M}) + \sum_{i=r+1}^{r+k} \sigma_i^2(\mathbf{M}),$$

in time  $T_{\text{SVD}} + O\left(\left(\frac{2}{\epsilon}\right)^{r-k} \cdot k \cdot m\right)$ .

*Proof.* Recall that given a real nonnegative  $m \times k$  matrix  $\mathbf{W} \in \mathcal{W}_k$ ,  $\mathbf{H}^\top = \mathbf{W}_*^\top \mathbf{M}$  minimizes the Frobenius error  $\|\mathbf{M} - \mathbf{W}\mathbf{H}^\top\|_F^2$  over the set of real nonnegative  $n \times k$  matrices (Proof of Lemma 2). In turn, for any  $\mathbf{W} \in \mathcal{W}_k$ , and  $\mathbf{H}$  selected as above,

$$\|\mathbf{M} - \mathbf{W}\mathbf{H}^\top\|_F^2 = \|\mathbf{M}\|_F^2 - \|\mathbf{M}^\top \mathbf{W}\|_F^2. \quad (25)$$

Let  $\overline{\mathbf{W}}$  be the output of Algorithm 1, for input matrix  $\mathbf{M}_r$  the best rank- $r$  approximation of  $\mathbf{M}$ . That is, in the sequel of this proof,  $\overline{\mathbf{M}} = \mathbf{M}_r$ . By the guarantees of Algorithm 1, (Lemma 1), the output  $\overline{\mathbf{W}} \in \mathcal{W}_k$  of Algorithm 1 is such that

$$\|\mathbf{M}_r^\top \overline{\mathbf{W}}\|_F^2 \geq (1 - \epsilon) \cdot \|\mathbf{M}_r^\top \overline{\mathbf{W}}_*\|_F^2,$$

where  $\overline{\mathbf{W}}_* \triangleq \arg \max_{\mathbf{W} \in \mathcal{W}_k} \|\mathbf{M}_r^\top \mathbf{W}\|_F^2$ . In turn, by Lemma 3 (and in particular taking into account the remark 1 whose conditions are satisfied since  $\mathbf{M}\mathbf{M}^\top - \mathbf{M}_r\mathbf{M}_r^\top$  is PSD), we have

$$\begin{aligned} \|\mathbf{M}^\top \overline{\mathbf{W}}\|_F^2 &\geq (1 - \epsilon) \cdot \|\mathbf{M}^\top \overline{\mathbf{W}}_*\|_F^2 - \sum_{i=1}^k \lambda_i(\mathbf{M}\mathbf{M}^\top - \mathbf{M}_r\mathbf{M}_r^\top) \\ &= (1 - \epsilon) \cdot \|\mathbf{M}^\top \overline{\mathbf{W}}_*\|_F^2 - \sum_{i=r+1}^{r+k} \sigma_i^2(\mathbf{M}). \end{aligned} \quad (26)$$

Given the output  $\overline{\mathbf{W}}$  of Algorithm 1, Algorithm 2 outputs the pair  $\overline{\mathbf{W}}, \overline{\mathbf{H}}^\top \triangleq \overline{\mathbf{W}}^\top \mathbf{M}$ . By (25), for this choice of  $\overline{\mathbf{H}}$ , taking into account (26),

$$\begin{aligned} \|\mathbf{M} - \overline{\mathbf{W}}\overline{\mathbf{H}}^\top\|_F^2 &= \|\mathbf{M}\|_F^2 - \|\mathbf{M}^\top \overline{\mathbf{W}}\|_F^2 \\ &\leq \|\mathbf{M}\|_F^2 - (1 - \epsilon) \cdot \|\mathbf{M}^\top \overline{\mathbf{W}}_*\|_F^2 + \sum_{i=r+1}^{r+k} \sigma_i^2(\mathbf{M}) \\ &= \|\mathbf{M}\|_F^2 - \|\mathbf{M}^\top \overline{\mathbf{W}}_*\|_F^2 + \epsilon \cdot \|\mathbf{M}^\top \overline{\mathbf{W}}_*\|_F^2 + \sum_{i=r+1}^{r+k} \sigma_i^2(\mathbf{M}) \\ &= \|\mathbf{M} - \overline{\mathbf{W}}_*\overline{\mathbf{H}}_*^\top\|_F^2 + \epsilon \cdot \|\mathbf{M}^\top \overline{\mathbf{W}}_*\|_F^2 + \sum_{i=r+1}^{r+k} \sigma_i^2(\mathbf{M}) \\ &= \|\mathbf{M} - \overline{\mathbf{W}}_*\overline{\mathbf{H}}_*^\top\|_F^2 + \epsilon \cdot \sum_{i=1}^k \sigma_i^2(\mathbf{M}) + \sum_{i=r+1}^{r+k} \sigma_i^2(\mathbf{M}), \end{aligned}$$

where the last inequality follows by Lemma 10. This completes the proof.  $\square$

**Theorem 2.** For any  $m \times n$  real nonnegative matrix  $\mathbf{M}$ , target dimension  $k$ , and desired accuracy  $0 < \epsilon < 1$ , Algorithm 2 with parameters  $\epsilon$  and  $r = \lceil k/\epsilon \rceil$  outputs an ONMF pair  $\overline{\mathbf{W}}, \overline{\mathbf{H}}$ , such that

$$\|\mathbf{M} - \overline{\mathbf{W}}\overline{\mathbf{H}}^\top\|_F^2 \leq \mathcal{E}_* + \epsilon \cdot \|\mathbf{M}\|_F^2,$$

in time  $T_{\text{SVD}} + \left(\frac{1}{\epsilon}\right)^{(k^2/\epsilon)} \cdot (k \cdot m)$ .

*Proof.* By Lemma 4, Algorithm 2 with parameters  $r$  and  $\epsilon$ , outputs an ONMF pair  $\overline{\mathbf{W}}, \overline{\mathbf{H}}$ , such that

$$\|\mathbf{M} - \overline{\mathbf{W}}\overline{\mathbf{H}}^\top\|_F^2 \leq \mathcal{E}_* + \epsilon \cdot \sum_{i=1}^k \sigma_i^2(\mathbf{M}) + \sum_{i=r+1}^{r+k} \sigma_i^2(\mathbf{M}). \quad (27)$$

Noting that for  $k < r$  the  $k$  (squared) singular values  $\sigma_i^2(\mathbf{M})$ ,  $i = r + 1, \dots, r + k$  are the smallest among the  $r$  (squared) singular values  $\sigma_i^2(\mathbf{M})$ ,  $i = k + 1, \dots, r + k$ , the last term in the right-hand side can be upper bounded as follows:

$$\sum_{i=r+1}^{r+k} \sigma_i^2(\mathbf{M}) \leq \frac{k}{r} \cdot \sum_{i=k+1}^{r+k} \sigma_i^2(\mathbf{M}). \quad (28)$$

For  $r = \lceil k/\epsilon \rceil$ , and combining the (28) and (27), we have

$$\begin{aligned} \|\mathbf{M} - \overline{\mathbf{W}}\overline{\mathbf{H}}^\top\|_{\mathbb{F}}^2 &\leq \mathcal{E}_* + \epsilon \cdot \sum_{i=1}^k \sigma_i^2(\mathbf{M}) + \epsilon \cdot \sum_{i=k+1}^{r+k} \sigma_i^2(\mathbf{M}) \\ &= \mathcal{E}_* + \epsilon \cdot \sum_{i=1}^{r+k} \sigma_i^2(\mathbf{M}) \\ &\leq \mathcal{E}_* + \epsilon \cdot \|\mathbf{M}\|_{\mathbb{F}}^2, \end{aligned}$$

which is the desired guarantee. The time complexity readily follows from the steps of Algorithm 2 and that of Algorithm 1, which concludes the proof.  $\square$

### A.5 Correctness of Algorithm 3

**Lemma 5.** *For any  $m \times k$  matrix  $\mathbf{A}$ , Algorithm 3 outputs the  $m \times k$  nonnegative matrix*

$$\widehat{\mathbf{W}} = \arg \max_{\mathbf{W} \in \mathcal{W}_k} \sum_{j=1}^k \langle \mathbf{w}_j, \mathbf{a}_j \rangle^2,$$

in time  $O(2^k \cdot k \cdot m)$ .

*Proof.* Let  $\mathcal{I}_j \subseteq [m]$ ,  $j = 1 \dots, k$  denote the supports of the  $k$  columns of optimal solution  $\widehat{\mathbf{W}}$ . The orthogonality requirements (in conjunction with nonnegativity) imply that the supports  $\mathcal{I}_j$ ,  $j = 1, \dots, k$  are disjoint. Further, it is straightforward to verify that due to the nonnegativity constrains in  $\mathcal{W}_k$ , the support of the  $j$ th column,  $\mathcal{I}_j$ , must contain only indices corresponding to nonnegative or nonpositive entries of  $\mathbf{a}_j$ , but not a combination of both. Algorithm 3 considers all  $2^k$  sign combinations for the support sets, e.g.,  $\mathcal{I}_1$  containing positive entries,  $\mathcal{I}_2$  negative, etc., by equivalently solving the maximization on all  $2^k$  matrices  $\hat{\mathbf{A}} = \mathbf{A} \cdot \text{diag}(\mathbf{s})$ ,  $\mathbf{s} \in \{\pm 1\}^k$  and returning the solution that performs best on the original input  $\mathbf{A}$ . Therefore, without loss of generality, in the sequel we assume that all support sets correspond to nonnegative entries of  $\mathbf{A}$ .

If an oracle reveals the supports  $\mathcal{I}_j$ ,  $j = 1, \dots, k$ , the exact value of  $\widehat{\mathbf{X}}$  can be readily determined, according to the Cauchy-Schwarz inequality: the  $j$ th column,  $\widehat{\mathbf{x}}_j$ , is supported only on  $\mathcal{I}_j$ , and its nonzero sub-vector is set to  $(\widehat{\mathbf{x}}_j)_{\mathcal{I}_j} = [\mathbf{a}_j]_{\mathcal{I}_j} / \|\mathbf{a}_j\|$ , which maximizes the inner product with the corresponding sub-vector of  $\mathbf{a}_j$ . In turn, the objective function attains value equal to

$$\sum_{j=1}^k (\widehat{\mathbf{x}}_j^\top \mathbf{a}_j)^2 = \sum_{j=1}^k \sum_{i \in \mathcal{I}_j} A_{ij}^2, \quad (29)$$

where the first equality stems from the fact that  $[\mathbf{a}_j]_{\mathcal{I}_j} \geq \mathbf{0}$ . It suffices to show that Alg. 3 correctly determines the collection of support sets  $\mathcal{I}_j$ ,  $j = 1, \dots, k$ .

Alg. 3 constructs the collection of support sets  $\mathcal{I}_j$ ,  $j = 1, \dots, k$ , according to the following rule:

$$i \in \mathcal{I}_j \Leftrightarrow A_{ij} > \max\{0, A_{iw}\}, \quad \forall w \in [k] \setminus \{j\}, \quad (30)$$

i.e., index  $i \in [m]$  is assigned to the support of the  $j$ th column if and only if  $A_{ij}$  is positive and the largest entry in the  $i$ th row of  $\mathbf{A}$ . Note that any procedure to construct supports that satisfy the requirements described in the beginning of this proof would assign each index  $i \in [m]$  to at most one of the sets  $\mathcal{I}_j$ ,  $j \in \{1, \dots, k\}$ , while it would need to ensure that  $i \in \mathcal{I}_j$  if and only if  $A_{ij} > 0$ . The rule in (30) additionally requires that index  $i \in [m]$  is assigned to  $\mathcal{I}_j$  if and only if  $A_{ij}$  is the largest (positive) entry in the  $i$ th row of  $\mathbf{A}$ .

Assume, for the sake of contradiction, that there exists a set of optimal supports  $\mathcal{I}_j$ ,  $j = 1, \dots, k$  which does not adhere to the rule in (30), i.e., there exist  $u \in [k]$  and  $q \in [m]$ , such  $q \in \mathcal{I}_u$ , while  $0 < A_{qu} < A_{qv}$  for some  $v \in [k]$ ,  $v \neq u$ . Consider a collection of supports  $\tilde{\mathcal{I}}_j$ ,  $j = 1, \dots, k$ , with

$$\tilde{\mathcal{I}}_j = \mathcal{I}_j, \quad \forall j \in [k] \setminus \{u, v\}, \quad \tilde{\mathcal{I}}_u = \mathcal{I}_u \setminus \{q\} \quad \text{and} \quad \tilde{\mathcal{I}}_v = \mathcal{I}_v \cup \{q\}. \quad (31)$$

Note that the collection of supports in (31) satisfies the desired constraints. Further, the objective value achieved for the new supports (according to (29)) is equal to

$$\sum_{j=1}^k \sum_{i \in \tilde{\mathcal{I}}_j} A_{ij}^2 = \sum_{j=1}^k \sum_{i \in \mathcal{I}_j} A_{ij}^2 - A_{qu}^2 + A_{qv}^2 > \sum_{j=1}^k \sum_{i \in \mathcal{I}_j} A_{ij}^2,$$

contradicting the optimality of the collection  $\mathcal{I}_j$ ,  $j = 1, \dots, k$ . We conclude that the collection of optimal support sets must satisfy (30).

The construction of the support sets according to (30) requires determining the largest entry of each of the  $m$  rows of  $\mathbf{A}$ , which can be done in  $O(km)$ . Once the supports are determined, each of the  $k$  columns of  $\hat{\mathbf{X}}$  is constructed in  $O(m)$ . Taking into account that the above procedure is repeated  $2^k$  times for each of the sign patters, the desired result follows.  $\square$

## B Auxiliary Technical Results

**Lemma 6.** For any real  $m \times n$  matrix  $\mathbf{M}$ , and any  $r, k \leq \min\{m, n\}$ ,

$$\sum_{i=r+1}^{r+k} \sigma_i(\mathbf{M}) \leq \frac{k}{\sqrt{r+k}} \cdot \|\mathbf{M}\|_{\text{F}},$$

where  $\sigma_i(\mathbf{M})$  is the  $i$ th largest singular value of  $\mathbf{M}$ .

*Proof.* By the Cauchy-Schwartz inequality,

$$\sum_{i=r+1}^{r+k} \sigma_i(\mathbf{M}) = \sum_{i=r+1}^{r+k} |\sigma_i(\mathbf{M})| \leq \left( \sum_{i=r+1}^{r+k} \sigma_i^2(\mathbf{M}) \right)^{1/2} \cdot \|\mathbf{1}_k\|_2 = \sqrt{k} \cdot \left( \sum_{i=r+1}^{r+k} \sigma_i^2(\mathbf{M}) \right)^{1/2}.$$

Note that  $\sigma_{r+1}(\mathbf{M}), \dots, \sigma_{r+k}(\mathbf{M})$  are the  $k$  smallest among the  $r+k$  largest singular values. Hence,

$$\sum_{i=r+1}^{r+k} \sigma_i^2(\mathbf{M}) \leq \frac{k}{r+k} \sum_{i=1}^{r+k} \sigma_i^2(\mathbf{M}) \leq \frac{k}{r+k} \sum_{i=1}^{\min\{m,n\}} \sigma_i^2(\mathbf{M}) = \frac{k}{r+k} \|\mathbf{M}\|_{\text{F}}^2.$$

Combining the two inequalities, the desired result follows.  $\square$

**Corollary 1.** For any real  $m \times n$  matrix  $\mathbf{M}$  and  $k \leq \min\{m, n\}$ ,  $\sigma_k(\mathbf{M}) \leq k^{-1/2} \cdot \|\mathbf{M}\|_{\text{F}}$ .

*Proof.* It follows immediately from Lemma 6.  $\square$

**Lemma 7.** Let  $a_1, \dots, a_n$  and  $b_1, \dots, b_n$  be  $2n$  real numbers and let  $p$  and  $q$  be two numbers such that  $1/p + 1/q = 1$  and  $p > 1$ . We have

$$\left| \sum_{i=1}^n a_i b_i \right| \leq \left( \sum_{i=1}^n |a_i|^p \right)^{1/p} \cdot \left( \sum_{i=1}^n |b_i|^q \right)^{1/q}.$$

**Lemma 8.** For any  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times k}$ ,

$$|\langle \mathbf{A}, \mathbf{B} \rangle| \triangleq |\text{Tr}(\mathbf{A}^\top \mathbf{B})| \leq \|\mathbf{A}\|_{\text{F}} \|\mathbf{B}\|_{\text{F}}.$$

*Proof.* The inequality follows from Lemma 7 for  $p = q = 2$ , treating  $\mathbf{A}$  and  $\mathbf{B}$  as vectors.  $\square$

**Lemma 9.** For any two real matrices  $\mathbf{A}$  and  $\mathbf{B}$  of appropriate dimensions,

$$\|\mathbf{A}\mathbf{B}\|_{\text{F}} \leq \min\{\|\mathbf{A}\|_2 \|\mathbf{B}\|_{\text{F}}, \|\mathbf{A}\|_{\text{F}} \|\mathbf{B}\|_2\}.$$

*Proof.* Let  $\mathbf{b}_i$  denote the  $i$ th column of  $\mathbf{B}$ . Then,

$$\|\mathbf{AB}\|_{\text{F}}^2 = \sum_i \|\mathbf{Ab}_i\|_2^2 \leq \sum_i \|\mathbf{A}\|_2^2 \|\mathbf{b}_i\|_2^2 = \|\mathbf{A}\|_2^2 \sum_i \|\mathbf{b}_i\|_2^2 = \|\mathbf{A}\|_2^2 \|\mathbf{B}\|_{\text{F}}^2.$$

Similarly, using the previous inequality,

$$\|\mathbf{AB}\|_{\text{F}}^2 = \|\mathbf{B}^{\top} \mathbf{A}^{\top}\|_{\text{F}}^2 \leq \|\mathbf{B}^{\top}\|_2^2 \|\mathbf{A}^{\top}\|_{\text{F}}^2 = \|\mathbf{B}\|_2^2 \|\mathbf{A}\|_{\text{F}}^2.$$

Combining the two upper bounds, the desired result follows.  $\square$

**Corollary 2.** Let  $\mathbf{M}_r$  denote the best rank- $r$  approximation of  $\mathbf{M}$ , obtained by the truncated singular value decomposition of  $\mathbf{M}$ . Then, for any  $d > r$ ,  $\sigma_d \leq \|\mathbf{M} - \mathbf{M}_r\|_{\text{F}} / \sqrt{d - r}$ .

*Proof.* By definition

$$\|\mathbf{M} - \mathbf{M}_r\|_{\text{F}}^2 = \sum_{i=r+1}^n \sigma_i^2 \geq \sum_{i=r+1}^d \sigma_i^2 \geq (d - r) \cdot \sigma_d^2,$$

from which the desired result follows.  $\square$

**Lemma 10.** For any real  $m \times n$  matrix  $\mathbf{A}$ , and any  $k \leq \min\{m, n\}$ ,

$$\max_{\substack{\mathbf{Y} \in \mathbb{R}^{n \times k} \\ \mathbf{Y}^{\top} \mathbf{Y} = \mathbf{I}_k}} \|\mathbf{AY}\|_{\text{F}} = \left( \sum_{i=1}^k \sigma_i^2(\mathbf{A}) \right)^{1/2}.$$

Equality is realized for  $\mathbf{Y}$  coinciding with the  $k$  leading right singular vectors of  $\mathbf{A}$ .

*Proof.* Let  $\mathbf{U}\Sigma\mathbf{V}^{\top}$  be the singular value decomposition of  $\mathbf{A}$ ;  $\mathbf{U}$  and  $\mathbf{V}$  are  $m \times m$  and  $n \times n$  unitary matrices respectively, while  $\Sigma$  is a diagonal matrix with  $\Sigma_{jj} = \sigma_j$ , the  $j$ th largest singular value of  $\mathbf{A}$ ,  $j = 1, \dots, d$ , where  $d \triangleq \min\{m, n\}$ . Due to the invariance of the Frobenius norm under unitary multiplication,

$$\|\mathbf{AY}\|_{\text{F}}^2 = \|\mathbf{U}\Sigma\mathbf{V}^{\top} \mathbf{Y}\|_{\text{F}}^2 = \|\Sigma\mathbf{V}^{\top} \mathbf{Y}\|_{\text{F}}^2. \quad (32)$$

Continuing from (32),

$$\|\Sigma\mathbf{V}^{\top} \mathbf{Y}\|_{\text{F}}^2 = \text{TR}(\mathbf{Y}^{\top} \mathbf{V} \Sigma^2 \mathbf{V}^{\top} \mathbf{Y}) = \sum_{i=1}^k \mathbf{y}_i^{\top} \left( \sum_{j=1}^d \sigma_j^2 \cdot \mathbf{v}_j \mathbf{v}_j^{\top} \right) \mathbf{y}_i = \sum_{j=1}^d \sigma_j^2 \cdot \sum_{i=1}^k (\mathbf{v}_j^{\top} \mathbf{y}_i)^2.$$

Let  $z_j \triangleq \sum_{i=1}^k (\mathbf{v}_j^{\top} \mathbf{y}_i)^2$ ,  $j = 1, \dots, d$ . Note that each individual  $z_j$  satisfies

$$0 \leq z_j \triangleq \sum_{i=1}^k (\mathbf{v}_j^{\top} \mathbf{y}_i)^2 \leq \|\mathbf{v}_j\|^2 = 1,$$

where the last inequality follows from the fact that the columns of  $\mathbf{Y}$  are orthonormal. Further,

$$\sum_{j=1}^d z_j = \sum_{j=1}^d \sum_{i=1}^k (\mathbf{v}_j^{\top} \mathbf{y}_i)^2 = \sum_{i=1}^k \sum_{j=1}^d (\mathbf{v}_j^{\top} \mathbf{y}_i)^2 = \sum_{i=1}^k \|\mathbf{y}_i\|^2 = k.$$

Combining the above, we conclude that

$$\|\mathbf{AY}\|_{\text{F}}^2 = \sum_{j=1}^d \sigma_j^2 \cdot z_j \leq \sigma_1^2 + \dots + \sigma_k^2. \quad (33)$$

Finally, it is straightforward to verify that if  $\mathbf{y}_i = \mathbf{v}_i$ ,  $i = 1, \dots, k$ , then (33) holds with equality.  $\square$

**Corollary 3.** For any real  $m \times m$  PSD matrix  $\mathbf{A}$ , and  $k \times m$  matrix  $\mathbf{X}$  with  $k \leq m$  orthonormal columns,

$$\mathrm{Tr}(\mathbf{X}^\top \mathbf{A} \mathbf{X}) = \sum_{i=1}^k \lambda_i(\mathbf{A})$$

where  $\lambda_i(\mathbf{A})$  is the  $i$ th largest eigenvalue of  $\mathbf{A}$ . Equality is achieved for  $\mathbf{X}$  coinciding with the  $k$  leading eigenvectors of  $\mathbf{A}$ .

*Proof.* Let  $\mathbf{A} = \mathbf{V}\mathbf{V}^\top$  be a factorization of the PSD matrix  $\mathbf{A}$ . Then,  $\mathrm{Tr}(\mathbf{X}^\top \mathbf{A} \mathbf{X}) = \mathrm{Tr}(\mathbf{X}^\top \mathbf{V}\mathbf{V}^\top \mathbf{X}) = \|\mathbf{V}^\top \mathbf{X}\|_F^2$ . The desired result follows by Lemma 10 and the fact that  $\lambda_i(\mathbf{A}) = \sigma_i^2(\mathbf{V})$ ,  $i = 1, \dots, m$ .  $\square$

## C Net of the $\ell_2$ -unit sphere

In this section, we provide a simple probabilistic construction of an  $\epsilon$ -net of the  $\ell_2$ -unit sphere  $\mathbb{S}_2^{d-1}$ , i.e., the set of points  $\mathbf{x}$  such that  $\|\mathbf{x}\|_2 = 1$ .

**Lemma 11** ([34], Lemma 5.2). For any  $\epsilon > 0$ , there exists an  $\epsilon$ -net  $\mathcal{N}_\epsilon$  of the unit Euclidean sphere  $\mathbb{S}_2^{d-1}$  equipped with the Euclidean metric, such that

$$m_\epsilon \triangleq |\mathcal{N}_\epsilon| \leq (1 + 2/\epsilon)^d.$$

*Proof.* Let  $\mathcal{N}_\epsilon$  be a maximal  $\epsilon$ -separated subset of  $\mathbb{S}_2^{d-1}$ . In other words,  $d(x, y) \geq \epsilon$  for all  $x, y \in \mathcal{N}_\epsilon$ ,  $x \neq y$ , and no set containing  $\mathcal{N}_\epsilon$  has this property.

Such a set can be constructed iteratively: select an arbitrary point on the sphere and at each subsequent step select a point that is at distance at least  $\epsilon$  from all previously selected points. By the compactness of the sphere, the iterative construction process will terminate after a finite number of steps, and the resulting set will satisfy the above properties.

The maximality property, implies that  $\mathcal{N}_\epsilon$  is an  $\epsilon$ -net of  $\mathbb{S}_2^{d-1}$ . If this was not the case, then there would exist an  $x \in \mathbb{S}_2^{d-1}$  such that  $d(x, y) > \epsilon$ ,  $\forall y \in \mathcal{N}_\epsilon$ . The  $\mathcal{N}_\epsilon \cup \{x\}$  would be an  $\epsilon$ -separated set, that contains  $\mathcal{N}_\epsilon$ , contradicting the maximality of the latter.

By the separation property, we infer that the balls of radius  $\epsilon/2$  centered at the points of  $\mathcal{N}_\epsilon$  are disjoint. This follows from the triangle inequality. Further, all such balls lie in the ball  $(1 + \epsilon/2) \mathbb{B}_2^d$ , where  $\mathbb{B}_2^d$  denotes the unit Euclidean ball centered at the origin. Comparing the volumes, we have

$$\mathrm{Vol}\left(\frac{\epsilon}{2} \mathbb{B}_2^d\right) \cdot |\mathcal{N}_\epsilon| \leq \mathrm{Vol}\left(\left(1 + \frac{\epsilon}{2}\right) \mathbb{B}_2^d\right).$$

Taking into account that  $\mathrm{Vol}(r \cdot \mathbb{B}_2^d) = r^d \cdot \mathrm{Vol}(\mathbb{B}_2^d)$ ,

$$|\mathcal{N}_\epsilon| \leq \left(1 + \frac{\epsilon}{2}\right)^d / \left(\frac{\epsilon}{2}\right)^d = \left(1 + \frac{2}{\epsilon}\right)^d,$$

which is the desired result.  $\square$

Lemma 11 regards the unit Euclidean sphere. However, the sequence of arguments used in the proof essentially hold for the case of the unit ball  $\mathbb{B}_2^d$ , i.e., there exists an  $\epsilon$ -net of  $\mathbb{B}_2^d$ , with at most  $(1 + 2/\epsilon)^d$  points.

**Constructing an  $\epsilon$ -net of the unit sphere.** There are many constructions for  $\epsilon$ -nets on the sphere, both deterministic and randomized. In the following we review a simple randomized construction, initially studied by Wyner [35] in the asymptotic  $d \rightarrow \infty$  regime.

By Lemma 11, there exists an  $\epsilon$ -net  $\mathcal{N}_\epsilon$  of  $\mathbb{S}_2^{d-1}$ . Consider the balls of radii  $\epsilon$  centered at the points of  $\mathcal{N}_\epsilon$ . The balls cover all points of  $\mathbb{S}_2^{d-1}$ ; if there existed a point  $x$  on  $\mathbb{S}_2^{d-1}$  not included in any ball, it would imply that this point is at distance at least  $\epsilon$  from all points of  $\mathcal{N}_\epsilon$  contradicting the fact that  $\mathcal{N}_\epsilon$  is a  $\epsilon$ -net.

The intersection of each of the previous balls with  $\mathbb{S}_2^{d-1}$  is a spherical cap, and hence, according to the above, the collection of spherical caps covers  $\mathbb{S}_2^{d-1}$ . (Note that the spherical caps, as well as the balls, overlap.)

Consider a set  $\mathcal{Q}$ , containing at least one point from each spherical cap. Then,  $\mathcal{Q}$  is a  $2\epsilon$ -net of  $\mathbb{S}_2^{d-1}$ . To verify that, note the following. Consider a point  $x \in \mathbb{S}_2^{d-1}$ . By construction,  $\mathcal{N}_\epsilon$  contains a point  $y$  such that  $d(x, y) \leq \epsilon$ . Consider the spherical cap centered at  $y$ . By definition,  $\mathcal{Q}$  contains a point  $\tilde{y}$  in that spherical cap, and hence  $d(y, \tilde{y}) \leq \epsilon$ . By triangle inequality, it follows that  $d(x, \tilde{y}) \leq 2\epsilon$ . Since the point  $x$  is arbitrary, we conclude that  $\mathcal{Q}$  is a  $2\epsilon$ -net.

We draw points randomly and independently, uniformly distributed on  $\mathbb{S}_2^{d-1}$ . This can be accomplished, for instance, by randomly and independently generating vectors in  $\mathbb{R}^d$  distributed according to the multivariate normal distribution  $N(\mathbf{0}, \mathbf{I})$  and normalizing their length to 1. A randomly selected point lies in a specific spherical cap with probability  $p \geq 1/m_\epsilon$ . By a standard probability arguments (Coupon collector’s problem),  $O(m_\epsilon \ln(m_\epsilon/\delta))$  points uniformly distributed over  $\mathbb{S}_2^{d-1}$  suffice for at least one random point to lie in each sphere cap with probability at least  $1 - \delta$ . Substituting the value of  $m_\epsilon$  from Lemma 11, we find that  $O(d\epsilon^{-d} \cdot \ln \frac{1}{\epsilon\delta})$  suffice to form a  $2\epsilon$ -net, with probability at least  $1 - \delta$ . Note that  $\delta$  can be chosen to scale with the dimension  $n$  of the problem.

**Lemma 12.** *A set of  $O(d(\epsilon/2)^{-d} \cdot \ln \frac{2}{\epsilon\delta})$  randomly and independently drawn points uniformly distributed on  $\mathbb{S}_2^{d-1}$  suffices to construct an  $\epsilon$ -net of  $\mathbb{S}_2^{d-1}$  with probability at least  $1 - \delta$ .*

## D Additional Experimental Results

Component 1	Component 2	Component 3	Component 4	Component 5
american	coach	add	ago	billion
attack	game	cup	called	business
campaign	games	food	com	companies
country	guy	hour	family	company
government	hit	large	help	cost
group	left	makes	high	deal
leader	night	minutes	home	industry
official	play	oil	look	market
political	player	pepper	need	million
president	point	serving	part	money
zzz_al_gore	run	small	problem	number
zzz_bush	season	sugar	program	percent
zzz_george_bush	team	tablespoon	right	plan
zzz_u_s	win	teaspoon	school	stock
zzz_united_states	won	water	show	system

Table 3: ONMF with  $r=5$  orthogonal components ( $102 \cdot 10^3$ -dimensional vectors) on the words-by-document matrix of the NY Times bag-of-words dataset [33]. The table depicts the words corresponding to the 15 largest entries of each component. The 5 retrieved components are extremely sparse: 90% of their mass is concentrated in 134, 35, 65, 269 and 59 entries, respectively.

**Large-scale text analysis: clustering words.** We evaluate the performance of ONMF S as a clustering algorithm on the NY Times bag-of-words dataset [33]. The dataset is represented by a  $102\text{K} \times 300\text{K}$  words-by-articles matrix. Given an approximate ONMF of that matrix, the  $102\text{K} \times k$  nonnegative, orthogonal factor  $\mathbf{W}$  induces an assignment of words to  $r$  clusters, which in this case can be interpreted as *topics*. That is, each column of  $\mathbf{W}$  suggests a topic, defined by the words corresponding to its nonzero entries.

We run ONMF S with target dimension  $k = 5$  topics, and accuracy parameter  $r = 5$ , while we configure it to stop if no progress is observed after  $T = 300$  consecutive candidate solutions. Table 3 lists the words corresponding to the 15 largest entries of each orthogonal component (column of  $\mathbf{W}$ ). Arguably, each component can be interpreted as a distinct topic, illustrating the potential of ONMF in text analysis. Further, we note that although the components of  $\mathbf{W}$  are not explicitly restricted to

be sparse, they tend to be: 90% of the  $\ell_2$  mass of each component is concentrated in approximately 100-200 entries (words) out of the roughly 102K present in the dataset.

## Appendix References

- [34] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [35] Aaron D Wyner. Random packings and coverings of the unit n-sphere. *Bell System Technical Journal*, 46(9):2111–2118, 1967.