MAT 167: Applied Linear Algebra Lecture 23: Text Mining II

Naoki Saito

Department of Mathematics University of California, Davis

November 19 & 21, 2025

Outline

Clustering

2 Nonnegative Matrix Factorization

Outline

Clustering

2 Nonnegative Matrix Factorization

- Instead of using the left singular vectors as a basis to approximate a term-document matrix, let's examine the cluster centers (centroids) obtained by k-means algorithm as a basis.
- Let $C_k = [c_1 \dots c_k] \in \mathbb{R}^{m \times k}$ be the k cluster centroids obtained by the k-means algorithm.
- c_j 's are non-orthogonal; hence it is more convenient to obtain a set of orthonormal vectors that spans range(C_k).
- To do so, we can use the reduced QR factorization: $C_k = \widehat{Q}_k \widehat{R}_k$ where $\widehat{Q}_k \in \mathbb{R}^{m \times k}$, and $\widehat{R}_k \in \mathbb{R}^{k \times k}$.
- Now, let's approximate A using \widehat{Q}_k in the sense of the least squares as:

$$\min_{G_k \in \mathbb{R}^{k \times n}} \|A - \widehat{Q}_k G_k\|_F.$$

$$\min_{\mathbf{r}_i \in \mathbb{R}^k} \|\mathbf{a}_j - \widehat{Q}_k \mathbf{g}_j\|_2, \quad j = 1 : n.$$

- Instead of using the left singular vectors as a basis to approximate a term-document matrix, let's examine the cluster centers (centroids) obtained by k-means algorithm as a basis.
- Let $C_k = [c_1 \dots c_k] \in \mathbb{R}^{m \times k}$ be the k cluster centroids obtained by the k-means algorithm.
- c_j 's are non-orthogonal; hence it is more convenient to obtain a set of orthonormal vectors that spans range(C_k).
- To do so, we can use the reduced QR factorization: $C_k = \widehat{Q}_k \widehat{R}_k$ where $\widehat{Q}_k \in \mathbb{R}^{m \times k}$, and $\widehat{R}_k \in \mathbb{R}^{k \times k}$.
- Now, let's approximate A using \widehat{Q}_k in the sense of the least squares as:

$$\min_{G_k \in \mathbb{R}^{k \times n}} \|A - \widehat{Q}_k G_k\|_F.$$

$$\min_{\boldsymbol{x}:\in\mathbb{R}^k} \|\boldsymbol{a}_j - \widehat{Q}_k \boldsymbol{g}_j\|_2, \quad j = 1:n.$$

- Instead of using the left singular vectors as a basis to approximate a term-document matrix, let's examine the cluster centers (centroids) obtained by k-means algorithm as a basis.
- Let $C_k = [c_1 \dots c_k] \in \mathbb{R}^{m \times k}$ be the k cluster centroids obtained by the k-means algorithm.
- c_j 's are non-orthogonal; hence it is more convenient to obtain a set of orthonormal vectors that spans range(C_k).
- To do so, we can use the reduced QR factorization: $C_k = \hat{Q}_k \hat{R}_k$ where $\hat{Q}_k \in \mathbb{R}^{m \times k}$, and $\hat{R}_k \in \mathbb{R}^{k \times k}$.
- Now, let's approximate A using \hat{Q}_k in the sense of the least squares as

$$\min_{G_k \in \mathbb{R}^{k \times n}} \|A - \widehat{Q}_k G_k\|_F.$$

$$\min_{\mathbf{g}_i \in \mathbb{R}^k} \|\mathbf{a}_j - \widehat{Q}_k \mathbf{g}_j\|_2, \quad j = 1 : n.$$

- Instead of using the left singular vectors as a basis to approximate a term-document matrix, let's examine the cluster centers (centroids) obtained by k-means algorithm as a basis.
- Let $C_k = [c_1 \dots c_k] \in \mathbb{R}^{m \times k}$ be the k cluster centroids obtained by the k-means algorithm.
- c_j 's are non-orthogonal; hence it is more convenient to obtain a set of orthonormal vectors that spans range(C_k).
- To do so, we can use the reduced QR factorization: $C_k = \widehat{Q}_k \widehat{R}_k$ where $\widehat{Q}_k \in \mathbb{R}^{m \times k}$, and $\widehat{R}_k \in \mathbb{R}^{k \times k}$.
- ullet Now, let's approximate A using \widehat{Q}_k in the sense of the least squares as

$$\min_{G_k \in \mathbb{R}^{k \times n}} \|A - \widehat{Q}_k G_k\|_F.$$

$$\min_{\mathbf{g}_{i} \in \mathbb{R}^{k}} \|\mathbf{a}_{j} - \widehat{Q}_{k}\mathbf{g}_{j}\|_{2}, \quad j = 1 : n.$$

- Instead of using the left singular vectors as a basis to approximate a term-document matrix, let's examine the cluster centers (centroids) obtained by k-means algorithm as a basis.
- Let $C_k = [c_1 \dots c_k] \in \mathbb{R}^{m \times k}$ be the k cluster centroids obtained by the k-means algorithm.
- c_j 's are non-orthogonal; hence it is more convenient to obtain a set of orthonormal vectors that spans range(C_k).
- To do so, we can use the reduced QR factorization: $C_k = \widehat{Q}_k \widehat{R}_k$ where $\widehat{Q}_k \in \mathbb{R}^{m \times k}$, and $\widehat{R}_k \in \mathbb{R}^{k \times k}$.
- ullet Now, let's approximate A using \widehat{Q}_k in the sense of the least squares as:

$$\min_{G_k \in \mathbb{R}^{k \times n}} \|A - \widehat{Q}_k G_k\|_F.$$

$$\min_{\boldsymbol{\sigma}:\in\mathbb{R}^k}\|\boldsymbol{a}_j-\widehat{Q}_k\boldsymbol{g}_j\|_2, \quad j=1:n.$$

- Instead of using the left singular vectors as a basis to approximate a term-document matrix, let's examine the cluster centers (centroids) obtained by k-means algorithm as a basis.
- Let $C_k = [c_1 \dots c_k] \in \mathbb{R}^{m \times k}$ be the k cluster centroids obtained by the k-means algorithm.
- c_j 's are non-orthogonal; hence it is more convenient to obtain a set of orthonormal vectors that spans range(C_k).
- To do so, we can use the reduced QR factorization: $C_k = \widehat{Q}_k \widehat{R}_k$ where $\widehat{Q}_k \in \mathbb{R}^{m \times k}$, and $\widehat{R}_k \in \mathbb{R}^{k \times k}$.
- ullet Now, let's approximate A using \widehat{Q}_k in the sense of the least squares as:

$$\min_{G_k \in \mathbb{R}^{k \times n}} \|A - \widehat{Q}_k G_k\|_F.$$

$$\min_{\boldsymbol{g}_j \in \mathbb{R}^k} \|\boldsymbol{a}_j - \widehat{Q}_k \boldsymbol{g}_j\|_2, \quad j = 1 : n.$$

- Since the columns of \widehat{Q}_k are orthonormal, we can get the following LS solution: $\mathbf{g}_j = \widehat{Q}_k^{\mathsf{T}} \mathbf{a}_j$, j = 1 : n. Hence $G_k = \widehat{Q}_k^{\mathsf{T}} A$.
- The inner product between the query vector q and the document vector a_i can be approximated as:

$$\mathbf{q}^{\mathsf{T}} \mathbf{a}_{j} \approx \mathbf{q}^{\mathsf{T}} \widehat{Q}_{k} \mathbf{g}_{j} = (\widehat{Q}_{k}^{\mathsf{T}} \mathbf{q})^{\mathsf{T}} \mathbf{g}_{j} = \mathbf{q}_{k}^{\mathsf{T}} \mathbf{g}_{j}, \ \mathbf{q}_{k} := \widehat{Q}_{k}^{\mathsf{T}} \mathbf{q}$$

Hence, the cosine similarity can be approximated as:

$$\frac{\mathbf{q}^{\mathsf{T}} \mathbf{a}_{j}}{\|\mathbf{q}\|_{2} \|\mathbf{a}_{j}\|_{2}} \approx \frac{\mathbf{q}_{k}^{\mathsf{T}} \mathbf{g}_{j}}{\|\mathbf{q}\|_{2} \|\mathbf{g}_{j}\|_{2}}$$

- Since the columns of \widehat{Q}_k are orthonormal, we can get the following LS solution: $\mathbf{g}_i = \widehat{Q}_k^{\mathsf{T}} \mathbf{a}_i$, j = 1 : n. Hence $G_k = \widehat{Q}_k^{\mathsf{T}} A$.
- The inner product between the query vector \mathbf{q} and the document vector \mathbf{a}_j can be approximated as:

$$\boldsymbol{q}^{\mathsf{T}}\boldsymbol{a}_{j}\approx\boldsymbol{q}^{\mathsf{T}}\widehat{Q}_{k}\boldsymbol{g}_{j}=(\widehat{Q}_{k}^{\mathsf{T}}\boldsymbol{q})^{\mathsf{T}}\boldsymbol{g}_{j}=\boldsymbol{q}_{k}^{\mathsf{T}}\boldsymbol{g}_{j},\ \boldsymbol{q}_{k}:=\widehat{Q}_{k}^{\mathsf{T}}\boldsymbol{q}.$$

Hence, the cosine similarity can be approximated as

$$\frac{\mathbf{q}^{\mathsf{T}} \mathbf{a}_{j}}{\|\mathbf{q}\|_{2} \|\mathbf{a}_{j}\|_{2}} \approx \frac{\mathbf{q}_{k}^{\mathsf{T}} \mathbf{g}_{j}}{\|\mathbf{q}\|_{2} \|\mathbf{g}_{j}\|_{2}}$$

- Since the columns of \widehat{Q}_k are orthonormal, we can get the following LS solution: $\mathbf{g}_j = \widehat{Q}_k^\mathsf{T} \mathbf{a}_j$, j = 1 : n. Hence $G_k = \widehat{Q}_k^\mathsf{T} A$.
- The inner product between the query vector \mathbf{q} and the document vector \mathbf{a}_i can be approximated as:

$$\boldsymbol{q}^{\mathsf{T}}\boldsymbol{a}_{j}\approx\boldsymbol{q}^{\mathsf{T}}\widehat{Q}_{k}\boldsymbol{g}_{j}=(\widehat{Q}_{k}^{\mathsf{T}}\boldsymbol{q})^{\mathsf{T}}\boldsymbol{g}_{j}=\boldsymbol{q}_{k}^{\mathsf{T}}\boldsymbol{g}_{j},\;\boldsymbol{q}_{k}:=\widehat{Q}_{k}^{\mathsf{T}}\boldsymbol{q}.$$

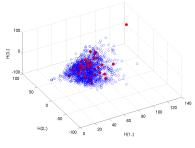
• Hence, the cosine similarity can be approximated as:

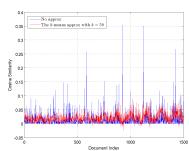
$$\frac{\boldsymbol{q}^{\mathsf{T}}\boldsymbol{a}_{j}}{\|\boldsymbol{q}\|_{2}\|\boldsymbol{a}_{j}\|_{2}} \approx \frac{\boldsymbol{q}_{k}^{\mathsf{T}}\boldsymbol{g}_{j}}{\|\boldsymbol{q}\|_{2}\|\boldsymbol{g}_{j}\|_{2}}.$$

- k = 50; the same query vector ('entropy', 'minimum', 'maximum').
- The approximation error between $\widehat{Q}_k G_k$ and A was $\|A \widehat{Q}_k G_k\|_F / \|A\|_F \approx 0.7227$, which was worse than that using the top 100 SVD basis.

- k = 50; the same query vector ('entropy', 'minimum', 'maximum').
- The approximation error between $\widehat{Q}_k G_k$ and A was $\|A \widehat{Q}_k G_k\|_F / \|A\|_F \approx 0.7227$, which was worse than that using the top 100 SVD basis.

- k = 50; the same query vector ('entropy', 'minimum', 'maximum').
- The approximation error between $\widehat{Q}_k G_k$ and A was $\|A \widehat{Q}_k G_k\|_F / \|A\|_F \approx 0.7227$, which was worse than that using the top 100 SVD basis.





Documents in U_{100} [:,1:3]

Cosine Similarity

With the 50-means based approximation, tol=0.2, 0.1, 0.05 correspond to 0, 0, 81 returned documents; Compare these with the no approximation case: 4, 15, 89; and with the best 100 approximation using SVD: 0, 4, 72.

- ullet Running the k-means algorithm with large m and n is slow in general.
- If your document set really consists of k different topics (or categories), then this k-means-based approach should work well.
 Example: The Science News Dataset consisting of articles in the area of Anthropology, Astronomy, Behavioral Sciences, Earth Sciences, Life Sciences, Math & CS, Medicine, Physics. Which value of k should be used is still a question though.
- However, in the case of the NIPS data where there is not much clustering structure, it may not worth trying this approach considering the computational cost.

- Running the k-means algorithm with large m and n is slow in general.
- If your document set really consists of k different topics (or categories), then this k-means-based approach should work well.
 Example: The Science News Dataset consisting of articles in the area of Anthropology, Astronomy, Behavioral Sciences, Earth Sciences, Life Sciences, Math & CS, Medicine, Physics. Which value of k should be used is still a question though.
- However, in the case of the NIPS data where there is not much clustering structure, it may not worth trying this approach considering the computational cost.

- Running the k-means algorithm with large m and n is slow in general.
- If your document set really consists of k different topics (or categories), then this k-means-based approach should work well.
 Example: The Science News Dataset consisting of articles in the area of Anthropology, Astronomy, Behavioral Sciences, Earth Sciences, Life Sciences, Math & CS, Medicine, Physics. Which value of k should be used is still a question though.
- However, in the case of the NIPS data where there is not much clustering structure, it may not worth trying this approach considering the computational cost.

Outline

Clustering

Nonnegative Matrix Factorization

- Consider the NNMF of a term-document matrix $A \approx WH$ where $W \in \mathbb{R}^{m \times k}$, $H \in \mathbb{R}^{k \times n}$, $1 < k \le \min(m, n)$.
- We want to represent (or approximate) both query vectors and the term-document matrix using the basis vectors $\{\mathbf{w}_1, ..., \mathbf{w}_k\}$ and do query task in that basis (or coordinates).
- a_j is already approximated using $\{w_1, ..., w_k\}$ with the coordinate vector h_j , j = 1 : n, i.e., $a_j \approx W h_j$.
- We need to approximate \boldsymbol{q} in the basis of W. To do so, we seek the LS approximation of \boldsymbol{q} in range(W), i.e., $\min_{\widehat{\boldsymbol{q}} \in \mathbb{R}^k} \|\boldsymbol{q} W\widehat{\boldsymbol{q}}\|_2$.
- Hence we need to solve the normal equation: $W^{\mathsf{T}}W\widehat{\boldsymbol{q}}=W^{\mathsf{T}}\boldsymbol{q}$.
- To do so, we use the reduced QR factorization of $W = \widehat{Q}\widehat{R}$.
- Then, using the argument of Lecture 10, the normal equation above is equivalent to $\widehat{R}\widehat{q} = \widehat{Q}^T q$, i.e., $\widehat{q} = \widehat{R}^{-1}\widehat{Q}^T q$.
- The cosine similarity in the basis of $\{w_1, ..., w_k\}$ can be written as:

$$\frac{\widehat{\mathbf{q}}^{\mathsf{T}} \mathbf{h}_j}{\|\widehat{\mathbf{q}}\|_2 \|\mathbf{h}_i\|_2}, \quad j = 1 : n.$$

- Consider the NNMF of a term-document matrix $A \approx WH$ where $W \in \mathbb{R}^{m \times k}$, $H \in \mathbb{R}^{k \times n}$, $1 < k \le \min(m, n)$.
- We want to represent (or approximate) both query vectors and the term-document matrix using the basis vectors $\{\boldsymbol{w}_1, ..., \boldsymbol{w}_k\}$ and do query task in that basis (or coordinates).
- a_j is already approximated using $\{w_1, ..., w_k\}$ with the coordinate vector h_j , j = 1 : n, i.e., $a_j \approx W h_j$.
- We need to approximate \boldsymbol{q} in the basis of W. To do so, we seek the LS approximation of \boldsymbol{q} in range(W), i.e., $\min_{\widehat{\boldsymbol{q}} \in \mathbb{R}^k} \|\boldsymbol{q} W\widehat{\boldsymbol{q}}\|_2$.
- Hence we need to solve the normal equation: $W^{\mathsf{T}}W\widehat{\boldsymbol{q}}=W^{\mathsf{T}}\boldsymbol{q}$.
- To do so, we use the reduced QR factorization of $W = \widehat{Q}\widehat{R}$
- Then, using the argument of Lecture 10, the normal equation above is equivalent to $\widehat{R}\widehat{\boldsymbol{q}} = \widehat{Q}^{\mathsf{T}}\boldsymbol{q}$, i.e., $\widehat{\boldsymbol{q}} = \widehat{R}^{-1}\widehat{Q}^{\mathsf{T}}\boldsymbol{q}$.
- The cosine similarity in the basis of $\{w_1, ..., w_k\}$ can be written as:

$$\frac{\widehat{\boldsymbol{q}}^{\mathsf{T}}\boldsymbol{h}_{j}}{\|\widehat{\boldsymbol{q}}\|_{2}\|\boldsymbol{h}_{i}\|_{2}}, \quad j=1:n.$$

- Consider the NNMF of a term-document matrix $A \approx WH$ where $W \in \mathbb{R}^{m \times k}$, $H \in \mathbb{R}^{k \times n}$, $1 < k \le \min(m, n)$.
- We want to represent (or approximate) both query vectors and the term-document matrix using the basis vectors $\{\boldsymbol{w}_1, ..., \boldsymbol{w}_k\}$ and do query task in that basis (or coordinates).
- a_j is already approximated using $\{w_1, ..., w_k\}$ with the coordinate vector h_j , j = 1 : n, i.e., $a_j \approx W h_j$.
- We need to approximate \boldsymbol{q} in the basis of W. To do so, we seek the LS approximation of \boldsymbol{q} in range(W), i.e., $\min_{\widehat{\boldsymbol{q}} \in \mathbb{R}^k} \|\boldsymbol{q} W\widehat{\boldsymbol{q}}\|_2$.
- Hence we need to solve the normal equation: $W^{\mathsf{T}}W\widehat{\boldsymbol{q}}=W^{\mathsf{T}}\boldsymbol{q}$.
- To do so, we use the reduced QR factorization of $W = \widehat{Q}\widehat{R}$
- Then, using the argument of Lecture 10, the normal equation above is equivalent to $\widehat{R}\widehat{\boldsymbol{q}} = \widehat{Q}^{\mathsf{T}}\boldsymbol{q}$, i.e., $\widehat{\boldsymbol{q}} = \widehat{R}^{-1}\widehat{Q}^{\mathsf{T}}\boldsymbol{q}$.
- The cosine similarity in the basis of $\{w_1, ..., w_k\}$ can be written as:

$$\frac{\widehat{\boldsymbol{q}}^{\mathsf{T}}\boldsymbol{h}_{j}}{\|\widehat{\boldsymbol{q}}\|_{2}\|\boldsymbol{h}_{i}\|_{2}}, \quad j=1:n.$$

- Consider the NNMF of a term-document matrix $A \approx WH$ where $W \in \mathbb{R}^{m \times k}$, $H \in \mathbb{R}^{k \times n}$, $1 < k \le \min(m, n)$.
- We want to represent (or approximate) both query vectors and the term-document matrix using the basis vectors $\{\boldsymbol{w}_1, ..., \boldsymbol{w}_k\}$ and do query task in that basis (or coordinates).
- a_j is already approximated using $\{w_1, ..., w_k\}$ with the coordinate vector h_j , j = 1 : n, i.e., $a_j \approx W h_j$.
- We need to approximate \boldsymbol{q} in the basis of W. To do so, we seek the LS approximation of \boldsymbol{q} in range(W), i.e., $\min_{\widehat{\boldsymbol{q}} \in \mathbb{R}^k} \|\boldsymbol{q} W\widehat{\boldsymbol{q}}\|_2$.
- Hence we need to solve the normal equation: $W^{\mathsf{T}}W\widehat{\boldsymbol{q}}=W^{\mathsf{T}}\boldsymbol{q}$.
- To do so, we use the reduced QR factorization of $W = \widehat{Q}\widehat{R}$
- Then, using the argument of Lecture 10, the normal equation above is equivalent to $\widehat{R}\widehat{\boldsymbol{q}} = \widehat{Q}^{\mathsf{T}}\boldsymbol{q}$, i.e., $\widehat{\boldsymbol{q}} = \widehat{R}^{-1}\widehat{Q}^{\mathsf{T}}\boldsymbol{q}$.
- The cosine similarity in the basis of $\{w_1, ..., w_k\}$ can be written as:

$$\frac{\widehat{\boldsymbol{q}}^{\mathsf{T}}\boldsymbol{h}_{j}}{\|\widehat{\boldsymbol{q}}\|_{2}\|\boldsymbol{h}_{i}\|_{2}}, \quad j=1:n.$$

- Consider the NNMF of a term-document matrix $A \approx WH$ where $W \in \mathbb{R}^{m \times k}$, $H \in \mathbb{R}^{k \times n}$, $1 < k \le \min(m, n)$.
- We want to represent (or approximate) both query vectors and the term-document matrix using the basis vectors $\{w_1, ..., w_k\}$ and do query task in that basis (or coordinates).
- a_j is already approximated using $\{w_1, ..., w_k\}$ with the coordinate vector h_j , j = 1 : n, i.e., $a_j \approx W h_j$.
- We need to approximate \boldsymbol{q} in the basis of W. To do so, we seek the LS approximation of \boldsymbol{q} in range(W), i.e., $\min_{\widehat{\boldsymbol{q}} \in \mathbb{R}^k} \|\boldsymbol{q} W\widehat{\boldsymbol{q}}\|_2$.
- Hence we need to solve the normal equation: $W^T W \hat{q} = W^T q$.
- To do so, we use the reduced QR factorization of $W = \widehat{Q}\widehat{R}$
- Then, using the argument of Lecture 10, the normal equation above is equivalent to $\widehat{R}\widehat{q} = \widehat{Q}^T q$, i.e., $\widehat{q} = \widehat{R}^{-1}\widehat{Q}^T q$.
- The cosine similarity in the basis of $\{w_1, ..., w_k\}$ can be written as

$$\frac{\widehat{\boldsymbol{q}}^{\mathsf{T}}\boldsymbol{h}_{j}}{\|\widehat{\boldsymbol{q}}\|_{2}\|\boldsymbol{h}_{i}\|_{2}}, \quad j=1:n.$$

- Consider the NNMF of a term-document matrix $A \approx WH$ where $W \in \mathbb{R}^{m \times k}$, $H \in \mathbb{R}^{k \times n}$, $1 < k \le \min(m, n)$.
- We want to represent (or approximate) both query vectors and the term-document matrix using the basis vectors $\{\boldsymbol{w}_1, ..., \boldsymbol{w}_k\}$ and do query task in that basis (or coordinates).
- a_j is already approximated using $\{w_1, ..., w_k\}$ with the coordinate vector h_j , j = 1 : n, i.e., $a_j \approx W h_j$.
- We need to approximate \boldsymbol{q} in the basis of W. To do so, we seek the LS approximation of \boldsymbol{q} in range(W), i.e., $\min_{\widehat{\boldsymbol{q}} \in \mathbb{R}^k} \|\boldsymbol{q} W\widehat{\boldsymbol{q}}\|_2$.
- Hence we need to solve the normal equation: $W^T W \hat{q} = W^T q$.
- To do so, we use the reduced QR factorization of $W = \widehat{Q}\widehat{R}$.
- Then, using the argument of Lecture 10, the normal equation above is equivalent to $\widehat{R}\widehat{\boldsymbol{q}} = \widehat{Q}^{\mathsf{T}}\boldsymbol{q}$, i.e., $\widehat{\boldsymbol{q}} = \widehat{R}^{-1}\widehat{Q}^{\mathsf{T}}\boldsymbol{q}$.
- The cosine similarity in the basis of $\{w_1, ..., w_k\}$ can be written as

$$\frac{\widehat{\boldsymbol{q}}^{\mathsf{T}}\boldsymbol{h}_{j}}{\|\widehat{\boldsymbol{q}}\|_{2}\|\boldsymbol{h}_{i}\|_{2}}, \quad j=1:n.$$

- Consider the NNMF of a term-document matrix $A \approx WH$ where $W \in \mathbb{R}^{m \times k}$, $H \in \mathbb{R}^{k \times n}$, $1 < k \le \min(m, n)$.
- We want to represent (or approximate) both query vectors and the term-document matrix using the basis vectors $\{\boldsymbol{w}_1, ..., \boldsymbol{w}_k\}$ and do query task in that basis (or coordinates).
- a_j is already approximated using $\{w_1, ..., w_k\}$ with the coordinate vector h_j , j = 1 : n, i.e., $a_j \approx W h_j$.
- We need to approximate \boldsymbol{q} in the basis of W. To do so, we seek the LS approximation of \boldsymbol{q} in range(W), i.e., $\min_{\widehat{\boldsymbol{q}} \in \mathbb{R}^k} \|\boldsymbol{q} W\widehat{\boldsymbol{q}}\|_2$.
- Hence we need to solve the normal equation: $W^T W \hat{q} = W^T q$.
- To do so, we use the reduced QR factorization of $W = \widehat{Q}\widehat{R}$.
- Then, using the argument of Lecture 10, the normal equation above is equivalent to $\widehat{R}\widehat{q} = \widehat{Q}^T q$, i.e., $\widehat{q} = \widehat{R}^{-1}\widehat{Q}^T q$.
- The cosine similarity in the basis of $\{w_1, ..., w_k\}$ can be written as:

$$\frac{\widehat{\mathbf{q}}^{\mathsf{T}} \mathbf{h}_{j}}{\|\widehat{\mathbf{q}}\|_{2} \|\mathbf{h}_{i}\|_{2}}, \quad j = 1 : n$$

- Consider the NNMF of a term-document matrix $A \approx WH$ where $W \in \mathbb{R}^{m \times k}$, $H \in \mathbb{R}^{k \times n}$, $1 < k \le \min(m, n)$.
- We want to represent (or approximate) both query vectors and the term-document matrix using the basis vectors $\{\boldsymbol{w}_1, ..., \boldsymbol{w}_k\}$ and do query task in that basis (or coordinates).
- a_j is already approximated using $\{w_1, ..., w_k\}$ with the coordinate vector h_j , j = 1 : n, i.e., $a_j \approx W h_j$.
- We need to approximate \boldsymbol{q} in the basis of W. To do so, we seek the LS approximation of \boldsymbol{q} in range(W), i.e., $\min_{\widehat{\boldsymbol{q}} \in \mathbb{R}^k} \|\boldsymbol{q} W\widehat{\boldsymbol{q}}\|_2$.
- Hence we need to solve the normal equation: $W^T W \hat{q} = W^T q$.
- To do so, we use the reduced QR factorization of $W = \widehat{Q}\widehat{R}$.
- Then, using the argument of Lecture 10, the normal equation above is equivalent to $\widehat{R}\widehat{q} = \widehat{Q}^T q$, i.e., $\widehat{q} = \widehat{R}^{-1}\widehat{Q}^T q$.
- The cosine similarity in the basis of $\{w_1, ..., w_k\}$ can be written as:

$$\frac{\widehat{\boldsymbol{q}}^{\mathsf{T}}\boldsymbol{h}_{j}}{\|\widehat{\boldsymbol{q}}\|_{2}\|\boldsymbol{h}_{j}\|_{2}}, \quad j=1:n.$$

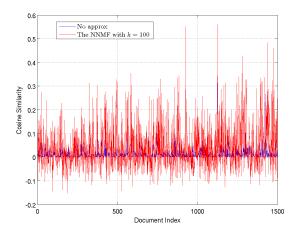
- k = 100 was used.
- $||A WH||_F / ||A||_F \approx 0.6302$, which was *slightly* worse than that using the top 100 SVD basis (0.6074).
- Each w_j concentrates on one term, and is close to the canonical vector $e_i \in \mathbb{R}^m$ for some i (recall: NNMF applied to the face database in Lecture 20).
- The peaks of w_j , j=1:10, correspond to: 'network', 'model', 'learning', 'function', 'unit', 'algorithm', 'input', 'data', 'neuron', 'cell', which are quite similar to the u_1 vector or the 10 most frequently used terms.
- On the other hand, because w_j 's are localized, the interpretation of the row vectors of H matrix becomes easy.

- k = 100 was used.
- $||A WH||_F / ||A||_F \approx 0.6302$, which was *slightly* worse than that using the top 100 SVD basis (0.6074).
- Each w_j concentrates on one term, and is close to the canonical vector $e_i \in \mathbb{R}^m$ for some i (recall: NNMF applied to the face database in Lecture 20).
- The peaks of \mathbf{w}_j , j=1:10, correspond to: 'network', 'model', 'learning', 'function', 'unit', 'algorithm', 'input', 'data', 'neuron', 'cell', which are quite similar to the \mathbf{u}_1 vector or the 10 most frequently used terms.
- On the other hand, because w_j 's are localized, the interpretation of the row vectors of H matrix becomes easy.

- k = 100 was used.
- $||A WH||_F / ||A||_F \approx 0.6302$, which was *slightly* worse than that using the top 100 SVD basis (0.6074).
- Each \mathbf{w}_i concentrates on one term, and is close to the canonical vector $\mathbf{e}_i \in \mathbb{R}^m$ for some i (recall: NNMF applied to the face database in Lecture 20).
- The peaks of \mathbf{w}_j , j=1:10, correspond to: 'network', 'model', 'learning', 'function', 'unit', 'algorithm', 'input', 'data', 'neuron', 'cell', which are quite similar to the \mathbf{u}_1 vector or the 10 most frequently used terms.
- On the other hand, because w_j 's are localized, the interpretation of the row vectors of H matrix becomes easy.

- k = 100 was used.
- $||A WH||_F / ||A||_F \approx 0.6302$, which was *slightly* worse than that using the top 100 SVD basis (0.6074).
- Each w_j concentrates on one term, and is close to the canonical vector $e_i \in \mathbb{R}^m$ for some i (recall: NNMF applied to the face database in Lecture 20).
- The peaks of \mathbf{w}_j , j=1:10, correspond to: 'network', 'model', 'learning', 'function', 'unit', 'algorithm', 'input', 'data', 'neuron', 'cell', which are quite similar to the \mathbf{u}_1 vector or the 10 most frequently used terms.
- On the other hand, because w_j 's are localized, the interpretation of the row vectors of H matrix becomes easy.

- k = 100 was used.
- $||A WH||_F / ||A||_F \approx 0.6302$, which was *slightly* worse than that using the top 100 SVD basis (0.6074).
- Each \mathbf{w}_i concentrates on one term, and is close to the canonical vector $\mathbf{e}_i \in \mathbb{R}^m$ for some i (recall: NNMF applied to the face database in Lecture 20).
- The peaks of \mathbf{w}_j , j=1:10, correspond to: 'network', 'model', 'learning', 'function', 'unit', 'algorithm', 'input', 'data', 'neuron', 'cell', which are quite similar to the \mathbf{u}_1 vector or the 10 most frequently used terms.
- On the other hand, because \mathbf{w}_j 's are localized, the interpretation of the row vectors of H matrix becomes easy.



With the NNMF-based approach using k = 100, tol=0.2, 0.1, 0.05 correspond to 101, 312, 535 returned documents; Compare with the no approximation case: 4, 15, 89. Changing the tol=0.4, 0.3, 0.2 with NNMF returns 5, 26, 101 documents.

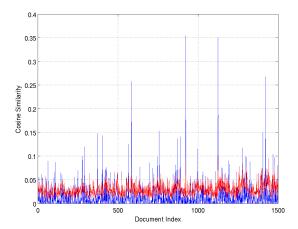
- Using the LS solution for the query saves computational cost given the NNMF is already obtained because one can avoid the explicit computation and storage of WH.
- If we can compute and store WH, then we could use the following approximation of the original cosine similarity:

$$\frac{\mathbf{q}^{\mathsf{T}} \mathbf{a}_{j}}{\|\mathbf{q}\|_{2} \|\mathbf{a}_{j}\|_{2}} \approx \frac{\mathbf{q}^{\mathsf{T}} W \mathbf{h}_{j}}{\|\mathbf{q}\|_{2} \|W \mathbf{h}_{j}\|_{2}}$$

- Using the LS solution for the query saves computational cost given the NNMF is already obtained because one can avoid the explicit computation and storage of WH.
- If we can compute and store WH, then we could use the following approximation of the original cosine similarity:

$$\frac{\boldsymbol{q}^{\mathsf{T}}\boldsymbol{a}_{j}}{\|\boldsymbol{q}\|_{2}\|\boldsymbol{a}_{j}\|_{2}} \approx \frac{\boldsymbol{q}^{\mathsf{T}}W\boldsymbol{h}_{j}}{\|\boldsymbol{q}\|_{2}\|W\boldsymbol{h}_{j}\|_{2}}.$$

My Reaction ...



With the NNMF-based approach using k = 100 using the above cosine similarity approximation, tol=0.2, 0.1, 0.05 correspond to 0, 1, 97 returned documents; Compare with the no approximation case: 4, 15, 89. Without using the LS query, some of the relevant documents do not stick out clearly.