

Stat models for Big Data Topic models and NMF

Purnamrita Sarkar Department of Statistics and Data Science The University of Texas at Austin

https://psarkar.github.io/teaching

- So, SVD returns directions or principal components
- But these are not interpretable.
- But what if we optimized the following?

$$\min_{\substack{U \in \mathbb{R}_{m \times k}^+ \\ V \in \mathbb{R}_{n \times k}^+}} \|A - UV^T\|_F^2$$

- So, SVD returns directions or principal components
- But these are not interpretable.
- But what if we optimized the following?

$$\min_{\substack{U \in \mathbb{R}^+_{m \times k} \\ V \in \mathbb{R}^+_{n \times k}}} \|A - UV^T\|_F^2$$

• Is this factorization unique?

- So, SVD returns directions or principal components
- But these are not interpretable.
- But what if we optimized the following?

$$\min_{\substack{U \in \mathbb{R}^+_{m \times k} \\ V \in \mathbb{R}^+_{n \times k}}} \|A - UV^T\|_F^2$$

- Is this factorization unique?
- No I could multiply U by a positive constant, and divide V by the same and that will give me the same UV^T

The non-negative matrix factorization angle

- Typically, the issues with uniqueness can be resolved by putting constraints on norm or sparsity.
- Despite that, we now have a non-convex loss. There a variety of algorithms, most of them based on alternating minimization type methods.

- Typically, the issues with uniqueness can be resolved by putting constraints on norm or sparsity.
- Despite that, we now have a non-convex loss. There a variety of algorithms, most of them based on alternating minimization type methods.
- Here is the loss function minimized by the buit-in NMF code in scikit-learn

$$\min_{\substack{\boldsymbol{U} \in \mathbb{R}^+_{m \times k} \\ \boldsymbol{V} \in \mathbb{R}^+_{n \times k}}} \|\boldsymbol{A} - \boldsymbol{U}\boldsymbol{V}^{\mathsf{T}}\|_F^2 + \alpha\beta \left(\|\operatorname{vec}(\boldsymbol{W})\|_1 + \|\operatorname{vec}(\boldsymbol{H})\|_1\right) \\ + \frac{1}{2}\alpha(1-\beta)\left(\|\boldsymbol{W}\|_F^2 + \|\boldsymbol{H}\|_F^2\right)$$

• α, β are regularization parameters

Why Non-negative matrix factorization

- Let us compare the basis vectors obtained using NMF and matrix factorization.
- Look at the right singular vectors or the V in the aforementioned optimization problem with k = 20.



- Take 1 minute to think how the two are different.
- Drumrolls

- The basis vectors from SVD are global, they are picking up a linear combination of the individual pixel values (which are the features)
- On the other hand, NMF is actually picking up the different parts of the threes, which can be thought of as pieces which are combined together in different ways to give many different handwritten 3's.

- The basis vectors from SVD are global, they are picking up a linear combination of the individual pixel values (which are the features)
- On the other hand, NMF is actually picking up the different parts of the threes, which can be thought of as pieces which are combined together in different ways to give many different handwritten 3's.
- So NMF is interpretable, and columns of *U* and *V* are not orthogonal.
- But we need conditions to make sure that algorithms return the global optima, and one needs to also think about uniqueness.

Matrix completion - NMF angle

| | Pide Pupulice Pupulice | TO KILL A Mockingbird | Right Ho, Cleaves | NOT COL | DEVIX AND CRAKE MARCARET ATWOOD | | PD James HILDREN OF NEN | MARGARET ATWOOD THIN HARGARES TALE |
|--------|------------------------------|--------------------------|----------------------|---------|---------------------------------------|---|----------------------------------|--|
| Alice | 4 | 3 | 5 | 4 | 1 | 1 | 1 | 2 |
| Bob | 4 | 5 | 4 | 5 | 1 | 2 | 2 | 1 |
| Meena | 4 | 5 | 4 | 4 | 4 | 5 | 5 | 3 |
| Asaf | 1 | 1 | 1 | 1 | 4 | 4 | 4 | 5 |
| Arthur | 2 | 1 | 1 | 1 | 5 | 4 | 4 | 4 |

- Consider a user-book rating matrix.
- We random pick 5 elements and set them to zero (think missing).

| 4 | 0 | 5 | 4 | 1 | 1 | 0 | 2 |
|---|---|---|---|---|---|---|---|
| 4 | 4 | 4 | 5 | 1 | 0 | 2 | 1 |
| 4 | 5 | 4 | 4 | 0 | 5 | 5 | 3 |
| 1 | 1 | 1 | 1 | 4 | 4 | 4 | 5 |
| 2 | 1 | 0 | 1 | 5 | 4 | 4 | 4 |

Matrix completion - NMF angle

- We will do SVD to get $Y = U_1 V_1^T$
- We will do NMF to get $Y = U_2 V_2^T$ Now we will use U_1 and U_2 to embed the users as we had before.



Table 1: (A) embedding with SVD, (B) embedding with NMF

• Take a few minutes to ponder over why these two are different and which one is more interpretable and why.

Matrix completion - NMF angle

- We will do SVD to get $Y = U_1 V_1^T$
- We will do NMF to get $Y = U_2 V_2^T$ Now we will use U_1 and U_2 to embed the users as we had before.



Table 2: (A) embedding with SVD, (B) embedding with NMF

• NMF is more interpretable, because Alice/Bob are placed on the X axis (approximately) and Arthur/Asaf on the Y axis, so its almost like the different directions are for the different genres of books, classics and dystopian fiction.



- You can take A as fixed
- *W* is stochastic and there are many models for generating documents as a mixture of topics.
- A notable such model is Latent Dirichlet Allocation, by Blei, Ng and Jordan (JMLR 2003). For a document,
 - Choose $N \sim Poisson(\xi)$
 - Choose θ ∼ Dir(α)
 - For each of the *N* words,
 - Choose topic $t \sim Multinomial(\theta)$
 - Choose word w_n from $p(w_n|z_n)$ specified by the columns of the fixed A matrix.

Prev work : It is NP-hard to compute NMF But if we make an assumption, then there is a simple polynomial time algorithm. Separability: A matriz An separable, if for every column of A, I a row of A whose only non-zero enty is in attrat column.



So, in previous example, the rows in W
appear as rows in M (upto scaling).
T Say I have the auchor words
T I know
$$W = \frac{M(2,)K_1}{M(3,:),K_2}$$

 $-\infty$ columps of W Sum to 1.





$$\overline{Q}_{i,j} = P(w_{2}=j|w_{1}=i) = \sum_{l} P(w_{2}=j, t_{1}=l|w_{1}=i) = \sum_{l} P(w_{2}=j|t_{1}=l,w_{1}=i) P(t_{1}=l|w_{1}=i) = \sum_{l} P(w_{2}=j|t_{1}=l,w_{1}=i) P(t_{1}=l|w_{1}=i) = \sum_{l} P(w_{2}=j|t_{1}=l) P(t_{1}=l|w_{1}=i) = \sum_{l} P(w_{2}=j|t_{1}=l) P(t_{1}=l|w_{1}=i) = \sum_{l} P(w_{2}=j|t_{1}=k) = \sum_{l} P(t_{1}=l|w_{1}=i) = \sum_{l} P(t_{1}=l|w_$$

If I know anchor mode set S,
how do we get
$$A_{ik} = P(w=i|t=k)$$
?
 $\overline{Q}_{ij} = \sum_{k} \frac{P(t_j=k|w_j=i)}{C_{ik}} \overline{Q}_{sk,j} \quad vx\tau \qquad \begin{array}{c} p_{i} = \Sigma Q_{i} \\ \vdots \\ Q = C \overline{Q}_{s}; 1 \\ C = \overline{Q} \overline{Q}_{s}^{T} (\overline{Q}_{s} \overline{Q}_{s}^{T})^{-1} \\ A_{ik} = \frac{P(t=k|w=i)}{\sum_{k} P(t=k|w=i)} P(w=i) = \frac{C_{ik} p_{i}}{\sum_{k} C_{ik} p_{i}} \end{array}$

17

• A Practical Algorithm for Topic Modeling with Provable Guarantees, Arora et al, ICML 2013