

B Online Appendix

B.1 Proof of Remark 4

Claim: Let $\|\cdot\|$ be an arbitrary matrix norm. For any column-stochastic matrix P of nonnegative rank K we have

$$\mathcal{C}_K(P) \equiv \mathcal{C}_K \cap \left\{ C \in \mathbb{R}^{V \times V} \mid CP^{\text{row}} = P^{\text{row}} \right\} \neq \emptyset$$

if and only if

$$\min_{C \in \mathcal{C}_K} \|CP^{\text{row}} - P^{\text{row}}\| = 0.$$

Proof. We first show the “ \implies ” direction. Since $\mathcal{C}_K(P) \neq \emptyset$, then there exists $C^* \in \mathcal{C}_K$ such that $C^*P^{\text{row}} = P^{\text{row}}$. Since,

$$0 \leq \inf_{C \in \mathcal{C}_K} \|CP^{\text{row}} - P^{\text{row}}\| \leq \|C^*P^{\text{row}} - P^{\text{row}}\| = 0,$$

then

$$\inf_{C \in \mathcal{C}_K} \|CP^{\text{row}} - P^{\text{row}}\| = \|C^*P^{\text{row}} - P^{\text{row}}\| = 0.$$

Thus, the infimum is attained and

$$\min_{C \in \mathcal{C}_K} \|CP^{\text{row}} - P^{\text{row}}\| = 0.$$

For the “ \impliedby ” we note that if

$$\min_{C \in \mathcal{C}_K} \|CP^{\text{row}} - P^{\text{row}}\| = 0,$$

then, by definition, there exists $C^* \in \mathcal{C}_K$ such that

$$\|C^*P^{\text{row}} - P^{\text{row}}\| = 0.$$

But since $\|\cdot\|$ is a norm, this implies $C^*P^{\text{row}} - P^{\text{row}} = 0$. □

B.2 Proof of Remark 5

Let P, Q be column-stochastic matrices of dimension $V \times D$. Define the total-variation distance between P and Q as

$$\|P - Q\|_{\text{TV}} = \frac{1}{2} \sum_{v=1}^V \sum_{d=1}^D |p_{v,d} - q_{v,d}|.$$

This extends the typical definition of the total-variation distance for discrete distributions; see p. 48, Proposition 4.2 in Levin & Peres (2017).

Claim: Suppose that P is a column-stochastic matrix of nonnegative rank $K \leq \min\{V, D\}$ that a) does not admit an anchor-word factorization in the sense of Definition 2, and b) there exists some $\epsilon > 0$

$$p_v \equiv \sum_{d=1}^D p_{v,d} > \epsilon, \quad \forall v = 1, \dots, V.$$

Then, there is no sequence of matrices $\{P_i\}_{i \in \mathbb{N}}$ for which $P_i = A_i W_i$, $(A_i, W_i) \in \Theta_0$ and $\|P - P_i\|_{TV} \rightarrow 0$.

Proof. We establish this result by contradiction. Suppose there is a sequence $\{P_i\}_{i \in \mathbb{N}}$ for which $P_i = A_i W_i$, $(A_i, W_i) \in \Theta_0$ and $\|P - P_i\|_{TV} \rightarrow 0$. Theorem 1 shows that for each $i \in \mathbb{N}$, there exists a matrix $C_i \in \mathcal{C}_K$ such that

$$C_i P_i^{\text{row}} = P_i^{\text{row}}.$$

Let $\|\cdot\|$ denote the Frobenius norm. For any C_i satisfying $C P_i = P_i$ we have

$$\begin{aligned} \|C_i P_i^{\text{row}} - P_i^{\text{row}}\| &= \|C_i P_i^{\text{row}} - C_i P_i^{\text{row}} + C_i P_i^{\text{row}} - P_i^{\text{row}} + P_i^{\text{row}} - P_i^{\text{row}}\|, \\ &\leq \|C_i(P_i^{\text{row}} - P_i^{\text{row}})\| + \|C_i P_i^{\text{row}} - P_i^{\text{row}}\| + \|P_i^{\text{row}} - P_i^{\text{row}}\|, \\ &= \|C_i(P_i^{\text{row}} - P_i^{\text{row}})\| + \|P_i^{\text{row}} - P_i^{\text{row}}\|, \\ &\leq (\|C_i\| + 1) \cdot \|P_i^{\text{row}} - P_i^{\text{row}}\|. \end{aligned}$$

Consequently,

$$\inf_{C \in \mathcal{C}_K} \|C P_i^{\text{row}} - P_i^{\text{row}}\| \leq (\|C_i\| + 1) \cdot \|P_i^{\text{row}} - P_i^{\text{row}}\| \quad (62)$$

for every $i \in \mathbb{N}$. Because \mathcal{C}_K is bounded (as the matrices $C \in \mathcal{C}_K$ have elements in $[0, 1]$), then the sequence $\{\|C_i\|\}_{i \in \mathbb{N}}$ is bounded. Moreover,

$$\begin{aligned} \|P_i^{\text{row}} - P_i^{\text{row}}\| &= \sqrt{\sum_{d=1}^D \sum_{v=1}^V (p_{v,d}^{\text{row}} - p_{i,(v,d)}^{\text{row}})^2} \\ &\leq \sum_{d=1}^D \sum_{v=1}^V |p_{v,d}^{\text{row}} - p_{i,(v,d)}^{\text{row}}| \\ &= \sum_{d=1}^D \sum_{v=1}^V \left| \frac{p_{v,d}}{p_v} - \frac{p_{i,(v,d)}}{p_{iv}} \right|, \end{aligned}$$

where \mathbf{p}_v and \mathbf{p}_{iv} represent the row sums of \mathbf{P} and \mathbf{P}_i , respectively. Since

$$\left| \frac{\mathbf{p}_{v,d}}{\mathbf{p}_v} - \frac{\mathbf{p}_{i,(v,d)}}{\mathbf{p}_{iv}} \right| = \left| \frac{\mathbf{p}_{v,d}}{\mathbf{p}_v} - \frac{\mathbf{p}_{i,(v,d)}}{\mathbf{p}_v} + \frac{\mathbf{p}_{i,(v,d)}}{\mathbf{p}_v} - \frac{\mathbf{p}_{i,(v,d)}}{\mathbf{p}_{iv}} \right|,$$

then

$$\begin{aligned} \|\mathbf{P}^{\text{row}} - \mathbf{P}_i^{\text{row}}\| &\leq \sum_{d=1}^D \sum_{v=1}^V \frac{1}{\mathbf{p}_v} \cdot |\mathbf{p}_{v,d} - \mathbf{p}_{i,(v,d)}| \\ &\quad + \sum_{d=1}^D \sum_{v=1}^V \frac{\mathbf{p}_{i,(v,d)}}{\mathbf{p}_v \cdot \mathbf{p}_{iv}} \cdot |\mathbf{p}_{iv} - \mathbf{p}_v|. \end{aligned}$$

Since $\|\mathbf{P}_i - \mathbf{P}\|_{\text{TV}} \rightarrow 0$ implies that $|\mathbf{p}_{i,(v,d)} - \mathbf{p}_{v,d}| \rightarrow 0$ for all $v = 1, \dots, V$ and $d = 1, \dots, D$ then

$$\|\mathbf{P}^{\text{row}} - \mathbf{P}_i^{\text{row}}\| \rightarrow 0,$$

and, because of (62)

$$\inf_{\mathbf{C} \in \mathcal{C}_k} \|\mathbf{C}\mathbf{P}^{\text{row}} - \mathbf{P}^{\text{row}}\| = 0.$$

This implies, by Theorem 1 that \mathbf{P} admits an anchor-word factorization. A contradiction. □

B.3 Proof that $\inf_{\mathbf{C} \in \mathcal{C}_k} \|\mathbf{C}\widehat{\mathbf{P}}^{\text{row}} - \widehat{\mathbf{P}}^{\text{row}}\|$ is always attained.

Claim: Let $\|\cdot\|$ denote the Frobenius norm. For any column-stochastic, row normalized matrix \mathbf{P}^{row} ,

$$\inf_{\mathbf{C} \in \mathcal{C}_k} \|\mathbf{C}\mathbf{P}^{\text{row}} - \mathbf{P}^{\text{row}}\| = \min_{\mathbf{C} \in \mathcal{C}_k} \|\mathbf{C}\mathbf{P}^{\text{row}} - \mathbf{P}^{\text{row}}\|.$$

Proof. We want to show the minimum of $\|\mathbf{C}\mathbf{P}^{\text{row}} - \mathbf{P}^{\text{row}}\|$ is attainable in \mathcal{C}_k when the norm is Frobenius. By the extreme value theorem—e.g., Munkres (2000) Theorem 27.4 on page 174—it is sufficient to show function $f_{\mathbf{P}}(\mathbf{C}) \equiv \|\mathbf{C}\mathbf{P}^{\text{row}} - \mathbf{P}^{\text{row}}\|$ is continuous in \mathbf{C} over \mathcal{C}_k and that \mathcal{C}_k is compact. For the rest of the proof, we work with the topology induced by the Euclidean metric in \mathbb{R}^{V^2} , and the topology over $\mathbb{R}^{V \times V}$ induced by the Frobenius norm.

First, we show that $f_{\mathbf{P}}(\mathbf{C})$ is continuous. For any $\varepsilon > 0$, there exists $\delta = \varepsilon/\|\mathbf{P}^{\text{row}}\|$ such that if $\|\mathbf{C} - \mathbf{C}_0\| < \delta$, then

$$|\|\mathbf{C}\mathbf{P}^{\text{row}} - \mathbf{P}^{\text{row}}\| - \|\mathbf{C}_0\mathbf{P}^{\text{row}} - \mathbf{P}^{\text{row}}\|| \leq \|\mathbf{C}\mathbf{P}^{\text{row}} - \mathbf{C}_0\mathbf{P}^{\text{row}}\| \leq \|\mathbf{C} - \mathbf{C}_0\| \cdot \|\mathbf{P}^{\text{row}}\| < \varepsilon.$$

The first inequality holds due to the reverse triangle inequality and the second inequality comes from

the submultiplicativity of the Frobenius norm; see Horn & Johnson (2012) page 340.

Second, we show that the set \mathcal{C}_K is compact. It is sufficient to show \mathcal{C}_K is closed since it is a subset of a compact space $[0, 1]^{K \times K}$; see Munkres (2000) Theorem 26.2 on page 165. For the compactness of the space $[0, 1]^{K \times K}$, we rely on facts that the space $[0, 1]^{K^2}$ is compact and the image of a compact space under a continuous map is compact—see, for example, Munkres (2000) Theorem 26.5 on page 166—where we depend on the continuous bijection $h_{ij}(\tilde{C}) = \tilde{C}_{V(i-1)+j}$ for any $\tilde{C} \in [0, 1]^{K^2}$.

For a sequence $\{C_n \in \mathcal{C}_K\}_{n \in \mathbb{N}}$ that converges, we want to show its limit C is in \mathcal{C}_K . Notice the matrix converges in the Frobenius norm is equivalent to entry-wise convergences in absolute values. That is, if $\lim_{n \rightarrow \infty} C_n = C$, for any $\varepsilon > 0$, there exists N such that if $n > N$, $|C_{n,ij} - C_{i,j}| \leq \|C_n - C\| \leq \varepsilon$. Also, if $\lim_{n \rightarrow \infty} C_{n,ij} = C_{ij}$ for all i and j , for any $\varepsilon/V > 0$, there exists $\{N_{ij}\}$ such that if $n > \sup\{N_{ij}\}$, $\|C_n - C\| \leq \sqrt{V^2(\frac{\varepsilon}{V})^2} = \varepsilon$. The last inequality is from the definition of the Frobenius norm.

Finally, by the definition of the convergence, the diagonal elements are bounded by 0 and 1, and the off-diagonal elements also share the same bounds because if $C_{n,ij} \leq C_{jj}$, $\lim C_{n,ij} \leq C_{jj}$. Therefore, C is in \mathcal{C}_K and \mathcal{C}_K is closed. □

B.4 An anchor word factorization always exists when $K = 2 \leq \min\{V, D\}$

B.4.1 Proof using condition (19) of Theorem 1

Let P be a nonnegative column-stochastic matrix of rank $K = 2 \leq \min\{V, D\}$. Thomas (1974) has shown that every rank two nonnegative matrix admits a nonnegative matrix factorization. Let (A, W) be the nonnegative matrices in $\mathbb{R}^{2 \times V} \times \mathbb{R}^{2 \times D}$ that factorize P ; that is $P = AW$.

Without loss of generality we can assume that A and W are column stochastic (that is, their columns add up to one). Also, suppose that the first term in the vocabulary solves the problem $c_1 \equiv \min_{v \in V} \alpha_{v2}/\alpha_{v1}$. That is, we assume that the first term of the vocabulary receives the lowest possible probability under topic two, relative to the probability that the same term receives under topic one. Analogously, suppose that the second term in the vocabulary solves $c_2 \equiv \min_{v \in V} \alpha_{v1}/\alpha_{v2}$. Note that if A were not organized in such a way, we could always permute the rows of A to achieve this structure. Note also that the ratios involving α_{v1} and α_{v2} are always well defined because none of the rows of P equal zero.

We will make use of the 2×2 matrix

$$T \equiv \begin{pmatrix} \frac{1}{1-c_2} & -\frac{c_1}{1-c_1} \\ -\frac{c_2}{1-c_2} & \frac{1}{1-c_1} \end{pmatrix},$$

where c_1 and c_2 are defined in the previous paragraph. Because A has rank two, both $c_1, c_2 \in (0, 1)$.

This implies that T is well defined; that its determinant is strictly positive, and that T^{-1} is a column-stochastic matrix.

In a slight abuse of notation, write A as the following block matrix

$$A = \begin{bmatrix} \underbrace{A^*}_{2 \times 2} \\ \underbrace{\tilde{A}}_{V-2 \times 2} \end{bmatrix}.$$

Consider then the $V \times V$ matrix given by

$$C \equiv \begin{bmatrix} \mathbb{I}_2 & \mathbf{0}_{2 \times V-2} \\ (\mathcal{R}_{\tilde{A}W})^{-1} \tilde{A} T \mathcal{R}_{T^{-1}W} & \mathbf{0}_{V-2 \times V-2} \end{bmatrix}. \quad (63)$$

We will show that this matrix satisfies the necessary and sufficient condition for anchor word factorization in Theorem 1.

We first show that C is an element of the set C_2 defined in Equation (17). Note first that $\text{Tr}(C) = 2$ and that the diagonal elements of the matrix C are either 0 or 1. Thus, we only need to show that the elements of the matrix

$$(\mathcal{R}_{\tilde{A}W})^{-1} \tilde{A} T \mathcal{R}_{T^{-1}W} \quad (64)$$

are nonnegative and bounded above by one.

We first show that the elements of (64) are nonnegative. Note that $\tilde{A}W$ (which corresponds to the lower $V - 2 \times D$ block of P) is a nonnegative matrix, which implies $\mathcal{R}_{\tilde{A}W}$ is nonnegative. Note also that because T^{-1} is column stochastic, then $T^{-1}W$ is a column-stochastic matrix. Finally, since \tilde{A} is column stochastic and $c_1, c_2 \in (0, 1)$, it follows that $\tilde{A}T$ is nonnegative.

We then show that the elements of (64) are bounded above by one. Since, by definition, \mathcal{R}_M is the diagonal matrix that contains the row sums of a matrix M , algebra shows that

$$\mathcal{R}_{\tilde{A}W} = \mathcal{R}_{(\tilde{A}T)(T^{-1}W)} = \mathcal{R}_{\tilde{A}T \mathcal{R}_{T^{-1}W}}.$$

Thus, the elements of the $V - 2 \times 2$ matrix (64) are bounded above by one. This shows that C is an element of the set C_2 .

Finally, we show that C satisfies the equation $CP^{\text{row}} = P^{\text{row}}$. Using the block matrix representation of A

$$P^{\text{row}} = \begin{pmatrix} (A^*W)^{\text{row}} \\ (\tilde{A}W)^{\text{row}} \end{pmatrix}.$$

The definition of C in Equation (63) implies

$$\begin{aligned} CP^{\text{row}} &= \begin{pmatrix} (A^*W)^{\text{row}} \\ (\mathcal{R}_{\tilde{A}W})^{-1} \tilde{A}T \mathcal{R}_{T^{-1}W} (A^*W)^{\text{row}} \end{pmatrix}, \\ &= \begin{pmatrix} (A^*W)^{\text{row}} \\ (\mathcal{R}_{\tilde{A}W})^{-1} \tilde{A}T \mathcal{R}_{T^{-1}W} \left((A^*T) (T^{-1}W) \right)^{\text{row}} \end{pmatrix}. \end{aligned}$$

By construction, A^*T is a diagonal matrix, which implies

$$\left((A^*T) (T^{-1}W) \right)^{\text{row}} = \left((T^{-1}W) \right)^{\text{row}} = \mathcal{R}_{T^{-1}W} T^{-1}W.$$

Thus, we conclude that $CP^{\text{row}} = P^{\text{row}}$, and thus $C \in \mathcal{C}_2(P)$. Theorem 1 thus implies that any matrix P of rank $K = 2$ admits an anchor word factorization.

B.4.2 Explicit anchor word factorization when $K = 2 \leq \min\{V, D\}$

The proof of Theorem 1 gives a simple formula to obtain the anchor word factorization of \mathbb{P} from $C \in \mathcal{C}_2(P)$. In particular, if we start out with the factors (A, W) that were used in the previous subsection, the proof of Theorem 1 implies that the column-normalized version of the $V \times K$ matrix

$$\begin{bmatrix} \mathbb{I}_K \\ \tilde{A}T \mathcal{R}_{T^{-1}W} \mathcal{R}_{A^*W}^{-1} \end{bmatrix} \quad (65)$$

provides an anchor word factorization of P . Since A^*T is diagonal and column stochastic, then the matrix in (65) equals

$$\begin{bmatrix} A^*T \\ \tilde{A}T \end{bmatrix} (A^*T)^{-1},$$

where we have used

$$\mathcal{R}_{A^*W} = \mathcal{R}_{A^*T T^{-1}W} = A^*T \mathcal{R}_{T^{-1}W}.$$

Thus,

$$A_0 = \begin{bmatrix} A^*T \\ \tilde{A}T \end{bmatrix}$$

and $W_0 \equiv T^{-1}W$ provide an anchor word factorization of P .

B.5 Anchor word factorization does not always exist $V = 4, K = D = 3$

B.5.1 Example

In this section we show that any matrix P of the form

$$P = \begin{pmatrix} \alpha & 0 & 0 \\ 0 & \gamma & 0 \\ 0 & 1 - \gamma & 1 - \beta \\ 1 - \alpha & 0 & \beta \end{pmatrix},$$

for $\alpha, \beta, \gamma \in (0, 1)$ does not admit an anchor word factorization.

The row-normalized version of P is given by:

$$P^{\text{row}} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & \frac{1-\gamma}{2-\gamma-\beta} & \frac{1-\beta}{2-\gamma-\beta} \\ \frac{1-\alpha}{1-\alpha+\beta} & 0 & \frac{\beta}{1-\alpha+\beta} \end{pmatrix},$$

We define the set $\tilde{\mathcal{C}}_K$ to be the set of $V \times V$ matrices of the form

$$\begin{bmatrix} \mathbb{I}_K & 0_{K \times V-K} \\ M & 0_{V-K \times K} \end{bmatrix},$$

where $M \geq 0$ is a row-normalized matrix (with rows different from zero, so that row-normalization is always well defined). From Lemma 1, we want to show there does not exist $C \in \tilde{\mathcal{C}}_K$ and a row permutation matrix Π such that $CP^{\text{row}} = \Pi P^{\text{row}}$.

Since $K = 3$ we can argue that it is only relevant to focus on four classes of permutations (which are indexed by the row of P^{row} that is placed at the bottom of the permuted matrix). Without loss of generality, we can focus on

$$P_1^{\text{row}} = \begin{pmatrix} \frac{1-\alpha}{1-\alpha+\beta} & 0 & \frac{\beta}{1-\alpha+\beta} \\ 0 & 1 & 0 \\ 0 & \frac{1-\gamma}{2-\gamma-\beta} & \frac{1-\beta}{2-\gamma-\beta} \\ 1 & 0 & 0 \end{pmatrix},$$

$$\begin{aligned}
\mathbf{P}_2^{\text{row}} &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & \frac{1-\gamma}{2-\gamma-\beta} & \frac{1-\beta}{2-\gamma-\beta} \\ \frac{1-\alpha}{1-\alpha+\beta} & 0 & \frac{\beta}{1-\alpha+\beta} \\ 0 & 1 & 0 \end{pmatrix}, \\
\mathbf{P}_3^{\text{row}} &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ \frac{1-\alpha}{1-\alpha+\beta} & 0 & \frac{\beta}{1-\alpha+\beta} \\ 0 & \frac{1-\gamma}{2-\gamma-\beta} & \frac{1-\beta}{2-\gamma-\beta} \end{pmatrix}, \\
\mathbf{P}_4^{\text{row}} &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & \frac{1-\gamma}{2-\gamma-\beta} & \frac{1-\beta}{2-\gamma-\beta} \\ \frac{1-\alpha}{1-\alpha+\beta} & 0 & \frac{\beta}{1-\alpha+\beta} \end{pmatrix}.
\end{aligned}$$

Note there is no $\mathbf{C} \in \tilde{\mathcal{C}}_{\mathbf{K}}$ such that $\mathbf{C}\mathbf{P}_i^{\text{row}} = \mathbf{P}_i^{\text{row}}$ for $i = 1, 2$, since this would require some elements of \mathbf{M} to be strictly above one.

Consider now the matrices $\mathbf{P}_3^{\text{row}}$ and $\mathbf{P}_4^{\text{row}}$. We can focus on $\mathbf{P}_3^{\text{row}}$, since the argument for the other matrix is entirely analogous. Let the elements of \mathbf{M} , which is a 1×3 matrix, be denoted as $[m_1, m_2, m_3]$. In order for the first element of the last row of $\mathbf{P}_3^{\text{row}}$ (which equals zero) to be a convex combination of the first three rows it is necessary to have $m_1 = m_3 = 0$. However, this implies that the last element of the fourth row of $\mathbf{P}_3^{\text{row}}$ (which equals $1 - \beta/2 - \gamma - \beta$) cannot be obtained as a convex combination of the first three rows, whenever $\beta \in (0, 1)$. Therefore there does not exist $\mathbf{C} \in \tilde{\mathcal{C}}_{\mathbf{K}}$ such that $\mathbf{C}\mathbf{P}_3^{\text{row}} = \mathbf{P}_3^{\text{row}}$. Since the argument for $\mathbf{P}_4^{\text{row}}$ is analogous, we conclude that the anchor word factorization does not exist for \mathbf{P} .

B.6 Upper bound for $q_{1-\alpha}^*(\mathbf{V}, \mathbf{K}, \mathbf{D}, \bar{\mathbf{N}}_{\mathbf{D}})$

Lemma 4. *Let $\|\cdot\|$ denote the Frobenius norm. For any $\alpha \in (0, 1)$*

$$q_{1-\alpha}^*(\mathbf{V}, \mathbf{D}, \mathbf{K}, \bar{\mathbf{N}}_{\mathbf{D}}) \leq \sup_{\mathbf{C} \in \mathcal{C}_{\mathbf{K}}} \|\mathbf{C} - \mathbb{I}_{\mathbf{V}}\| \cdot \tilde{q}_{1-\alpha}^*(\mathbf{V}, \mathbf{D}, \mathbf{K}, \bar{\mathbf{N}}_{\mathbf{D}}), \quad (66)$$

where

$$\tilde{q}_{1-\alpha}^*(\mathbf{V}, \mathbf{D}, \mathbf{K}, \bar{\mathbf{N}}_{\mathbf{D}}) = \sup_{(\mathbf{A}, \mathbf{W}) \in \Theta_0} \tilde{q}_{1-\alpha}(\mathbf{A}\mathbf{W}, \mathbf{V}, \mathbf{D}, \mathbf{K}, \bar{\mathbf{N}}_{\mathbf{D}})$$

and

$$\tilde{q}_{1-\alpha}(AW, V, D, K, \bar{N}_D) = \inf \left\{ q \in \mathbb{R}_+ \mid \mathbb{P}_{AW} \left(\|\hat{\mathbf{P}}^{\text{row}} - (AW)^{\text{row}}\| \leq q \right) \geq 1 - \alpha \right\}.$$

Proof. By definition—see Section 3.2— $q_{1-\alpha}(AW, V, D, K, \bar{N}_D)$ is the $1 - \alpha$ quantile of the test statistic $T(Y)$ under the distribution $\mathbf{P} = AW$, $(A, W) \in \Theta_0$. Thus:

$$q_{1-\alpha}(AW, V, D, K, \bar{N}_D) = \inf \left\{ q \in \mathbb{R}_+ \mid \mathbb{P}_{AW} (T(Y) < q) \geq 1 - \alpha \right\}$$

Let $C_P \in \mathcal{C}_K$ be the matrix for which $C_P \mathbf{P}^{\text{row}} - (AW)^{\text{row}} = \mathbf{0}$ (such a matrix exists by Theorem 1). Since the test statistic $T(Y)$ equals $\min_{C \in \mathcal{C}_K} \|C \hat{\mathbf{P}}^{\text{row}} - \hat{\mathbf{P}}^{\text{row}}\|$, it follows that

$$\begin{aligned} T(Y) &\leq \|C_P \hat{\mathbf{P}}^{\text{row}} - \hat{\mathbf{P}}^{\text{row}}\| \\ &= \|C_P \hat{\mathbf{P}}^{\text{row}} - C_P \mathbf{P}^{\text{row}} + C_P \mathbf{P}^{\text{row}} - \mathbf{P}^{\text{row}} + \mathbf{P}^{\text{row}} - \hat{\mathbf{P}}^{\text{row}}\| \\ &= \|(C_P - \mathbb{I}_V) (\hat{\mathbf{P}}^{\text{row}} - \mathbf{P}^{\text{row}})\| \\ &\leq \sup_{C \in \mathcal{C}_K} \|C - \mathbb{I}_V\| \cdot \|\hat{\mathbf{P}}^{\text{row}} - \mathbf{P}^{\text{row}}\|, \end{aligned}$$

where the last inequality follows from the submultiplicativity of Frobenius norm. This inequality implies that

$$Q_1 \equiv \left\{ q \in \mathbb{R}_+ \mid \mathbb{P}_{AW} \left(\sup_{C \in \mathcal{C}_K} \|C - \mathbb{I}_V\| \cdot \|\hat{\mathbf{P}}^{\text{row}} - \mathbf{P}^{\text{row}}\| \leq q \right) \geq 1 - \alpha \right\}$$

is a subset of

$$Q_0 \equiv \left\{ q \in \mathbb{R}_+ \mid \mathbb{P}_{AW} (T(Y) < q) \geq 1 - \alpha \right\}.$$

Therefore,

$$q_{1-\alpha}(AW, V, D, K, \bar{N}_D) = \inf Q_0 \leq \inf Q_1. \quad (67)$$

Define $C^*(V, K) \equiv \sup_{C \in \mathcal{C}_K} \|C - \mathbb{I}_V\|$. We want to show that

$$\inf Q_1 \leq C^*(V, K) \cdot \tilde{q}_{1-\alpha}(AW, V, D, K, \bar{N}_D).$$

Let

$$Q_2 \equiv \left\{ q \in \mathbb{R}_+ \mid \mathbb{P}_{AW} \left(\|\hat{\mathbf{P}}^{\text{row}} - \mathbf{P}^{\text{row}}\| \leq q \right) \geq 1 - \alpha \right\},$$

and note that, by definition,

$$\tilde{q}_{1-\alpha}(AW, V, D, K, \bar{N}_D) = \inf Q_2.$$

By definition of infimum, there exists a sequence $\{q_n\}_{n \in \mathbb{N}} \subseteq Q_2$ such that

$$\lim_{n \rightarrow \infty} q_n = \tilde{q}_{1-\alpha}(AW, V, D, K, \bar{N}_D). \quad (68)$$

For each q_n we have that

$$(C^*(V, K) \cdot q_n) \in Q_1.$$

Consequently,

$$\inf Q_1 \leq C^*(V, K) \cdot q_n$$

for all $n \in \mathbb{N}$. We thus conclude by (68) that

$$\inf Q_1 \leq C^*(V, K) \cdot \tilde{q}_{1-\alpha}(AW, V, D, K, \bar{N}_D)$$

and by (67) that

$$q_{1-\alpha}(AW, V, D, K, \bar{N}_D) \leq C^*(V, K) \cdot \tilde{q}_{1-\alpha}(AW, V, D, K, \bar{N}_D).$$

Taking the supremum on both sides over $(A, W) \in \Theta_0$ gives the desired result. □

B.7 Estimation error of different estimators

In this section we discuss two alternative estimators for P^{row} . Here is a description of the estimators and the results we derive:

1. *Nuclear-Norm Minimizer*: Let \hat{P}_{nuc} be the estimator suggested by McRae & Davenport (2021), Section 2.3, Theorem 2.2, p. 712. The following proposition follows from their Theorem 2.2:

Proposition 3. *Let $0 < \gamma < 1$ be an arbitrary scalar. For any (A, W) such that $p_v(A, W)/D \geq \gamma/V$*

$$\|\hat{P}_{\text{nuc}}^{\text{row}} - (AW)^{\text{row}}\|_F \leq 4 \sqrt{\frac{16}{\gamma^2} \cdot \frac{V^{3/2} \cdot \ln((D+V)/\epsilon) \cdot K}{N_{\text{min}}}} \quad (69)$$

with probability at least $1 - \epsilon$.

2. *Minimax Estimator for the columns*: Let \hat{P}_{min} the $V \times D$ matrix with (v, d) -entry given by

$(\sqrt{N_d}/V + n_{vd})/(\sqrt{N_d} + N_d)$. Let $\hat{P}_{\min}^{\text{row}}$ the row-normalized version of this estimator. In Section B.7.2 below we establish the following proposition:

Proposition 4. *Let $0 < \gamma < 1$ be arbitrary scalars. For any (A, W) such that $p_v(A, W)/D \geq \gamma/V$*

$$\|\hat{P}_{\min}^{\text{row}} - (AW)^{\text{row}}\|_F \leq \sqrt{\frac{8 \left(1 - \frac{1}{V}\right)}{\gamma^2 \cdot \epsilon} \cdot \frac{V^2}{N_{\min} + 2N_{\min}^{1/2} + 1}} \quad (70)$$

with probability at least $1 - \epsilon$.

The estimator that row-normalizes that minimax estimator is expected to satisfy the high-level assumption in (26) provided

$$\frac{V^2}{N_{\min} + 2N_{\min}^{1/2} + 1}$$

is small. Here, we rely on the same technique as Proposition 3 to derive the rate. We can also provide better rates with an order of

$$\frac{V^2}{D \cdot (N_{\min} + 2N_{\min}^{1/2} + 1)}$$

with other assumptions about probability design and other techniques.

Outline for this section: Let \hat{P} be an arbitrary estimator of the population term-document frequency matrix, P . Just as we did in the main body of the paper, define $\hat{P}^{\text{row}} \equiv \mathcal{R}_{\hat{P}}^{-1} \hat{P}$ and $P^{\text{row}} \equiv \mathcal{R}_P^{-1} P$. We establish a series of results that will allow us to provide finite-sample bounds for $\|\hat{P}^{\text{row}} - P^{\text{row}}\|_F$.

Lemma 5 below shows that in order to upper-bound the estimation error $\|\hat{P}^{\text{row}} - P^{\text{row}}\|_F$ we can analyze the terms

$$\|\mathcal{R}_P^{-1}(P - \hat{P})\|_F \quad (71)$$

and

$$\|(\mathcal{R}_{\hat{P}}^{-1} - \mathcal{R}_P^{-1})\hat{P}\|_F. \quad (72)$$

Lemma 6 uses Markov's inequality to provide an upper bound for the term in (71). Lemma 7 provides an upper bound for the term in (72). The bounds do not depend on the specific form of \hat{P} as long as the second moments of the estimator exist.

Lemma 5. *If $\|\mathcal{R}_P^{-1}(P - \hat{P})\|_F \leq \delta_1$ with probability at least $1 - \epsilon/2$, and $\|(\mathcal{R}_{\hat{P}}^{-1} - \mathcal{R}_P^{-1})\hat{P}\|_F \leq \delta_2$ with*

probability at least $1 - \epsilon/2$, then with probability at least $1 - \epsilon$,

$$\|\hat{\mathbf{P}}^{\text{row}} - \mathbf{P}^{\text{row}}\|_F \leq 2 \max\{\delta_1, \delta_2\}$$

Proof. Algebra shows that

$$\begin{aligned} \|\hat{\mathbf{P}}^{\text{row}} - \mathbf{P}^{\text{row}}\|_F &= \|\mathcal{R}_{\hat{\mathbf{p}}}^{-1} \hat{\mathbf{P}} - \mathcal{R}_{\mathbf{p}}^{-1} \mathbf{P}\|_F \\ &= \|\mathcal{R}_{\hat{\mathbf{p}}}^{-1} \hat{\mathbf{P}} - \mathcal{R}_{\mathbf{p}}^{-1} \hat{\mathbf{P}} + \mathcal{R}_{\mathbf{p}}^{-1} \hat{\mathbf{P}} - \mathcal{R}_{\mathbf{p}}^{-1} \mathbf{P}\|_F \\ &\leq \|\mathcal{R}_{\hat{\mathbf{p}}}^{-1} \hat{\mathbf{P}} - \mathcal{R}_{\mathbf{p}}^{-1} \hat{\mathbf{P}}\|_F + \|\mathcal{R}_{\mathbf{p}}^{-1} \hat{\mathbf{P}} - \mathcal{R}_{\mathbf{p}}^{-1} \mathbf{P}\|_F \\ &= \|\mathcal{R}_{\hat{\mathbf{p}}}^{-1} (\hat{\mathbf{P}} - \mathbf{P})\|_F + \|(\mathcal{R}_{\hat{\mathbf{p}}}^{-1} - \mathcal{R}_{\mathbf{p}}^{-1}) \hat{\mathbf{P}}\|_F, \end{aligned}$$

where the inequality comes from the triangle inequality.

The inequality above implies that for any constant c we have

$$\mathbb{P}(\|\hat{\mathbf{P}}^{\text{row}} - \mathbf{P}^{\text{row}}\|_F > c) \leq \mathbb{P}(\|\mathcal{R}_{\hat{\mathbf{p}}}^{-1} (\hat{\mathbf{P}} - \mathbf{P})\|_F + \|(\mathcal{R}_{\hat{\mathbf{p}}}^{-1} - \mathcal{R}_{\mathbf{p}}^{-1}) \hat{\mathbf{P}}\|_F > c).$$

Moreover, the right-hand side of the equation above is upper-bounded by

$$\mathbb{P}(\|\mathcal{R}_{\hat{\mathbf{p}}}^{-1} (\hat{\mathbf{P}} - \mathbf{P})\|_F > c/2 \text{ or } \|(\mathcal{R}_{\hat{\mathbf{p}}}^{-1} - \mathcal{R}_{\mathbf{p}}^{-1}) \hat{\mathbf{P}}\|_F > c/2).$$

The subadditivity of probability measures then implies

$$\begin{aligned} \mathbb{P}(\|\hat{\mathbf{P}}^{\text{row}} - \mathbf{P}^{\text{row}}\|_F > c) &\leq \mathbb{P}(\|\mathcal{R}_{\hat{\mathbf{p}}}^{-1} (\hat{\mathbf{P}} - \mathbf{P})\|_F > c/2) \\ &\quad + \mathbb{P}(\|(\mathcal{R}_{\hat{\mathbf{p}}}^{-1} - \mathcal{R}_{\mathbf{p}}^{-1}) \hat{\mathbf{P}}\|_F > c/2). \end{aligned}$$

Take $c = 2 \max\{\delta_1, \delta_2\}$ and note that

$$\mathbb{P}(\|\mathcal{R}_{\hat{\mathbf{p}}}^{-1} (\hat{\mathbf{P}} - \mathbf{P})\|_F > \max\{\delta_1, \delta_2\}) \leq \mathbb{P}(\|\mathcal{R}_{\hat{\mathbf{p}}}^{-1} (\hat{\mathbf{P}} - \mathbf{P})\|_F > \delta_1) < \epsilon/2,$$

and analogously $\mathbb{P}(\|(\mathcal{R}_{\hat{\mathbf{p}}}^{-1} - \mathcal{R}_{\mathbf{p}}^{-1}) \hat{\mathbf{P}}\|_F > \max\{\delta_1, \delta_2\}) < \epsilon/2$. \square

Lemma 6. *Suppose that the second moments of $\hat{\mathbf{p}}_{\mathbf{v}\mathbf{d}}$ exist for $\mathbf{v} = 1, \dots, V$ and $\mathbf{d} = 1, \dots, D$. Then with probability at least $1 - \epsilon$*

$$\|\mathcal{R}_{\hat{\mathbf{p}}}^{-1} (\hat{\mathbf{P}} - \mathbf{P})\|_F \leq \frac{1}{\mathbf{p}_{\mathbf{v}\mathbf{min}}} \sqrt{\frac{\sum_{\mathbf{v}=1}^V \sum_{\mathbf{d}=1}^D \mathbb{E} [(\hat{\mathbf{p}}_{\mathbf{v}\mathbf{d}} - \mathbf{p}_{\mathbf{v}\mathbf{d}})^2]}{\epsilon}},$$

where the expectation \mathbb{E} is taken under the true data generating process \mathbf{P} .

Proof. The definition of Frobenius norm implies that for any $x > 0$

$$\begin{aligned} \mathbb{P}(\|\mathcal{R}_{\hat{\mathbf{P}}}^{-1}(\hat{\mathbf{P}} - \mathbf{P})\|_{\text{F}} > x) &= \mathbb{P}\left(\sum_{\mathbf{v}} \sum_{\mathbf{d}} \frac{1}{p_{\mathbf{v}}^2} (p_{\mathbf{v}\mathbf{d}} - \hat{p}_{\mathbf{v}\mathbf{d}})^2 > x^2\right) \\ &\leq \mathbb{P}\left(\frac{1}{p_{\mathbf{v}\min}^2} \sum_{\mathbf{v}} \sum_{\mathbf{d}} (p_{\mathbf{v}\mathbf{d}} - \hat{p}_{\mathbf{v}\mathbf{d}})^2 > x^2\right) \\ &\leq \frac{\sum_{\mathbf{v}} \sum_{\mathbf{d}} \mathbb{E}(p_{\mathbf{v}\mathbf{d}} - \hat{p}_{\mathbf{v}\mathbf{d}})^2}{p_{\mathbf{v}\min}^2 x^2}, \end{aligned}$$

where the last step follows from Markov's inequality. Taking x to be

$$\sqrt{\frac{\sum_{\mathbf{v}=1}^{\mathbf{V}} \sum_{\mathbf{d}=1}^{\mathbf{D}} \mathbb{E}[(\hat{p}_{\mathbf{v}\mathbf{d}} - p_{\mathbf{v}\mathbf{d}})^2]}{p_{\mathbf{v}\min}^2 \epsilon}}$$

completes the proof. □

Lemma 7. *Suppose that the second moments of $\hat{p}_{\mathbf{v}\mathbf{d}}$ exist for $\mathbf{v} = 1, \dots, \mathbf{V}$ and $\mathbf{d} = 1, \dots, \mathbf{D}$. Then with probability at least $1 - \epsilon$*

$$\|(\mathcal{R}_{\hat{\mathbf{P}}}^{-1} - \mathcal{R}_{\mathbf{P}}^{-1})\hat{\mathbf{P}}\|_{\text{F}} \leq \frac{1}{p_{\mathbf{v}\min}} \sqrt{\frac{\sum_{\mathbf{v}=1}^{\mathbf{V}} \mathbb{E}[(p_{\mathbf{v}} - \hat{p}_{\mathbf{v}})^2]}{\epsilon}}$$

where the expectation \mathbb{E} is taken under the true data generating process \mathbf{P} , and $p_{\mathbf{v}} \equiv \sum_{\mathbf{d}=1}^{\mathbf{d}} p_{\mathbf{v}\mathbf{d}}$, $\hat{p}_{\mathbf{v}} \equiv \sum_{\mathbf{d}=1}^{\mathbf{d}} \hat{p}_{\mathbf{v}\mathbf{d}}$.

Proof.

$$\begin{aligned} \|(\mathcal{R}_{\hat{\mathbf{P}}}^{-1} - \mathcal{R}_{\mathbf{P}}^{-1})\hat{\mathbf{P}}\|_{\text{F}} &= \left[\sum_{\mathbf{v}} \sum_{\mathbf{d}} \left(\frac{1}{p_{\mathbf{v}}} - \frac{1}{\hat{p}_{\mathbf{v}}}\right)^2 \hat{p}_{\mathbf{v}\mathbf{d}}^2 \right]^{1/2} \\ &= \left[\sum_{\mathbf{v}} \sum_{\mathbf{d}} \frac{(\hat{p}_{\mathbf{v}} - p_{\mathbf{v}})^2}{p_{\mathbf{v}}^2 \hat{p}_{\mathbf{v}}^2} \hat{p}_{\mathbf{v}\mathbf{d}}^2 \right]^{1/2} \\ &= \left[\sum_{\mathbf{v}} \frac{(\hat{p}_{\mathbf{v}} - p_{\mathbf{v}})^2}{p_{\mathbf{v}}^2 \hat{p}_{\mathbf{v}}^2} \sum_{\mathbf{d}} \hat{p}_{\mathbf{v}\mathbf{d}}^2 \right]^{1/2} \\ &\leq \left[\sum_{\mathbf{v}} \frac{(\hat{p}_{\mathbf{v}} - p_{\mathbf{v}})^2}{p_{\mathbf{v}}^2 \hat{p}_{\mathbf{v}}^2} \hat{p}_{\mathbf{v}}^2 \right]^{1/2} \end{aligned}$$

$$\leq \left[\frac{1}{p_{v\min}^2} \sum_v (\hat{p}_v - p_v)^2 \right]^{1/2}.$$

The inequality above holds since $(\sum_d \hat{p}_{vd}^2)^{1/2} \leq \sum_d \hat{p}_{vd} = \hat{p}_v$.

Then, for any $x > 0$

$$\begin{aligned} \mathbb{P}(\|(\mathcal{R}_{\hat{p}}^{-1} - \mathcal{R}_p^{-1})\hat{P}\|_F > x) &\leq \mathbb{P}\left(\frac{1}{p_{v\min}^2} \sum_v (\hat{p}_v - p_v)^2 > x^2\right) \\ &\leq \frac{\sum_v \mathbb{E}((\hat{p}_v - p_v)^2)}{p_{v\min}^2 x^2}, \end{aligned}$$

where the last line follows by Markov's inequality. Taking

$$x = \frac{1}{p_{v\min}} \sqrt{\frac{\sum_v \mathbb{E}(p_v - \hat{p}_v)^2}{\epsilon}},$$

yields the desired result. □

B.7.1 Estimation error of $P_{\text{freq}}^{\text{row}}$

Proof of Proposition 2. In a slight abuse of notation, let \hat{P} denote the $V \times D$ matrix with (v, d) -entry given by n_{vd}/N_d . Let \hat{P}^{row} the row-normalized version of this estimator.

Note that

$$\begin{aligned} \sum_v \sum_d \mathbb{E}[(\hat{p}_{vd} - p_{vd})^2] &= \sum_v \sum_d \frac{p_{vd}(1 - p_{vd})}{N_d} \\ &\leq \sum_v \sum_d \frac{p_{vd}(1 - p_{vd})}{N_{\min}} \\ &= \sum_d \frac{1 - \sum_v p_{vd}^2}{N_{\min}} \\ &\leq \frac{D(1 - \frac{1}{V})}{N_{\min}}. \end{aligned}$$

The first equality holds because n_{vd} is a binomial distribution with parameter N_d and p_{vd} . The second equality holds since the $\sum_v p_{vd} = 1$. The second inequality comes from the fact that

$$\min_{p_{1d}, \dots, p_{Vd}} \sum_v p_{vd}^2 \quad \text{s.t.} \quad \sum_v p_{vd} = 1$$

equals $1/V$. Therefore, by Lemma 6 with probability at least $1 - \epsilon/2$

$$\|\mathcal{R}_P^{-1}(P - \hat{P})\|_F \leq \frac{1}{p_{v\min}} \sqrt{\frac{2D(1 - \frac{1}{V})}{N_{\min}\epsilon}}.$$

Moreover, since by assumption, $p_{v\min}/D \geq \gamma/V$, we have that

$$\|\mathcal{R}_P^{-1}(P - \hat{P})\|_F \leq \sqrt{\frac{2V^2(1 - \frac{1}{V})}{\gamma^2 D N_{\min}\epsilon}}.$$

Lemma 7 implies that with probability at least $1 - \epsilon/2$

$$\begin{aligned} \|(\mathcal{R}_{\hat{P}}^{-1} - \mathcal{R}_P^{-1})\hat{P}\|_F &\leq \frac{1}{p_{v\min}} \sqrt{\frac{2 \sum_v \mathbb{E}(\mathbf{p}_v - \hat{\mathbf{p}}_v)^2}{\epsilon}} \\ &= \frac{1}{p_{v\min}} \sqrt{\frac{2 \sum_v \sum_d \mathbb{E}[(\hat{p}_{vd} - p_{vd})]^2}{\epsilon}} \\ &= \frac{1}{p_{v\min}} \sqrt{\frac{2D(1 - \frac{1}{V})}{N_{\min}\epsilon}} \\ &\leq \sqrt{\frac{2V^2(1 - \frac{1}{V})}{\gamma^2 D N_{\min}\epsilon}}, \end{aligned}$$

where the second equality holds because the estimators \hat{p}_{vd} are unbiased and they are also independent across documents.

Finally, Lemma 5, implies that if \hat{P}^{row} is based on the row-normalization of the empirical frequencies then

$$\|\hat{P}^{\text{row}} - (AW)^{\text{row}}\|_F \leq \sqrt{\frac{8 \left(1 - \frac{1}{V}\right)}{\gamma^2 \cdot \epsilon} \cdot \frac{V^2}{N_{\min} \cdot D}}$$

with probability at least $1 - \epsilon$. □

B.7.2 Estimation error of P_{\min}^{row}

Proof of Proposition 4. In a slight abuse of notation, let \hat{P} denote the $V \times D$ matrix with (v, d) -entry given by $(\sqrt{N_d}/V + n_{vd})/(\sqrt{N_d} + N_d)$. Let \hat{P}^{row} be the row-normalized version of this estimator.

As above, we show that

$$\sum_v \sum_d \mathbb{E}[(\hat{p}_{vd} - p_{vd})]^2 = \sum_v \sum_d \frac{N_d p_{vd} - \frac{2N_d p_{vd}}{V} + \frac{N_d}{V^2}}{(\sqrt{N_d} + N_d)^2}$$

$$\begin{aligned}
&\leq \sum_{\mathbf{v}} \sum_{\mathbf{d}} \frac{\mathbf{p}_{\mathbf{v}\mathbf{d}} - \frac{2\mathbf{p}_{\mathbf{v}\mathbf{d}}}{\mathbf{V}} + \frac{1}{\mathbf{V}^2}}{\mathbf{N}_{\min} + 2\mathbf{N}_{\min}^{1/2} + 1} \\
&= \sum_{\mathbf{d}} \sum_{\mathbf{v}} \frac{\mathbf{p}_{\mathbf{v}\mathbf{d}} - \frac{2\mathbf{p}_{\mathbf{v}\mathbf{d}}}{\mathbf{V}} + \frac{1}{\mathbf{V}^2}}{\mathbf{N}_{\min} + 2\mathbf{N}_{\min}^{1/2} + 1} \\
&= \frac{\mathbf{D}(1 - \frac{1}{\mathbf{V}})}{\mathbf{N}_{\min} + 2\mathbf{N}_{\min}^{1/2} + 1}
\end{aligned}$$

The first equality holds because $n_{\mathbf{v}\mathbf{d}}$ is a binomial distribution with parameter $\mathbf{N}_{\mathbf{d}}$ and $\mathbf{p}_{\mathbf{v}\mathbf{d}}$. The third equality holds since the $\sum_{\mathbf{v}} \mathbf{p}_{\mathbf{v}\mathbf{d}} = 1$.

Therefore, by Lemma 6 with probability at least $1 - \epsilon/2$

$$\|\mathcal{R}_{\mathbf{P}}^{-1}(\mathbf{P} - \hat{\mathbf{P}})\|_{\text{F}} \leq \frac{1}{\mathbf{p}_{\mathbf{v}\min}} \sqrt{\frac{2\mathbf{D}(1 - \frac{1}{\mathbf{V}})}{(\mathbf{N}_{\min} + 2\mathbf{N}_{\min}^{1/2} + 1)\epsilon}}.$$

Moreover, since by assumption, $\mathbf{p}_{\mathbf{v}\min}/\mathbf{D} \geq \gamma/\mathbf{V}$, we have that

$$\|\mathcal{R}_{\mathbf{P}}^{-1}(\mathbf{P} - \hat{\mathbf{P}})\|_{\text{F}} \leq \sqrt{\frac{2\mathbf{V}^2(1 - \frac{1}{\mathbf{V}})}{\gamma^2\mathbf{D}(\mathbf{N}_{\min} + 2\mathbf{N}_{\min}^{1/2} + 1)\epsilon}}.$$

Note that

$$\sum_{\mathbf{v}} \mathbb{E} \left[\sum_{\mathbf{d}} (\hat{\mathbf{p}}_{\mathbf{v}\mathbf{d}} - \mathbf{p}_{\mathbf{v}\mathbf{d}})^2 \right] = \sum_{\mathbf{v}} \sum_{\mathbf{d}} \mathbb{E}(\hat{\mathbf{p}}_{\mathbf{v}\mathbf{d}} - \mathbf{p}_{\mathbf{v}\mathbf{d}})^2 + \sum_{\mathbf{v}} \sum_{\mathbf{d} \neq \mathbf{d}'} \mathbb{E}(\hat{\mathbf{p}}_{\mathbf{v}\mathbf{d}} - \mathbf{p}_{\mathbf{v}\mathbf{d}}) \mathbb{E}(\hat{\mathbf{p}}_{\mathbf{v}\mathbf{d}'} - \mathbf{p}_{\mathbf{v}\mathbf{d}'}).$$

We use the bound for the first term again and for the second term, we know

$$\mathbb{E}(\hat{\mathbf{p}}_{\mathbf{v}\mathbf{d}} - \mathbf{p}_{\mathbf{v}\mathbf{d}}) = \frac{\frac{1}{\mathbf{V}} - \mathbf{p}_{\mathbf{v}\mathbf{d}}}{\sqrt{\mathbf{N}_{\mathbf{d}} + 1}}.$$

So

$$\begin{aligned}
\sum_{\mathbf{v}} \sum_{\mathbf{d} \neq \mathbf{d}'} \mathbb{E}(\hat{\mathbf{p}}_{\mathbf{v}\mathbf{d}} - \mathbf{p}_{\mathbf{v}\mathbf{d}}) \mathbb{E}(\hat{\mathbf{p}}_{\mathbf{v}\mathbf{d}'} - \mathbf{p}_{\mathbf{v}\mathbf{d}'}) &= \sum_{\mathbf{v}} \sum_{\mathbf{d} \neq \mathbf{d}'} \frac{1}{(\sqrt{\mathbf{N}_{\mathbf{d}} + 1})^2} \left(\frac{1 - \mathbf{V}(\mathbf{p}_{\mathbf{v}\mathbf{d}} + \mathbf{p}_{\mathbf{v}\mathbf{d}'})}{\mathbf{V}^2} + \mathbf{p}_{\mathbf{v}\mathbf{d}}\mathbf{p}_{\mathbf{v}\mathbf{d}'} \right) \\
&= \sum_{\mathbf{d} \neq \mathbf{d}'} \frac{1}{(\sqrt{\mathbf{N}_{\mathbf{d}} + 1})^2} \sum_{\mathbf{v}} \left(\frac{1 - \mathbf{V}(\mathbf{p}_{\mathbf{v}\mathbf{d}} + \mathbf{p}_{\mathbf{v}\mathbf{d}'})}{\mathbf{V}^2} + \mathbf{p}_{\mathbf{v}\mathbf{d}}\mathbf{p}_{\mathbf{v}\mathbf{d}'} \right) \\
&= \sum_{\mathbf{d} \neq \mathbf{d}'} \frac{1}{(\sqrt{\mathbf{N}_{\mathbf{d}} + 1})^2} \left(\sum_{\mathbf{v}} \mathbf{p}_{\mathbf{v}\mathbf{d}}\mathbf{p}_{\mathbf{v}\mathbf{d}'} - \frac{1}{\mathbf{V}} \right)
\end{aligned}$$

$$\begin{aligned}
&\leq \sum_{d \neq d'} \frac{1}{(\sqrt{N_d} + 1)^2} \left(1 - \frac{1}{V}\right) \\
&\leq \frac{D^2 \left(1 - \frac{1}{V}\right)}{N_{\min} + 2N_{\min}^{1/2} + 1}
\end{aligned}$$

The third equality holds since the $\sum_v p_{vd} = 1$. The first inequality comes from the fact that

$$\max \sum_v p_{vd} p_{vd'} \quad \text{s.t.} \quad \sum_v p_{vj} = 1 \quad \text{and} \quad p_{vj} \geq 0 \quad \text{for } j = d \text{ or } d'$$

equals to 1 by Kuhn-Tucker conditions. Therefore,

$$\sum_v \mathbb{E} \left[\sum_d (\hat{p}_{vd} - p_{vd})^2 \right] \leq \frac{D(D+1) \left(1 - \frac{1}{V}\right)}{N_{\min} + 2N_{\min}^{1/2} + 1}.$$

Lemma 7 implies that with probability at least $1 - \epsilon/2$

$$\begin{aligned}
\|(\mathcal{R}_{\hat{p}}^{-1} - \mathcal{R}_p^{-1})\hat{P}\|_F &\leq \frac{1}{p_{v\min}} \sqrt{\frac{2 \sum_v \mathbb{E}(p_v - \hat{p}_v)^2}{\epsilon}} \\
&\leq \frac{1}{p_{v\min}} \sqrt{2 \frac{D(D+1) \left(1 - \frac{1}{V}\right)}{N_{\min} + 2N_{\min}^{1/2} + 1}} \\
&\leq \sqrt{\frac{2(D+1)V^2 \left(1 - \frac{1}{V}\right)}{\gamma^2 D \left(N_{\min} + 2N_{\min}^{1/2} + 1\right) \epsilon}},
\end{aligned}$$

Finally, Lemma 5, implies that if \hat{P}^{row} is based on the row-normalization of the minimax estimator then

$$\|\hat{P}^{\text{row}} - (AW)^{\text{row}}\|_F \leq \sqrt{\frac{8 \left(1 - \frac{1}{V}\right)}{\gamma^2 \cdot \epsilon} \cdot \frac{V^2}{N_{\min} + 2N_{\min}^{1/2} + 1}}$$

with probability at least $1 - \epsilon$. □

B.8 Topic estimation of FOMC1 corpus using Arora, Ge, Kannan & Moitra (2012), Ke & Wang (2022) and LDA



(a) Topic 1: wage



(b) Topic 2: recoveri



(c) Topic 3: kind



(d) Topic 4: uncertainti

Figure 12: Arora, Ge & Moitra (2012)'s estimator of A in the FOMC1 corpus. Each panel shows the word cloud of words of a topic (column in A matrix), where the font size is proportional to term's weight in the topic, and the top 5 terms with largest weights are colored. The estimated anchor words for each topic is in the caption.



(a) Topic 1



(b) Topic 2



(c) Topic 3



(d) Topic 4

Figure 13: Ke & Wang (2022)’s estimator of A in the FOMC1 corpus. Each panel shows the word cloud of words of a topic (column in A matrix), where the font size is proportional to term’s weight in the topic, and the top 5 terms with largest weights are colored. The estimated anchor words for each topic is in the caption.



(a) Topic 1



(b) Topic 2



(c) Topic 3



(d) Topic 4

Figure 14: Latent Dirichlet Allocation estimator of A in the FOMC1 corpus with uniform priors. Each panel shows the word cloud of words of a topic (column in A matrix), where the font size is proportional to term's weight in the topic, and the top 5 terms with largest weights are colored. The estimated anchor words for each topic is in the caption.