



SVD Application: Categorizing Movie Data

Laura Lyman Margot Gerritsen

Stanford University

ICME Summer Workshops
Fundamentals of Data Science

August 2, 2021

Simple example with movie data

Let A be a rectangular data matrix, whose

- rows correspond to users,
- columns correspond to movies, and
- entries are ratings that users have given movies.

Specifically, $A_{ij} = \text{rating user } i \text{ has given movie } j$. For example,

$$A = \begin{bmatrix} & \text{The Matrix} & \text{Alien} & \text{Serenity} & \text{Casablanca} & \text{Amelie} \\ & 1 & 1 & 1 & 0 & 0 \\ & 2 & 2 & 2 & 0 & 0 \\ & 1 & 1 & 1 & 0 & 0 \\ & 5 & 5 & 5 & 0 & 0 \\ & 0 & 0 & 0 & 2 & 2 \\ & 0 & 0 & 0 & 3 & 3 \\ & 0 & 0 & 0 & 1 & 1 \end{bmatrix}.$$

Simple example with movie data

Let A be a rectangular data matrix, whose

- rows correspond to users,
- columns correspond to movies, and
- entries are ratings that users have given movies.

Specifically, $A_{ij} = \text{rating user } i \text{ has given movie } j$. For example,

User 3 rated "Alien" 1 out of 5

$$A = \begin{bmatrix} & \text{The Matrix} & \text{Alien} & \text{Serenity} & \text{Casablanca} & \text{Amelie} \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 & 0 \\ 0 & 0 & 0 & 3 & 3 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 \end{bmatrix}$$

User 5 saw Casablanca and Amelie and rated them 2/5

Simple example with movie data

Let A be a rectangular data matrix, whose

- rows correspond to users,
- columns correspond to movies, and
- entries are ratings that users have given movies.

Specifically, $A_{ij} = \text{rating user } i \text{ has given movie } j$. For example,

$A_{ij} = 0 \iff$
user i has not
rated movie j

$$A = \left[\begin{array}{ccccc} \text{The Matrix} & \text{Alien} & \text{Serenity} & \text{Casablanca} & \text{Amelie} \\ 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{array} \right].$$

e.g. first 4 users
didn't rate
Casablanca or
Amelie

last 3 users {

Simple example with movie data

Let A be a rectangular data matrix, whose

- rows correspond to users,
- columns correspond to movies, and
- entries are ratings that users have given movies.

Specifically, $A_{ij} =$ rating user i has given movie j . For example,

$$A = \left[\begin{array}{ccc|cc} & \text{The Matrix} & \text{Alien} & \text{Serenity} & \text{Casablanca} & \text{Amelie} \\ \text{Sci Fi category } \left\{ & 1 & 1 & 1 & 0 & 0 \\ & 2 & 2 & 2 & 0 & 0 \\ & 1 & 1 & 1 & 0 & 0 \\ & 5 & 5 & 5 & 0 & 0 \\ \hline & 0 & 0 & 0 & 2 & 2 \\ & 0 & 0 & 0 & 3 & 3 \\ & 0 & 0 & 0 & 1 & 1 \end{array} \right] . \quad \text{ROMANCE } \heartsuit \text{ category } \right\}$$

SVD for movie data

Using the compact SVD,

$$\left[\begin{array}{ccc|cc} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ \hline 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{array} \right] = \underbrace{\begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix}}_U \times \underbrace{\begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix}}_{\Sigma} \times \underbrace{\begin{bmatrix} 0.58 & 0 \\ 0.58 & 0 \\ 0.58 & 0 \\ 0 & 0.71 \\ 0 & 0.71 \end{bmatrix}}_{V^T}.$$

Romance SCI FI

- Σ : # of categories and their relative strength in the data set

SVD for movie data

Using the compact SVD,

1	1	1	0	0	
2	2	2	0	0	
1	1	1	0	0	
5	5	5	0	0	
0	0	0	2	2	
0	0	0	3	3	
0	0	0	1	1	

Romance SCI FI

if user 3
rates
higher

\hat{u}_1	\hat{u}_2	\hat{u}_3	v_1	v_2
4	4	4	0	0
5	5	5	0	0
0	0	0	2	2
0	0	0	3	3
0	0	0	1	1

$$\Sigma = \begin{bmatrix} \dots & \dots \\ 0 & 5.29 \end{bmatrix} v'$$

- Σ : # of categories and their relative strength in the data set

SVD for movie data

Using the compact SVD,

1	1	1	0	0
2	2	2	0	0
1	1	1	0	0
<u>5</u>	<u>5</u>	<u>5</u>	0	0
0	0	0	<u>2</u>	<u>2</u>
0	0	0	<u>3</u>	<u>3</u>
0	0	0	1	1

Romance SCI FI

if user 3
rates
higher

$\hat{4}$	$\hat{4}$	$\hat{4}$	v	v
5	5	5	0	0
0	0	0	2	2
0	0	0	3	3
0	0	0	1	1

$$\Sigma = \begin{bmatrix} \text{---} & & & & \\ 0 & & & & 5.29 \\ \text{---} & & & & \\ & & & & v' \end{bmatrix}$$

increased

$$\Sigma = \begin{bmatrix} & & & & \\ 0 & & & & 4.47 \\ & & & & \\ & & & & \downarrow \text{decreased} \end{bmatrix}$$

if user 5
changes to
Sci Fi

1	1	1	0	0
5	5	5	0	0
<u>2</u>	<u>2</u>	<u>2</u>	0	0
0	0	0	3	3
0	0	0	1	1

- Σ : # of categories and their relative strength in the data set

SVD for movie data

Using the compact SVD,

$$\left[\begin{array}{ccc|cc} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ \hline 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{array} \right] = \underbrace{\begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix}}_U \times \underbrace{\begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix}}_{\Sigma} \times \underbrace{\begin{bmatrix} 0.58 & 0 \\ 0.58 & 0 \\ 0.58 & 0 \\ 0 & 0.71 \\ 0 & 0.71 \end{bmatrix}}_{V^T}.$$

- Σ : # of categories and their relative strength in the data set
- U : “user to category (genre)” similarity matrix

SVD for movie data

Using the compact SVD,

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ \underline{5} & \underline{5} & \underline{5} & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \underbrace{\begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ \textcolor{purple}{0.90} & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix}}_U \times \underbrace{\begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix}}_{\Sigma} \times \underbrace{\begin{bmatrix} 0.58 & 0 \\ 0.58 & 0 \\ 0.58 & 0 \\ 0 & 0.71 \\ 0 & 0.71 \end{bmatrix}}_{V^T}.$$

Strongest SciFi preference (highest ratings)

- Σ : # of categories and their relative strength in the data set
- U : “user to category (genre)” similarity matrix

SVD for movie data

Using the compact SVD,

$$\left[\begin{array}{ccc|cc} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ \hline 5 & 5 & 5 & 0 & 0 \\ \hline 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{array} \right] = \underbrace{\left[\begin{array}{cc} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ \hline 0.90 & 0 \\ \hline 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{array} \right]}_U \times \underbrace{\left[\begin{array}{cc} 9.64 & 0 \\ 0 & 5.29 \end{array} \right]}_{\Sigma} \times \underbrace{\left[\begin{array}{cc} 0.58 & 0 \\ 0.58 & 0 \\ 0.58 & 0 \\ 0 & 0.71 \\ 0 & 0.71 \end{array} \right]}_{V^T}.$$

Annotations:

- Strongest SciFi preference (highest ratings) points to the circled value 0.90 in the matrix U.
- User 6 fits best w/ romance category (highest ratings) points to the circled value 0.80 in the matrix U.

- Σ : # of categories and their relative strength in the data set
- U : "user to category (genre)" similarity matrix

SVD for movie data

Using the compact SVD,

each user gives identical ratings

$$\left[\begin{array}{ccc|cc} \underline{1} & \underline{1} & \underline{1} & 0 & 0 \\ \underline{2} & \underline{2} & \underline{2} & 0 & 0 \\ \underline{1} & \underline{1} & \underline{1} & 0 & 0 \\ \underline{5} & \underline{5} & \underline{5} & 0 & 0 \\ \hline 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{array} \right] = \underbrace{\begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix}}_U \times \underbrace{\begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix}}_{\Sigma} \times \underbrace{\begin{bmatrix} 0.58 & 0 \\ 0.58 & 0 \\ 0.58 & 0 \\ 0 & 0.71 \\ 0 & 0.71 \end{bmatrix}}_{V^T}.$$

first 3 movies belong to SciFi to the same extent (and are separate from romance)

- Σ : # of categories and their relative strength in the data set
- U : “user to category (genre)” similarity matrix
- V : “category (genre) to movie” similarity matrix

SVD for movie data

Using the compact SVD,

each user gives identical ratings

$$\left[\begin{array}{ccc|cc} \underline{1} & \underline{1} & \underline{1} & 0 & 0 \\ \underline{2} & \underline{2} & \underline{2} & 0 & 0 \\ \underline{1} & \underline{1} & \underline{1} & 0 & 0 \\ \underline{5} & \underline{5} & \underline{5} & 0 & 0 \\ \hline 0 & 0 & 0 & \underline{2} & \underline{2} \\ 0 & 0 & 0 & \underline{3} & \underline{3} \\ 0 & 0 & 0 & \underline{1} & \underline{1} \end{array} \right] = \underbrace{\begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix}}_U \times \underbrace{\begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix}}_{\Sigma} \times \underbrace{\begin{bmatrix} 0.58 & 0 \\ 0.58 & 0 \\ 0.58 & 0 \\ 0 & 0.71 \\ 0 & 0.71 \end{bmatrix}}_{V^T}$$

all rate Casablanca and Amelie equally

first 3 movies belong to SciFi to the same extent (and are separate from)

belong equally to romance ❤

- Σ : # of categories and their relative strength in the data set
- U : "user to category (genre)" similarity matrix
- V : "category (genre) to movie" similarity matrix

Diversifying within a category

Suppose instead the second user — let's call him Pete — has non-homogenous ratings; he really prefers Alien to the other Sci Fi movies and rates it a 5.

Diversifying within a category

Suppose instead the second user — let's call him Pete — has non-homogenous ratings; he really prefers Alien to the other Sci Fi movies and rates it a 5.

- Now 3 nonzero singular values
- In compact SVD, we now need an extra column in U and V^T to span the Sci-fi columns of A (which are no longer identical)

$$\left[\begin{array}{ccc|cc} 1 & 1 & 1 & 0 & 0 \\ 2 & 5 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ \hline 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{array} \right] = \underbrace{\left[\begin{array}{ccc} 0.16 & 0 & 0.10 \\ 0.52 & 0 & -0.85 \\ 0.16 & 0 & 0.10 \\ 0.82 & 0 & 0.50 \\ \hline 0 & 0.53 & 0 \\ 0 & 0.80 & 0 \\ 0 & 0.27 & 0 \end{array} \right]}_U \left[\begin{array}{ccc} 10.47 & 0 & 0 \\ 0 & 5.29 & 0 \\ 0 & 0 & 2.11 \end{array} \right] \underbrace{\left[\begin{array}{ccccc} 0.52 & 0.67 & 0.52 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \\ 0.48 & -0.74 & 0.48 & 0 & 0 \end{array} \right]}_{V^T}$$

Diversifying within a category

Suppose instead the second user — let's call him Pete — has non-homogenous ratings; he really prefers Alien to the other Sci Fi movies and rates it a 5.

- Now 3 nonzero singular values
- In compact SVD, we now need an extra column in U and V^T to span the Sci-fi columns of A (which are no longer identical)

$$\left[\begin{array}{ccc|cc} 1 & 1 & 1 & 0 & 0 \\ 2 & 5 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ \hline 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{array} \right] = \underbrace{\left[\begin{array}{ccc} 0.16 & 0 & 0.10 \\ 0.52 & 0 & -0.85 \\ 0.16 & 0 & 0.10 \\ 0.82 & 0 & 0.50 \\ \hline 0 & 0.53 & 0 \\ 0 & 0.80 & 0 \\ 0 & 0.27 & 0 \end{array} \right]}_U \underbrace{\left[\begin{array}{ccc|cc} 10.47 & 0 & 0 & 0 & 0 \\ 0 & 5.29 & 0 & 0 & 0 \\ 0 & 0 & 2.11 & 0 & 0 \\ \hline 0.52 & 0.67 & 0.52 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \\ 0.48 & -0.74 & 0.48 & 0 & 0 \end{array} \right]}_{V^T}$$

= unchanged values (romance category)

Mixing categories

Instead of updating ratings within a category, what if our second user (Pete) is interested in Sci-fi *and* romance? In fact, he *loves* Amelie.

Mixing categories

Instead of updating ratings within a category, what if our second user (Pete) is interested in Sci-fi *and* romance? In fact, he *loves* Amelie.

- Now there is “fill-in” where there were previously 0s when diversified ratings still fell into separate categories

Mixing categories

Instead of updating ratings within a category, what if our second user (Pete) is interested in Sci-fi *and* romance? In fact, he *loves* Amelie.

- Now there is “fill-in” where there were previously 0s when diversified ratings still fell into separate categories
- The magnitude of the “fill-in” corresponds to the strength of the new connection i.e. the mixed Sci-fi-romance grouping

Mixing categories

Instead of updating ratings within a category, what if our second user (Pete) is interested in Sci-fi *and* romance? In fact, he *loves* Amelie.

- Now there is “fill-in” where there were previously 0s when diversified ratings still fell into separate categories
- The magnitude of the “fill-in” corresponds to the strength of the new connection i.e. the mixed Sci-fi-romance grouping
- A weaker connection (say $A_{2,5} = 2$ instead of 4) gives a smaller corresponding Σ_{33} , along with fill-in entries smaller in magnitude

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & \color{blue}{4} \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.17 & \color{red}{0.06} & -0.06 \\ 0.43 & -0.44 & 0.79 \\ 0.17 & \color{red}{0.06} & -0.06 \\ 0.86 & \color{red}{0.29} & -0.31 \\ \color{red}{0.05} & -0.45 & -0.28 \\ \color{red}{0.08} & -0.68 & -0.42 \\ \color{red}{0.03} & -0.23 & -0.14 \end{bmatrix} \begin{bmatrix} 9.80 & 0 & 0 \\ 0 & 5.97 & 0 \\ 0 & 0 & 2.30 \end{bmatrix}^T$$

 = fill in
from mixing
categories