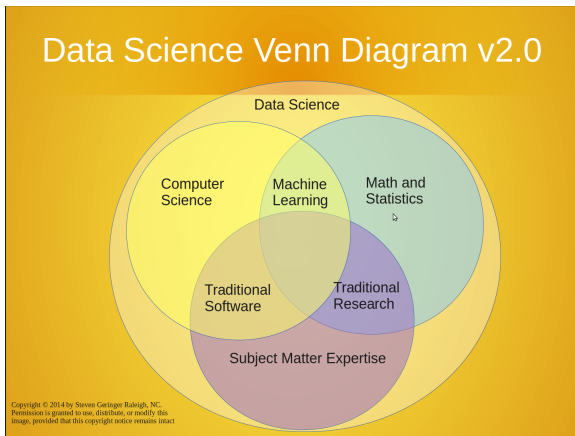


Non-Negative Matrix Factorization

Schwartz

August 14, 2017

How do you make a data scientist?



	Numerical Literacy	Scripting Coding	Inference Estimation	Predictive Modeling	Business Acumen	Creative Thinking	Problem Solving	?
--	--------------------	------------------	----------------------	---------------------	-----------------	-------------------	-----------------	---

Mathematics Degree								
Statistics Degree								
Economics Degree								
Computer Science Degree								
Data Science Immersive								
Independent Self Study								
Workshops and Lectures								
Community Engagement								
?								

Objectives

▶ NMF

- ▶ versus SVD
- ▶ non-negative
- ▶ parts-based model

▶ Uses

- ▶ learn/interpret latent reduced dimensionality features driving data
- ▶ soft cluster samples by latent features

▶ Estimating NMF

- ▶ gradient descent
- ▶ alternating least squares (ALS)
- ▶ multiplicative updating

What is NMF?

Singular Value Decomposition (SVD):

$$\begin{aligned} X_{n \times p} &= U_{n \times n} \Sigma_{n \times p} V_{p \times p}^T \\ &\approx U_{n \times k} \Sigma_{k \times k} V_{k \times p}^T \end{aligned}$$

What is NMF?

Singular Value Decomposition (SVD):

$$\begin{aligned} X_{n \times p} &= U_{n \times n} \Sigma_{n \times p} V_{p \times p}^T \\ &\approx U_{n \times k} \Sigma_{k \times k} V_{k \times p}^T \end{aligned}$$

Non-Negative Matrix Factorization (NMF):

$$\begin{aligned} X_{n \times p} &\approx W_{n \times k} H_{k \times p} \\ X_{ij}, W_{i'j'}, H_{i^*j^*} &\geq 0 \end{aligned}$$

What is NMF?

Singular Value Decomposition (SVD):

$$\begin{aligned} X_{n \times p} &= U_{n \times n} \Sigma_{n \times p} V_{p \times p}^T \\ &\approx U_{n \times k} \Sigma_{k \times k} V_{k \times p}^T \end{aligned}$$

Non-Negative Matrix Factorization (NMF):

$$\begin{aligned} X_{n \times p} &\approx W_{n \times k} H_{k \times p} \\ X_{ij}, W_{i'j'}, H_{i^*j^*} &\geq 0 \end{aligned}$$

So NMF is just SVD

– just drop the middle matrix and keep all the numbers positive

$$\geq 0$$

Keep all the numbers positive why?

Scores	item 1	item 2	item p
user 1				
user 2				
\vdots				
user n				

$$X_{n \times p}$$

$$= W_{n \times k} H_{k \times p}$$

$$\geq 0$$

Keep all the numbers positive why?

Scores	item 1	item 2	item p
user 1				
user 2				
\vdots				
user n				

$$= \begin{matrix} W \\ n \times k \end{matrix} \begin{matrix} H \\ k \times p \end{matrix}$$

$$\begin{matrix} X \\ n \times p \end{matrix}$$

What are some recommender systems you know about, and what kind of numbers are used in those ratings, typically?

$$\geq 0$$

Keep all the numbers positive why? Because we can

Scores	item 1	item 2	item p
user 1				
user 2				
\vdots				
user n				

$$= \begin{matrix} W \\ n \times k \end{matrix} \begin{matrix} H \\ k \times p \end{matrix}$$

$$\begin{matrix} X \\ n \times p \end{matrix}$$

What are some recommender systems you know about, and what kind of numbers are used in those ratings, typically?

$$\geq 0$$

Keep all the numbers positive why? Because we can

Scores	item 1	item 2	item p
user 1				
user 2				
\vdots				
user n				

$$= \begin{matrix} W \\ n \times k \end{matrix} \begin{matrix} H \\ k \times p \end{matrix}$$

$$X_{n \times p}$$

What are some recommender systems you know about, and what kind of numbers are used in those ratings, typically?

What about $W_{n \times k}$ and $H_{k \times p}$?

≥ 0

Keep all the numbers positive why? Because we can

Scores	item 1	item 2	item p
user 1				
user 2				
\vdots				
user n				

$$= \begin{matrix} W \\ n \times k \end{matrix} \begin{matrix} H \\ k \times p \end{matrix}$$

$$X_{n \times p}$$

What are some recommender systems you know about, and what kind of numbers are used in those ratings, typically?

What about $W_{n \times k}$ and $H_{k \times p}$?

“parts-based model”

≥ 0 : what is the NMF “parts based model”?

$$\underset{n \times p}{X} \approx \underset{n \times k}{W} \underset{k \times p}{H} \quad \hat{X}_{ij} = \underset{1 \times k}{W_{i \cdot}} \underset{k \times 1}{H_{\cdot j}}$$

$$X_{ij}, W_{i'j'}, H_{i^*j^*} \geq 0$$

| protein carbs fat fiber vitamins | (individuals diet implications)

(individuals food intake)

= | food 1 food 2 ... food k | ×

	protein	carbs	fat	fiber	vitamins
food 1	(food contributions)				
food 2					
⋮					
⋮					
food k					

≥ 0 : what is the NMF “parts based model”?

$$\underset{n \times p}{X} \approx \underset{n \times k}{W} \underset{k \times p}{H} \quad \hat{X}_{ij} = \underset{1 \times k}{W_i} \underset{k \times 1}{H_j}$$

$$X_{ij}, W_{i'j'}, H_{i^*j^*} \geq 0$$

protein carbs fat fiber vitamins					(individuals diet implications)
(individuals food intake)					
=	food 1	food 2	...	food k	×
					protein carbs fat fiber vitamins
					(food contributions)
					food 1
					food 2
					⋮
					food k

- ▶ every user i gets “amounts” of k factors (food): W 's i^{th} row
- ▶ factors k may contribute to item j (feature): H 's j^{th} column
- ▶ “agreement” in factors for user i and item j determines X_{ij}

≥ 0 : what is the NMF “parts based model”?

$$\underset{n \times p}{X} \approx \underset{n \times k}{W} \underset{k \times p}{H} \quad \hat{X}_{ij} = \underset{1 \times k}{W_{i \cdot}} \underset{k \times 1}{H_{\cdot j}}$$

$$X_{ij}, W_{i'j'}, H_{i^*j^*} \geq 0$$

protein carbs fat fiber vitamins					(individuals diet implications)
(individuals food intake)					
=	food 1	food 2	...	food k	×
					protein carbs fat fiber vitamins
					food 1
					food 2
					⋮
					⋮
					food k
					(food contributions)

- ▶ every user i gets “amounts” of k factors (food): W 's i^{th} row
- ▶ factors k may contribute to item j (feature): H 's j^{th} column
- ▶ “agreement” in factors for user i and item j determines X_{ij}
- ▶ everything being positive provides this “parts-based model”

Example: NMF Topic Modeling for NLP

	word 1	word 2	word p
doc 1				
doc 2				
⋮				
doc n				

$$\begin{array}{c|cccc} & \text{topic 1} & \dots & \text{topic } k & \\ \hline \text{doc 1} & & & & \\ \text{doc 2} & & & & \\ \vdots & & & & \\ \text{doc n} & & & & \end{array} = \begin{array}{c} \times \\ \text{topic 1} \\ \vdots \\ \text{topic } k \end{array} \begin{array}{c|cccc} & \text{word 1} & \dots & \text{word } p & \\ \hline & & & & \end{array}$$

- ▶ Identifies latent “topics” or features driving word appearance
- ▶ Says what words each of the topics are comprised of (cool!)
- ▶ Gives topic similarity (how?) for documents (soft clustering)

How good is this model?

$$\underset{n \times p}{X} \approx \underset{n \times k}{W} \underset{k \times p}{H}$$

How good is this model?

$$\underset{n \times p}{X} \approx \underset{n \times k}{W} \underset{k \times p}{H}$$

$$\sum_{i,j} (X_{ij} - \hat{X}_{ij})^2 = \sum_{i,j} (X_{ij} - W_{i \cdot} H_{\cdot j})^2$$

How good is this model?

$$\underset{n \times p}{X} \approx \underset{n \times k}{W} \underset{k \times p}{H}$$

$$\sum_{i,j} (X_{ij} - \hat{X}_{ij})^2 = \sum_{i,j} (X_{ij} - W_{i \cdot} H_{\cdot j})^2$$

- What i and j do we use to evaluate this on?

Scores	item 1	item 2	item p
user 1				
user 2				
\vdots				
user n				

How good is this model?

$$\underset{n \times p}{X} \approx \underset{n \times k}{W} \underset{k \times p}{H}$$

$$\sum_{i,j} (X_{ij} - \hat{X}_{ij})^2 = \sum_{i,j} (X_{ij} - W_i \cdot H_{.j})^2$$

- What i and j do we use to evaluate this on?

Scores	item 1	item 2	item p
user 1				
user 2				
\vdots				
user n				

- What's interesting about this compared to SVD?

How good is this model?

$$\underset{n \times p}{X} \approx \underset{n \times k}{W} \underset{k \times p}{H}$$

$$\sum_{i,j} (X_{ij} - \hat{X}_{ij})^2 = \sum_{i,j} (X_{ij} - W_i \cdot H_j)^2$$

- What i and j do we use to evaluate this on?

Scores	item 1	item 2	item p
user 1				
user 2				
\vdots				
user n				

- What's interesting about this compared to SVD?

Hint: what does SVD do with NA values?

How do we learn H and W ?

How do we learn H and W ?

► Gradient Descent

$$\frac{\partial}{\partial W_{i'k}} \sum_{i,j} (X_{ij} - W_{i'k} H_{kj})^2 = \sum_j -2(X_{i'j} - W_{i'k} H_{kj}) H_{kj}$$

$$\frac{\partial}{\partial H_{kj'}} \sum_{i,j} (X_{ij} - W_{i'k} H_{kj})^2 = \sum_i -2(X_{ij'} - W_{ik} H_{kj'}) W_{ik}$$

subject to $W_{ij}, H_{i'j'} \geq 0$

How do we learn H and W ?

► Gradient Descent

$$\frac{\partial}{\partial W_{i'k}} \sum_{i,j} (X_{ij} - W_{i'k} H_{kj})^2 = \sum_j -2(X_{i'j} - W_{i'k} H_{kj}) H_{kj}$$

$$\frac{\partial}{\partial H_{kj'}} \sum_{i,j} (X_{ij} - W_{i'k} H_{kj})^2 = \sum_i -2(X_{ij'} - W_{ik} H_{kj'}) W_{ik}$$

subject to $W_{ij}, H_{i'j'} \geq 0$

► Alternating Least Squares (ALS)

1. Initialize H and W with $H_{ij}, W_{i'j'} > 0$
1. Update $H_{.j}$ using OLS: $X_{.j} = \mathbf{W} H_{.j} + \epsilon_{.j}$

$$\begin{bmatrix} | \\ | \\ | \end{bmatrix} = \begin{bmatrix} \rightarrow \rightarrow \\ \rightarrow \rightarrow \\ \rightarrow \rightarrow \end{bmatrix} \begin{bmatrix} \downarrow \\ \downarrow \\ \downarrow \end{bmatrix}$$

2. Update $W_{i.}$ using OLS: $X_{i.}^T = \mathbf{H}^T W_{i.}^T + \epsilon_{i.}^T$

$$\begin{bmatrix} - & - \end{bmatrix}^T = \left(\begin{bmatrix} \rightarrow \rightarrow \end{bmatrix} \begin{bmatrix} \downarrow \downarrow \downarrow \downarrow \\ \downarrow \downarrow \downarrow \downarrow \\ \downarrow \downarrow \downarrow \downarrow \\ \downarrow \downarrow \downarrow \downarrow \end{bmatrix} \right)^T$$

3. If $H_{ij} < 0$ set $H_{ij} = 0$; if $W_{i'j'} < 0$ set $W_{i'j'} = 0$
4. Evaluate stopping criterion: return to step 1 if check fails
(What stopping criterion might we use?)

How do we learn H and W ?

► Lee and Seung's “multiplicative update rules”

0. Initialize H and W with $H_{ij}, W_{i'j'} > 0$

1. Update W and H with

$$W'_{ik} = W_{ik} \frac{(XH^T)_{ik}}{(WHH^T)_{ik}} \quad H'_{kj} = H_{kj} \frac{(W^T X)_{kj}}{(W^T WH)_{kj}}$$

2. Evaluate stopping criterion: return to step 1 if check fails
(What stopping criterion might we use?)

How do we learn H and W ?

► Lee and Seung's “multiplicative update rules”

0. Initialize H and W with $H_{ij}, W_{i'j'} > 0$

1. Update W and H with

$$W'_{ik} = W_{ik} \frac{(XH^T)_{ik}}{(WHH^T)_{ik}} \quad H'_{kj} = H_{kj} \frac{(W^T X)_{kj}}{(W^T WH)_{kj}}$$

2. Evaluate stopping criterion: return to step 1 if check fails
(What stopping criterion might we use?)

What is doing this?

► Recall the OLS estimate $\hat{\beta} = (X^T X)^{-1} X^T Y$

Which looks a lot like, e.g., $\hat{H}_{kj} = \frac{(W^T X)_{kj}}{(W^T W)_{kj}}$

► So the update looks a lot like $H'_{kj} = H_{kj} \frac{\hat{H}_{kj}}{H}$

so if the new estimate is increased/decreased relative to the current value then the estimate is increase/decrease by that proportion; but this change in H will result in a change in \hat{W} the next time, which will again change \hat{H} the next next time...

Wrap Up

Both do Unsupervised Dimensionality Reduction, but

	SVD/PCA	NMF
NA's	Nope	Yep
Estimation	Non-iterative	Iterative
Coefficients	Orthogonal, $+/-$ coefficients	only $+$ coefficients
Interpretation	Linear combination	"Parts-based"
k	Skree plot All of those as well \rightarrow	Cross-validation Permutation testing Cophenetic correlation Interpretability

Go look at all the other instructors
jupyter notebooks demoing this all