

Recommender Systems and Dimensionality Reduction

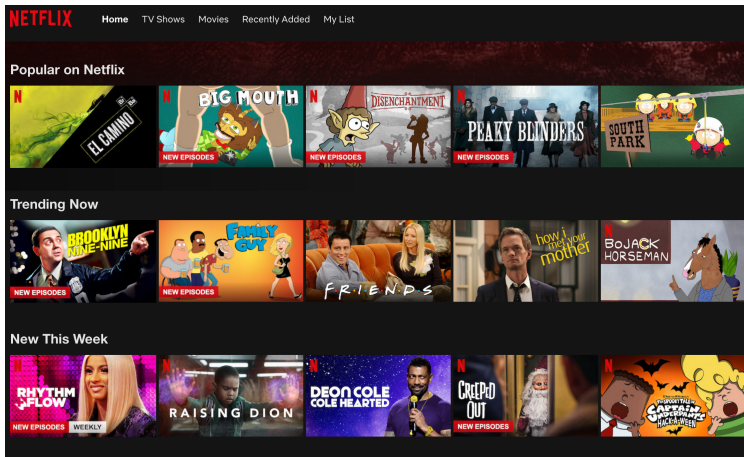
Sibylle Hess

1

Informal Problem Description



Recommending Movies like Netflix does



Who Would You Recommend What and Why?

	Star Wars	Interstellar	Blade Runner	Tron	2001: Space O.	Mars Attacks	Dune	Matrix	Robo Cop	Aliens	Terminator	Solaris	Avatar	12 Monkeys
Grace														
Carol														
Alice														
Bob														
Eve														
Chuck														

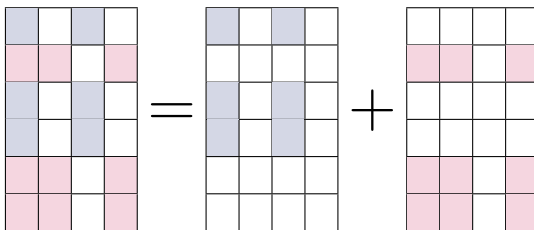


: Yeeey

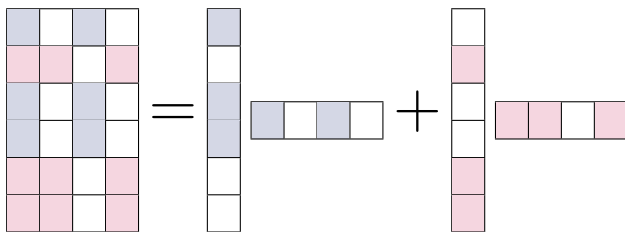


: Naaay

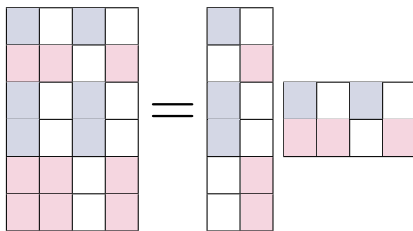
What is this Color Scheme in Math?



What is this Color Scheme in Math?



What is this Color Scheme in Math? A Matrix Product!



2

Derive the Formal Problem Definition

The Rank- r Matrix Factorization Problem

Given: a data matrix $D \in \mathbb{R}^{n \times d}$ and a rank $r < \min\{n, d\}$.

Find: matrices $X \in \mathbb{R}^{d \times r}$ and $Y \in \mathbb{R}^{n \times r}$ whose product approximates the data matrix:

$$\min_{X, Y} \|D - YX^T\|^2 \quad \text{s.t. } X \in \mathbb{R}^{d \times r}, Y \in \mathbb{R}^{n \times r}$$

Singular Value Decomposition

Theorem (SVD)

For every matrix $D \in \mathbb{R}^{n \times d}$ there exist orthogonal matrices $U \in \mathbb{R}^{n \times n}$, $V \in \mathbb{R}^{d \times d}$ and $\Sigma \in \mathbb{R}^{n \times d}$ such that

$$D = U\Sigma V^T, \text{ where}$$

- $U^T U = UU^T = I_n$, $V^T V = VV^T = I_d$
- Σ is a rectangular diagonal matrix, $\Sigma_{11} \geq \dots \geq \Sigma_{ll}$ where $l = \min\{n, d\}$

The column vectors $U_{\cdot s}$ and $V_{\cdot s}$ are called **left** and **right singular vectors** and the values $\sigma_i = \Sigma_{ii}$ are called **singular values** ($1 \leq i \leq l$).

Solutions to the Rank- r Matrix Factorization Problem

Theorem (Truncated SVD)

Let $D = U\Sigma V^T \in \mathbb{R}^{n \times d}$ be the singular decomposition of D .
Then the global minimizers X and Y of the rank- r MF problem

$$\min_{X,Y} \|D - YX^T\|^2 \text{ s.t. } X \in \mathbb{R}^{d \times r}, Y \in \mathbb{R}^{n \times r}.$$

satisfy

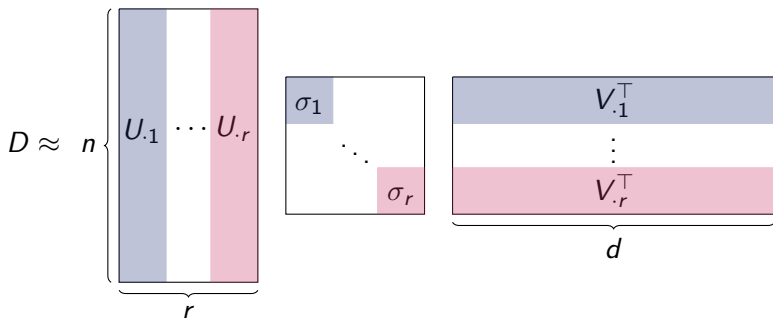
$$YX^T = U_{\mathcal{R}}\Sigma_{\mathcal{R}\mathcal{R}}V_{\mathcal{R}}^T, \text{ where } \mathcal{R} = \{1, \dots, r\}.$$

The proof follows from the orthogonal invariance of the Frobenius norm, yielding:

$$\min_{X,Y} \|D - YX^T\|^2 = \|\Sigma - U^T YX^T V\|^2$$

Truncated SVD

The approximation $D \approx U_{\cdot \mathcal{R}} \Sigma_{\mathcal{R} \mathcal{R}} V_{\cdot \mathcal{R}}^T$ is called **truncated SVD**.



Matrix Completion for Recommender Systems

		Movies			
		A	B	C	D
Users	1	★★★★★	?	★★☆☆☆	★★☆☆☆
	2	?	★★☆☆☆	★★★★★	?
	3	★★★★★	★★☆☆☆	★★★★★	★★☆☆☆
	4	★★★★★	?	★★★★★	★★★★☆
	5	★★★★★	★★★★★	?	?
	6	?	★★★★☆	★★★★★	★★★☆☆

Can we fill the ? with the rating which would be given by the user if (s)he had seen the movie?

Matrix Completion by SVD

Quick hack: replace the ? with the mean rating $\mu = 3$.

		Movies			
		A	B	C	D
Users	1	5	μ	2	1
	2	μ	1	5	μ
	3	5	1	5	2
	4	5	μ	5	3
	5	5	5	μ	μ
	6	μ	4	5	3

The Low-Rank Matrix Approximation Provides Recommendations

$$\begin{pmatrix} 5 & \mu & 1 & 1 \\ \mu & 1 & 5 & \mu \\ 2 & 1 & 5 & 3 \\ 4 & \mu & 4 & 2 \\ 5 & 5 & \mu & 1 \\ \mu & 1 & 5 & 3 \end{pmatrix} \approx \begin{pmatrix} 4.3 & 3.7 & 1.4 & 0.6 \\ 2.8 & 1.2 & 5.1 & 3.0 \\ 2.2 & 0.7 & 5.0 & 2.9 \\ 4.2 & 2.8 & 3.9 & 2.1 \\ 5.5 & 4.5 & 2.7 & 1.3 \\ 2.8 & 1.2 & 5.1 & 3.0 \end{pmatrix} \\
 = \begin{pmatrix} -0.3 & 0.5 \\ -0.4 & -0.4 \\ -0.4 & -0.4 \\ -0.4 & 0.1 \\ -0.5 & 0.5 \\ -0.4 & -0.4 \end{pmatrix} \begin{pmatrix} -9.0 & -5.8 & -9.5 & -5.3 \\ 2.6 & 3.3 & -3.3 & -2.2 \end{pmatrix}$$

Interpretation of MF for Recommender Systems

$$\begin{pmatrix} 5 & \mu & 1 & 1 \\ \mu & 1 & 5 & \mu \\ 2 & 1 & 5 & 3 \\ 4 & \mu & 4 & 2 \\ 5 & 5 & \mu & 1 \\ \mu & 1 & 5 & 3 \end{pmatrix} \approx \begin{pmatrix} -0.3 & 0.5 \\ -0.4 & -0.4 \\ -0.4 & -0.4 \\ -0.4 & 0.1 \\ -0.5 & 0.5 \\ -0.4 & -0.4 \end{pmatrix} \begin{pmatrix} -9.0 & -5.8 & -9.5 & -5.3 \\ 2.6 & 3.3 & -3.3 & -2.2 \end{pmatrix}$$

Every user's preferences are approximated by a linear combination of the rows in the second matrix:

$$\begin{aligned}
 (5 \quad \mu \quad 1 \quad 1) &\approx -0.3 \cdot (-9.0 \quad -5.8 \quad -9.5 \quad -5.3) \\
 &\quad + 0.5 \cdot (2.6 \quad 3.3 \quad -3.3 \quad -2.2)
 \end{aligned}$$

Matrix Completion by SVD

$$\begin{pmatrix} 5 & \mu & 1 & 1 \\ \mu & 1 & 5 & \mu \\ 2 & 1 & 5 & 3 \\ 4 & \mu & 4 & 2 \\ 5 & 5 & \mu & 1 \\ \mu & 1 & 5 & 3 \end{pmatrix} \approx \begin{pmatrix} 4.3 & 3.7 & 1.4 & 0.6 \\ 2.8 & 1.2 & 5.1 & 3.0 \\ 2.2 & 0.7 & 5.0 & 2.9 \\ 4.2 & 2.8 & 3.9 & 2.1 \\ 5.5 & 4.5 & 2.7 & 1.3 \\ 2.8 & 1.2 & 5.1 & 3.0 \end{pmatrix} \\
 = \begin{pmatrix} -0.3 & 0.5 \\ -0.4 & -0.4 \\ -0.4 & -0.4 \\ -0.4 & 0.1 \\ -0.5 & 0.5 \\ -0.4 & -0.4 \end{pmatrix} \begin{pmatrix} -9.0 & -5.8 & -9.5 & -5.3 \\ 2.6 & 3.3 & -3.3 & -2.2 \end{pmatrix}$$

Question: What happens if observations are sparse?

How can we prevent the approximation to the inserted mean values?

Adapt the objective to approximate only observed entries.

Making 3rd place in the Netflix Price 2009

Given: a data matrix $D \in \mathbb{R}^{n \times d}$ having observed entries D_{ik} for $(i, k) \in \mathcal{O} \subseteq \{1, \dots, n\} \times \{1, \dots, d\}$ the set of observed matrix entries, and a rank $r < \min\{n, d\}$.

Find: matrices $X \in \mathbb{R}^{d \times r}$ and $Y \in \mathbb{R}^{n \times r}$ whose product approximates the data matrix only on observed entries, indicated by $\mathbb{1}_{\mathcal{O}}$:

$$\min_{X, Y} \|\mathbb{1}_{\mathcal{O}} \circ (D - YX^T)\|^2 = \sum_{(i, k) \in \mathcal{O}} (D_{ik} - Y_i \cdot X_k^T)^2$$

s.t. $X \in \mathbb{R}^{d \times r}, Y \in \mathbb{R}^{n \times r}$

Optimization: Coordinate Descent

Truncated SVD solves the Rank- r Matrix Factorization Problem

Now something different:

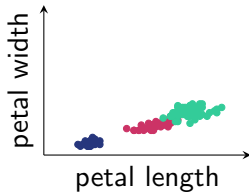
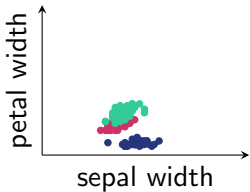
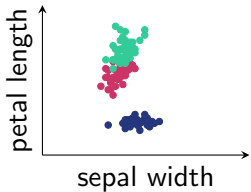
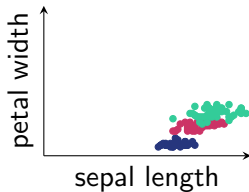
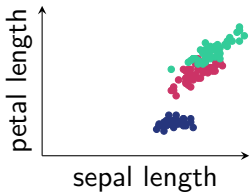
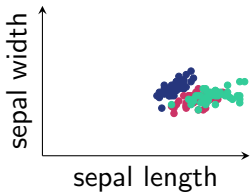
Finding low-dimensional
representations of the data by
truncated SVD

Exploring the Iris Dataset



sepal length	sepal width	petal length	petal width	class
5.1	3.5	1.4	0.2	setosa
6.4	3.5	4.5	1.2	versicolor
5.9	3.0	5.0	1.8	virginica
⋮	⋮	⋮	⋮	⋮

The First Step of Data Analysis: Visualization

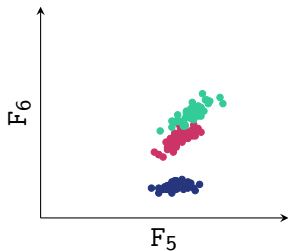


Which views are good?

We can also Generate our Own Features

$$F_5 = F_1 + F_2$$

$$F_6 = F_3 + F_4$$



- F_1 : sepal length
- F_2 : sepal width
- F_3 : petal length
- F_4 : petal width

How do we find good low-dimensional views on our data? How to create good new features?

Find the linear combination of features with highest variance.

2

Derive the Formal Problem Definition

Defining a new Feature by a Linear Combination

Given the $n \times d$ data matrix D gathering n observations of d features F_1, \dots, F_d , we define a new feature:

$$F_{d+1} = \sum_{k=1}^d \alpha_k F_k.$$

We have n observations of this new feature, given by

$$D_{\cdot,d+1} = \sum_{k=1}^d \alpha_k D_{\cdot,k} = D\alpha \in \mathbb{R}^n$$

The Sample Mean of the new Feature

Given observations $D_{:d+1} = D\alpha$ of the new feature $F_{d+1} = \sum_{k=1}^d \alpha_k F_k$, we compute the sample mean as

$$\mu_{F_{d+1}} = \frac{1}{n} \sum_{i=1}^n D_{id+1} = \boldsymbol{\mu}_F^\top \alpha, \quad \text{where } \boldsymbol{\mu}_F = \begin{pmatrix} \mu_{F_1} \\ \vdots \\ \mu_{F_d} \end{pmatrix}$$

is the vector gathering all sample means for the d features.

The Sample Variance of the new Feature

Given observations $D_{\cdot d+1} = D\alpha$ of the new feature

$$F_{d+1} = \sum_{k=1}^d \alpha_k F_k, \quad \text{with sample mean} \quad \mu_{F_{d+1}} = \mu_F^\top \alpha,$$

we compute the **sample variance** as

$$\sigma_{F_{d+1}}^2 = \frac{1}{n} \sum_{i=1}^n (D_{id+1} - \mu_{F_{d+1}})^2 = \frac{1}{n} \left\| \left(D - \mathbf{1} \mu_F^\top \right) \alpha \right\|^2$$

Sample Statistics of the new Feature

Given observations $D_{:d+1} = D\alpha$ of the new feature

$$F_{d+1} = \sum_{k=1}^d \alpha_k F_k,$$

the sample mean and variance is given by

$$\mu_{F_{d+1}} = \mu_F^\top \alpha, \quad \sigma_{F_{d+1}}^2 = \frac{1}{n} \left\| \left(D - \mathbf{1} \mu_F^\top \right) \alpha \right\|^2.$$

We are interested in the **direction** of maximal variance, so we can restrict the length of vector α : $\|\alpha\| = 1$

Finding the Direction of Maximal Sample Variance

The direction of largest variance α is the solution to the following optimization problem:

$$\begin{aligned}
 \max_{\|\alpha\|=1} \sigma_{d+1}^2 &= \max_{\|\alpha\|=1} \frac{1}{n} \left\| \left(D - \mathbf{1}\mu_F^\top \right) \alpha \right\|^2 \\
 &= \max_{\|\alpha\|=1} \frac{1}{n} \alpha^\top \left(D - \mathbf{1}\mu_F^\top \right)^\top \left(D - \mathbf{1}\mu_F^\top \right) \alpha \\
 &= \max_{\|\alpha\|=1} \frac{\alpha^\top C^\top C \alpha}{n},
 \end{aligned}$$

where $C = D - \mathbf{1}\mu_F^\top$ is the centered data matrix.

So, the direction of largest variance is given by the operator norm of the centered data matrix.

How can we derive a
low-dimensional representation
of the data?

Find the r orthogonal
directions of largest variance.

The Principal Components Analysis Task

Given: a data matrix $D \in \mathbb{R}^{n \times d}$ and a rank r .

Find: the r orthogonal direction of largest variance, given by the columns Z_s which are the solution to the following optimization problem:

$$\max_Z \text{tr}(Z^T C^T C Z) \quad \text{s.t. } Z \in \mathbb{R}^{n \times r}, Z^T Z = I$$

where $C = D - \mathbf{1}\mu_F^T$ is the centered data matrix.

3

Optimization

What is the solution Z of the objective of PCA?

The right singular vectors of C .

SVD Solves the Objective of PCA

Theorem (Value of the Operator Norm)

Let $C = U\Sigma V^T \in \mathbb{R}^{n \times d}$ be the SVD of the matrix C . The solution of the optimization problem

$$\max_Z \operatorname{tr}(Z^T C^T C Z) \quad \text{s.t. } Z \in \mathbb{R}^{n \times r}, Z^T Z = I$$

is given by $Z = V_{\mathcal{R}}$ for $\mathcal{R} = \{1, \dots, r\}$.

Proof (sketch): Show that the objective above is equivalent to

$$\min_Z \|C^T C - Z \Sigma_{\mathcal{R}\mathcal{R}}^2 Z^T\|^2 \quad \text{s.t. } Z \in \mathbb{R}^{n \times r}, Z^T Z = I.$$

Principal Components Analysis

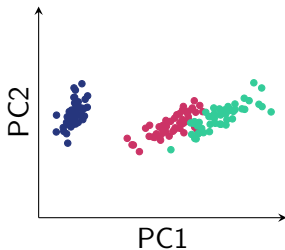
- 1: **function** PCA(D, r)
- 2: $C \leftarrow D - \mathbf{1}\mu_F^\top$ ▷ Center the data matrix
- 3: $(U_{\mathcal{R}}, \Sigma_{\mathcal{R}\mathcal{R}}, V_{\mathcal{R}}) \leftarrow \text{TRUNCATEDSVD}(C, r)$
- 4: **return** $CV_{\mathcal{R}}$ ▷ the low-dimensional view on the data
- 5: **end function**

PCA can be implemented such that the novel data representation is centered (returning $CV_{\mathcal{R}}$) or not (returning $DV_{\mathcal{R}}$).

Two-Dimensional PCA on the Iris Dataset

$$\text{PC1} = 0.36F_1 - 0.08F_2 + 0.85F_3 + 0.36F_4$$

$$\text{PC2} = 0.66F_1 + 0.73F_2 - 0.17F_3 - 0.07F_4$$



- F_1 : sepal length
- F_2 : sepal width
- F_3 : petal length
- F_4 : petal width

PCA enables

Dimensionality Reduction

Onto the Directions with

Maximal Variance