

Non-Negative Matrix Factorization (NMF)

(and Interpretable Topic Modeling)

Overview

- What is NMF?
 - Unbaking a cake
- Popular Applications
 - Document Clustering (i.e., Topic Modeling)
- Deeper Dive into Mechanics
 - Gradient Descent
 - Alternating Least Squares (ALS)
- NMF vs. PCA

Problem Motivation

- Soft clustering
 - Clustering where each observation can have partial membership in each cluster
- Identify latent features
 - Identify topics in a corpus of text
 - Identity parts of faces in facial recognition

Example: Unbaking a cake

(i.e., finding its recipe)

- Given a cake, we cannot directly observe its ingredients (flour (F) and sugar (S))
 - (the recipe is a closely guarded secret!)

recipe

a. k. a., $= W = (w_F \quad w_S)$

cake as **ingredients**

Example: Unbaking a cake

(i.e., finding its recipe)

- But we **can** directly measure its **macro nutrients (carbohydrates (*c*), proteins (*p*), fats (*f*))**
 - E.g., the cake is made of 50g of carbohydrates, 10g of proteins, and 3g of fats

cake as macro nutrients = $(50 \quad 10 \quad 3)$

$$V = (v_c \quad v_p \quad v_f)$$

Example: Unbaking a cake

(i.e., finding its recipe)

- Ingredients are themselves made of macro nutrients (for which we know the proportions):
 - 1 cup flour is made of 20g of carbohydrates, 5g of proteins, and 1g of fats
 - 1 cup sugar is made of 10g of carbohydrates, 0g of proteins, and 1g of fats

$$\text{ingredients as macro nutrients} = \begin{pmatrix} 20 & 5 & 1 \\ 10 & 0 & 1 \end{pmatrix}$$

$$H = \begin{pmatrix} h_{F,c} & h_{F,p} & h_{F,f} \\ h_{S,c} & h_{S,p} & h_{S,f} \end{pmatrix}$$

Example: Unbaking a cake (i.e., finding its recipe)

- Figuring out the recipe is about solving the following equation:

$$\begin{array}{ccccc} \text{cake} & & \text{ingredients} & & \text{cake} \\ \text{as} & & & & \text{as} \\ \text{ingredients} & \cdot & \text{as} & \approx & \text{as} \\ \text{(recipe)} & & \text{macro nutrients} & & \text{macro nutrients} \end{array}$$

$$W \cdot H \approx V$$

$$(w_F \quad w_S) \cdot \begin{pmatrix} 20 & 5 & 1 \\ 10 & 0 & 1 \end{pmatrix} \approx (50 \quad 10 \quad 3)$$

Generalization: Unbaking cakes (i.e., finding recipes)

$$\begin{array}{c}
 \text{cake} \mathbf{s} \\
 \text{as} \\
 \text{ingredients} \cdot \\
 \text{(recipe)}
 \end{array}
 \begin{array}{c}
 \text{ingredients} \\
 \text{as} \\
 \text{macro nutrients}
 \end{array}
 \approx
 \begin{array}{c}
 \text{cake} \mathbf{s} \\
 \text{as} \\
 \text{macro nutrients}
 \end{array}$$

$$W \cdot H \approx V$$

$$\begin{pmatrix}
 W_{\mathbf{1},F} & W_{\mathbf{1},S} \\
 W_{\mathbf{2},F} & W_{\mathbf{2},S} \\
 W_{\mathbf{3},F} & W_{\mathbf{3},S} \\
 \vdots & \vdots
 \end{pmatrix} \cdot \begin{pmatrix} 20 & 5 & 1 \\ 10 & 0 & 1 \end{pmatrix} \approx \begin{pmatrix} 50 & 10 & 3 \\ 30 & 5 & 2 \\ 25 & 3 & 3 \\ \vdots & \vdots & \vdots \end{pmatrix}$$

Generalization: Unbaking cakes

(i.e., finding recipes)

- Given H and V , can we estimate W ?

$$\begin{pmatrix} w_{1,F} & w_{1,S} \\ w_{2,F} & w_{2,S} \\ w_{3,F} & w_{3,S} \\ \vdots & \vdots \end{pmatrix} \cdot \begin{pmatrix} 20 & 5 & 1 \\ 10 & 0 & 1 \end{pmatrix} \approx \begin{pmatrix} 50 & 10 & 3 \\ 30 & 5 & 2 \\ 25 & 3 & 3 \\ \vdots & \vdots & \vdots \end{pmatrix}$$

- Can perform OLS minimization to find each recipe. E.g,

$$(w_{1,F} \quad w_{1,S}) \cdot \begin{pmatrix} 20 & 5 & 1 \\ 10 & 0 & 1 \end{pmatrix} \approx (50 \quad 10 \quad 3)$$

Generalization: Unbaking cakes

(i.e., finding recipes)

$$(w_{1,F} \quad w_{1,S}) \cdot \begin{pmatrix} 20 & 5 & 1 \\ 10 & 0 & 1 \end{pmatrix} \approx (50 \quad 10 \quad 3)$$

$$20w_{1,F} + 10w_{1,S} = \hat{v}_c \approx v_c = 50$$

$$5w_{1,F} + 0w_{1,S} = \hat{v}_p \approx v_p = 10$$

$$1w_{1,F} + 1w_{1,S} = \hat{v}_f \approx v_f = 3$$

$$\operatorname{argmin}_{w=(w_{1,F} \quad w_{1,S})} \|v - \hat{v}\|^2$$

Unbaking cakes and finding ingredients

- Given W and V , can we estimate H ?

$$\begin{pmatrix} 2 & 1 \\ 1 & 1 \\ 1 & 2 \\ \vdots & \vdots \end{pmatrix} \cdot \begin{pmatrix} h_{F,c} & h_{F,p} & h_{F,f} \\ h_{S,c} & h_{S,p} & h_{S,f} \end{pmatrix} \approx \begin{pmatrix} 50 & 10 & 3 \\ 30 & 5 & 2 \\ 25 & 3 & 3 \\ \vdots & \vdots & \vdots \end{pmatrix}$$

$$2h_{F,c} + 1h_{F,c} = \hat{v}_{1,c} \approx v_{1,c} = 50$$

$$1w_{1,F} + 1w_{1,S} = \hat{v}_{1,p} \approx v_{1,p} = 10$$

$$1w_{1,F} + 2w_{1,S} = \hat{v}_{1,f} \approx v_{1,f} = 3$$

(...continued for all observations)

$$\underset{H=\begin{pmatrix} h_{F,c} & h_{F,p} & h_{F,f} \\ h_{S,c} & h_{S,p} & h_{S,f} \end{pmatrix}}{\operatorname{argmin}} \|V - \hat{V}\|^2$$

Unbaking cakes and finding ingredients

- What is you don't know W and H ?

$$\begin{pmatrix} w_{1,F} & w_{1,S} \\ w_{2,F} & w_{2,S} \\ w_{3,F} & w_{3,S} \\ \vdots & \vdots \end{pmatrix} \cdot \begin{pmatrix} h_{F,c} & h_{F,p} & h_{F,f} \\ h_{S,c} & h_{S,p} & h_{S,f} \end{pmatrix} \approx \begin{pmatrix} 50 & 10 & 3 \\ 30 & 5 & 2 \\ 25 & 3 & 3 \\ \vdots & \vdots & \vdots \end{pmatrix}$$

- NMF can learn what the ingredients and the proportions are from all these cakes

What is Non-Negative Matrix Factorization (NMF)?

Matrix $V_{m \times n}$ where each entry $v_{i,j} \geq 0$

$$\begin{array}{ccccc} W & \times & H & = & V \\ m \times \mathbf{r} & & \mathbf{r} \times n & & m \times n \end{array}$$

(also, $w_{i,j} \geq 0$ and $h_{i,j} \geq 0$)

(non-negativity constraint, unlike the decompositions we've looked at so far)

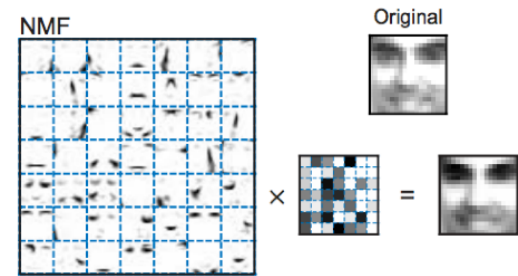
What is Non-Negative Matrix Factorization (NMF)?

$$\begin{array}{ccccc} W & \times & H & \approx & V \\ m \times r & & r \times n & & m \times n \\ \begin{array}{c} W \\ \left[\begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\ \hline \square & \square \\ \hline \square & \square \\ \hline \end{array} \right] & \times & \begin{array}{c} H \\ \left[\begin{array}{|c|c|c|c|c|c|} \hline \square & \square & \square & \square & \square & \square \\ \hline \square & \square & \square & \square & \square & \square \\ \hline \end{array} \right] & \approx & \begin{array}{c} V \\ \left[\begin{array}{|c|c|c|c|c|c|} \hline \square & \square & \square & \square & \square & \square \\ \hline \square & \square & \square & \square & \square & \square \\ \hline \square & \square & \square & \square & \square & \square \\ \hline \square & \square & \square & \square & \square & \square \\ \hline \end{array} \right] \end{array} \end{array}$$

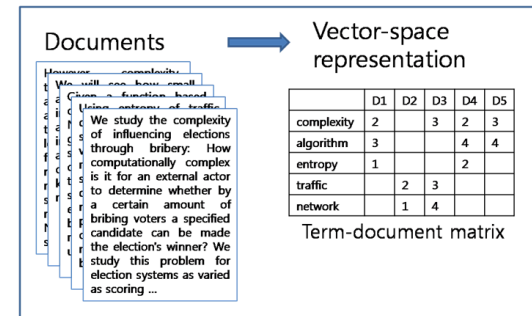
- r is set by you ($r < \min(m, n)$)
- Cannot be solved analytically, so approximated numerically
- When $r > \min(m, n)$, we can perfectly recreate V as $W \cdot H$
- Otherwise, some sort of data compression is happening

Popular Applications

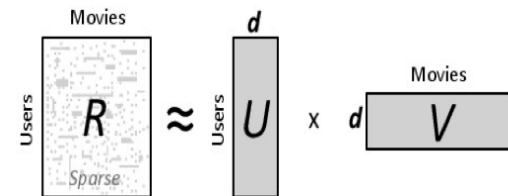
- Computer Visioning
 - Identify / classifying objects
 - Generally reducing feature space of images



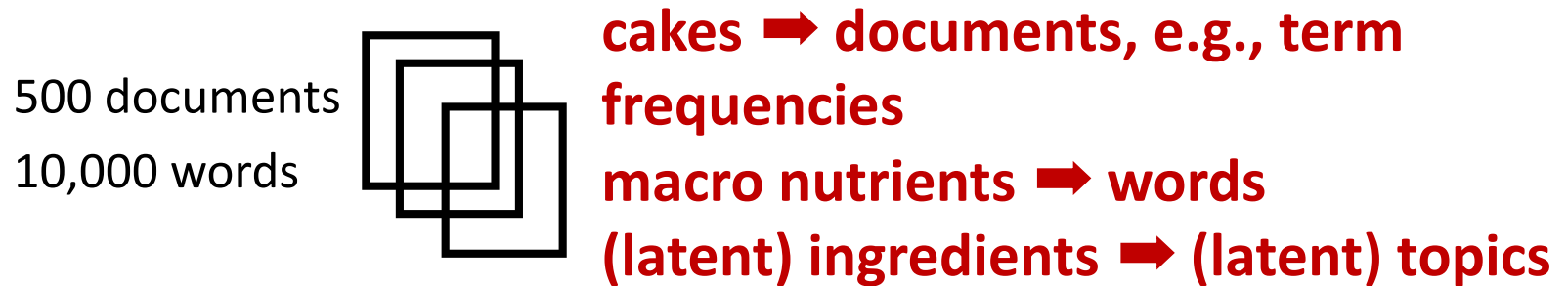
- Topic Modeling
 - A.k.a., Document Clustering
 - Today!



- Recommender systems
 - In just a few days!



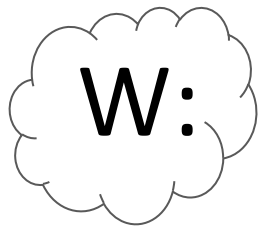
Topic Modeling with NMF



$$V \approx W \cdot H$$

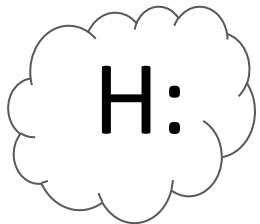
cake as **ingredients** as **cake**
as \cdot as \approx as
ingredients macro nutrients macro nutrients

document as **topics** as **document**
as \cdot as \approx as
topics words words



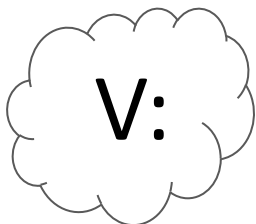
documents

latent topics



latent topics

words



documents

words



- In looking at the dimensions of W , we notice that the number of rows remain the same, as we would expect. Thus, each row of W must represent some information about the corresponding row in V .
- Think of the column W as latent topics where the higher the word's cell value, the higher the word's rank for that latent feature.
- The same can be said about H : the number of columns between V and H are the same. Thus, each column of H must represent some information about the corresponding column in V .

Choosing k ?

- Choosing k is more of an art than a science. We can look at how "good" of an approximation WH is for V and try to find the smallest k that makes it suitably small.
- At the end of the day, though, k , is likely going to be chosen based on intuition that you derive from inspecting the topics and possibly from some domain knowledge.

Example H Matrix

	"president"	"coach"	...	"team"	
$H =$ $k \times m$	0.3	5.1	...	4.2	Topic 1
	10.3	1.07	...	0.08	Topic 2
	⋮
	
	2.03	0.3	...	0.001	Topic k

Topic 1?

Topic 2?

Example H Matrix

	"president"	"coach"	...	"team"	
$H =$ $k \times m$	0.3	5.1	...	4.2	Topic 1
	10.3	1.07	...	0.08	Topic 2
	⋮
	
	2.03	0.3	...	0.001	Topic k

Topic 1: Sports?

Topic 2: Politics?

Mechanics

Minimize $\|V - WH\|^2$

with respect to W and H subject to $W, H \geq 0$

Steps

(1) Start with some random W and H (random, small positive values)

(2) Repeatedly adjust W and H to make RMSE smaller

$$H_{a\mu} \leftarrow H_{a\mu} \frac{(W^T V)_{a\mu}}{(W^T W H)_{a\mu}} \quad W_{ia} \leftarrow W_{ia} \frac{(V H^T)_{ia}}{(W H H^T)_{ia}}$$

($H_{a\mu}$ and W_{ia} are single entries in W and H)

- This is gradient descent! It's simple but can be slow. The convergence is sensitive to choice of step size.

(3) Stop when some threshold is met

- Decrease in RMSE, # of iterations, etc.

Mechanics – Alternating Least Squares (ALS)

$$\text{Minimize } \|V - WH\|^2$$

with respect to W and H subject to $W, H \geq 0$

Steps

(1) Randomly initialize^(*) W and H to the **appropriate shapes**

^(*) Initialization methods vary. Random initialization and k-means are common. One way: Initialize W to small, positive, random values.

(2) Repeat following

- Holding W fixed, update H by minimizing sum of squared errors. **Ensure all $H \geq 0$. (Clip negative values)**
- Holding H fixed, update W by minimizing sum of squared errors. **Ensure all $W \geq 0$. (Clip negative values)**

(3) Stop until convergence (i.e., some threshold is met)

- Decrease in RMSE, # of iterations, etc.

Note: ALS is “biconvex”, i.e., it is not globally convex, so it will not necessarily find local optima (hence the strategic initializations)

NMF vs. PCA

PCA

Unsupervised dimensionality reduction

Orthogonal vectors with positive and negative coefficients

“Holistic”; difficult to interpret

Non-iterative

NMF

Unsupervised dimensionality reduction

Non-negative coefficients

“Parts-based”; easier to interpret

Iterative (the presented algorithm)