

CMSC5741 Big Data Tech. & Apps.

Lecture 6: MapReduce

Assignment Project Exam Help

<https://eduassistpro.github.io/>

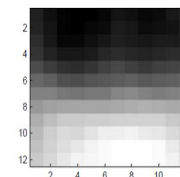
Add WeChat edu_assist_pro

Prof. Michael R. Lyu

Computer Science & Engineering Dept.

The Chinese University of Hong Kong

A Compression Example



Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Outline

- Motivation

- SVD

Assignment Project Exam Help

- CUR

– Application o <https://eduassistpro.github.io/>

- PCA

Add WeChat edu_assist_pro

– Extension to robust PCA

Dimensionality Reduction Motivation

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- **Assumption:** Data lie on or near a low d -dimensional subspace
- **Axes of this subspace are effective representation of the data**

Dimensionality Reduction Motivation

- **Compress / reduce dimensionality:**

- 10^6 rows; 10^3 columns; no updates

- Random access

Assignment Project Exam Help

error: OK

<https://eduassistpro.github.io/>

16

16

Add WeChat edu_assist_pro



The above matrix is really “2-dimensional.” All rows can be reconstructed by scaling $[1 \ 1 \ 1 \ 0 \ 0]$ or $[0 \ 0 \ 0 \ 1 \ 1]$

Rank of a Matrix

- **Q:** What is **rank** of a matrix **A**?
- **A:** No. of **linearly independent** rows/columns of **A**
- **For example:**
 - Matrix **A** =
$$\begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 1 & 1 & 2 \end{bmatrix}$$
 - **Why?** The first two rows are linearly independent, so the rank is at least 2, but all three rows are linearly dependent (the first is equal to the sum of the second and third) so the rank must be less than 3.
- **Why do we care about low rank?**
 - We can write **A** as two “basis” vectors: $\begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \end{bmatrix}$
 - And new coordinates of : $\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & -1 \end{bmatrix}$

Rank is “Dimensionality”

- **Cloud of points 3D space:**

- Think of point positions

- as a matrix: **Assignment Project Exam Help**

- 1 row per point:**

<https://eduassistpro.github.io/>

- **We can rewrite coordinates** **Add WeChat edu_assist_pro ntly!**

- Old basis vectors: $[1\ 0\ 0]$ $[0\ 1\ 0]$ $[0\ 0\ 1]$

- **New basis vectors:** $[1\ 2\ 1]$ $[-2\ -3\ 1]$

- Then **A** has new coordinates: $[1\ 0]$. **B:** $[0\ 1]$, **C:** $[1\ -1]$

- **Notice: We reduced the number of coordinates!**

Dimensionality Reduction

- Goal of dimensionality reduction is to discover the axis of data!

Assignment Project Exam Help

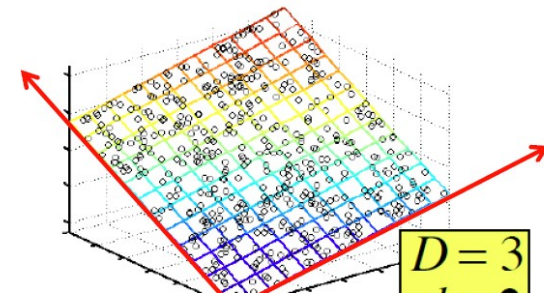
than representing
<https://eduassistpro.github.io/>oint with 2 coordinates
ent each point with
Add WeChat edu_assist_pro ate (corresponding to
the position of the point on
the red line).

By doing this we incur a bit of **error** as the points do not exactly lie on the line

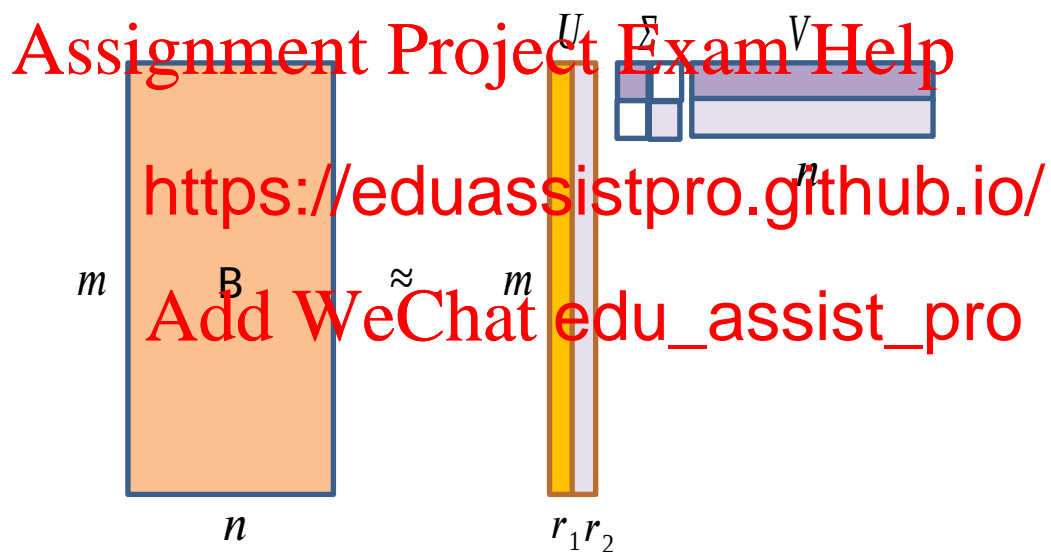
Why Reduce Dimensions?

Why reduce dimensions?

- **Discover hidden correlations/topics**
 - Words that o er
- **Remove redundant features**
 - Not all words are useful
- **Interpretation and visualization**
- **Easier storage and processing of the data**



SVD: Dimensionality Reduction



SVD: Singular Value Decomposition

- For an $m \times n$ matrix A , we can decompose it as $A = U \Sigma V^T$, where
 - U is an $m \times m$ real or complex orthonormal matrix (i.e., $U^H U = I$)
 - Σ is an $m \times n$ real diagonal matrix with non-negative real entries
 - V^T (the conjugate transpose of V if V is complex, the transpose of V if V is real) is an $n \times n$ real or complex orthonormal matrix.

SVD: Singular Value Decomposition

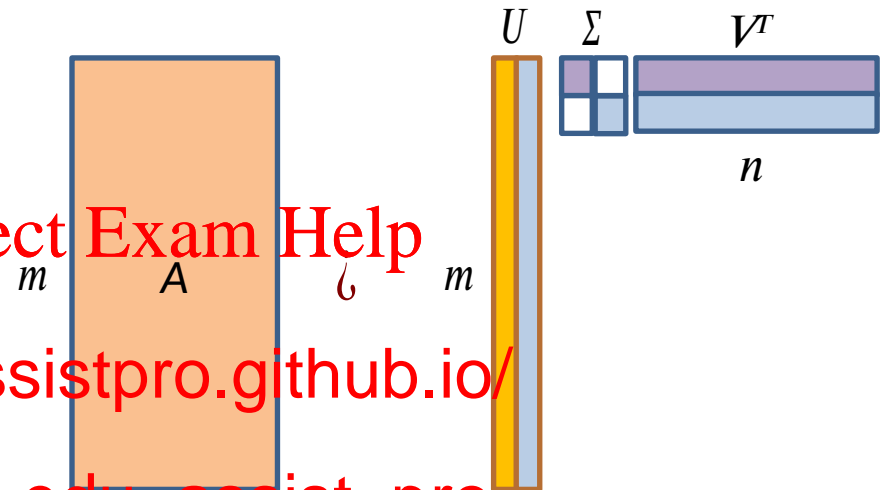
- When $\text{rank}(A) = r$:

- : input data matrix
 - matrix (e.g., document

- : left singular vectors
 - matrix (documents, topics)

- : singular values
 - diagonal matrix (strength of each “topic”)
 - rank of matrix

- : right singular vectors
 - matrix (terms, topics)



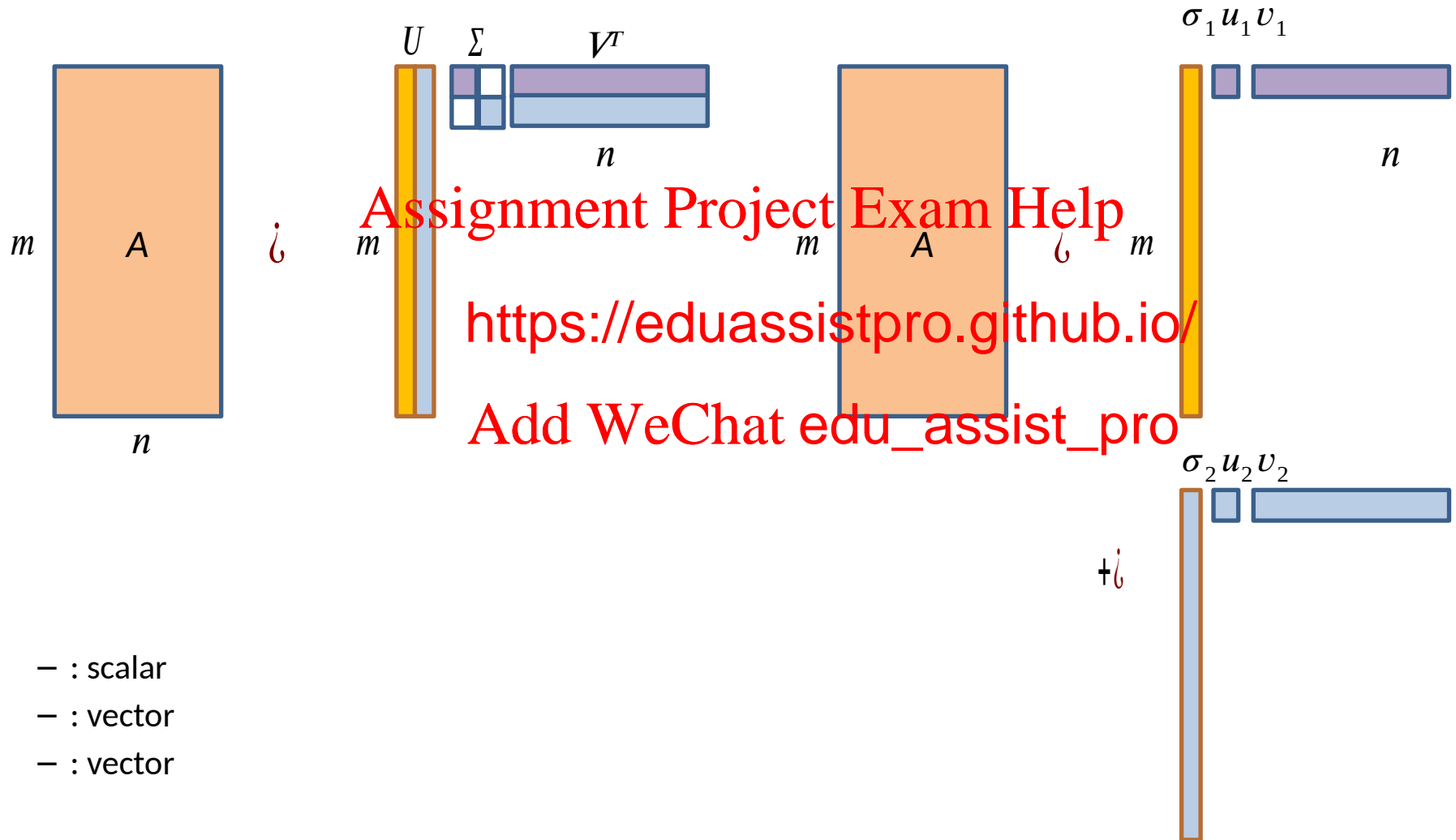
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- : scalar
- : vector
- : vector

SVD: Singular Value Decomposition



SVD Properties

- It is always possible to do SVD, i.e. decompose a matrix A into $U \Sigma V^T$, where
- U, Σ, V : unique
- U, V : column orth
– $U^T U = I, V^T V = I$ (I : identity matrix)
- Σ : diagonal
 - Entries (singular values) are non-negative,
 - Sorted in decreasing order ($\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$).

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

SVD Example

- We give an example of a simple SVD decomposition

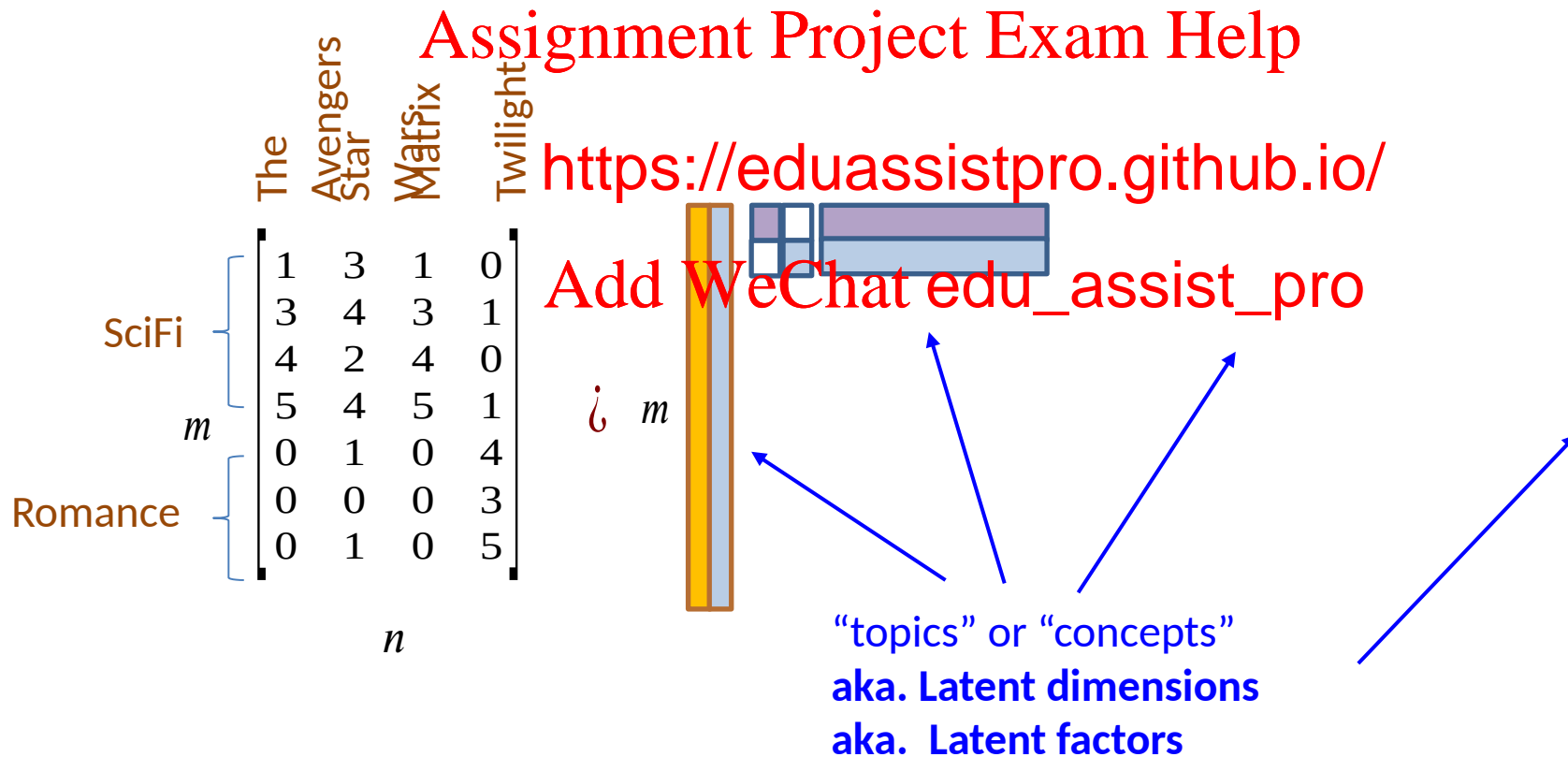
Assignment Project Exam Help

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 2 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ \sqrt{0.2} & 0 & 0 & 0 & \sqrt{0.8} \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{bmatrix}$$

A
 U
 Σ
 T

SVD: Example – Users-to-Movies

- example: Users to Movies



SVD: Example – Users-to-Movies

- example

SciFi-concept Romance-concept

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

$$\begin{matrix}
 \text{SciFi} \\
 \text{Romance}
 \end{matrix}
 \begin{matrix}
 \left[\begin{array}{cccc}
 1 & 3 & 1 & 0 \\
 3 & 4 & 3 & 1 \\
 4 & 2 & 4 & 0 \\
 5 & 4 & 5 & 1 \\
 0 & 1 & 0 & 4 \\
 0 & 0 & 0 & 3 \\
 0 & 1 & 0 & 5
 \end{array} \right]
 \end{matrix}
 \begin{matrix}
 \text{The Avengers Star Wars Twilight} \\
 m \\
 n
 \end{matrix}
 =
 \begin{matrix}
 \left[\begin{array}{ccc}
 0.48 & -0.02 & 0.43 \\
 0.49 & -0.08 & -0.5 \\
 0.68 & -0.07 & -0.16 \\
 0.06 & 0.57 & 0.06 \\
 0.01 & 0.41 & -0.20 \\
 0.07 & 0.70 & -0.01
 \end{array} \right]
 \begin{matrix}
 m \\
 n
 \end{matrix}
 \begin{matrix}
 \left[\begin{array}{cccc}
 0 & 0 & 0 & 0 \\
 7.1 & 0 & 0 & 2.5
 \end{array} \right]
 \end{matrix}
 \begin{matrix}
 X \\
 n
 \end{matrix}
 \begin{matrix}
 \left[\begin{array}{cccc}
 0.59 & 0.54 & 0.59 & 0.05 \\
 -0.10 & 0.12 & -0.10 & 0.98 \\
 -0.37 & 0.83 & -0.37 & -0.17
 \end{array} \right]
 \end{matrix}
 \begin{matrix}
 n
 \end{matrix}$$

SVD: Example – Users-to-Movies

- example

U is “user-to-concept”
similarity matrix

Assignment Project Exam Help

SciFi-concept Romance-concept

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

$$\begin{matrix}
 \text{SciFi} \\
 \text{Romance}
 \end{matrix}
 \begin{matrix}
 \left[\begin{array}{c}
 \text{The Avengers Star Wars Twilight} \\
 \begin{matrix} 1 & 3 & 1 & 0 \\ 3 & 4 & 3 & 1 \\ 4 & 2 & 4 & 0 \\ 5 & 4 & 5 & 1 \\ 0 & 1 & 0 & 4 \\ 0 & 0 & 0 & 3 \\ 0 & 1 & 0 & 5 \end{matrix}
 \end{array} \right]
 \end{matrix}
 \begin{matrix}
 m \\
 n
 \end{matrix}
 \begin{matrix}
 \left[\begin{array}{ccc}
 0.48 & -0.02 & 0.43 \\
 0.49 & -0.08 & -0.5 \\
 0.68 & -0.07 & -0.16 \\
 0.06 & 0.57 & 0.06 \\
 0.01 & 0.41 & -0.20 \\
 0.07 & 0.70 & -0.01
 \end{array} \right]
 \end{matrix}
 \begin{matrix}
 m \\
 n
 \end{matrix}
 \begin{matrix}
 \left[\begin{array}{ccc}
 0 & 0 & 0 \\
 7.1 & 0 & 0 \\
 0 & 0 & 2.5
 \end{array} \right]
 \end{matrix}
 \begin{matrix}
 n \\
 n
 \end{matrix}
 \begin{matrix}
 \left[\begin{array}{ccc}
 0.59 & 0.54 & 0.59 & 0.05 \\
 -0.10 & 0.12 & -0.10 & 0.98 \\
 -0.37 & 0.83 & -0.37 & -0.17
 \end{array} \right]
 \end{matrix}
 \begin{matrix}
 n \\
 n
 \end{matrix}
 \begin{matrix}
 X
 \end{matrix}$$

SVD: Example – Users-to-Movies

- example

Assignment Project Exam Help

SciFi concept

gth" of the SciFi-concept

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

SciFi

Romance

m

n

m

n

X

$$\begin{bmatrix}
 1 & 3 & 1 & 0 \\
 3 & 4 & 3 & 1 \\
 4 & 2 & 4 & 0 \\
 5 & 4 & 5 & 1 \\
 0 & 1 & 0 & 4 \\
 0 & 0 & 0 & 3 \\
 0 & 1 & 0 & 5
 \end{bmatrix}$$

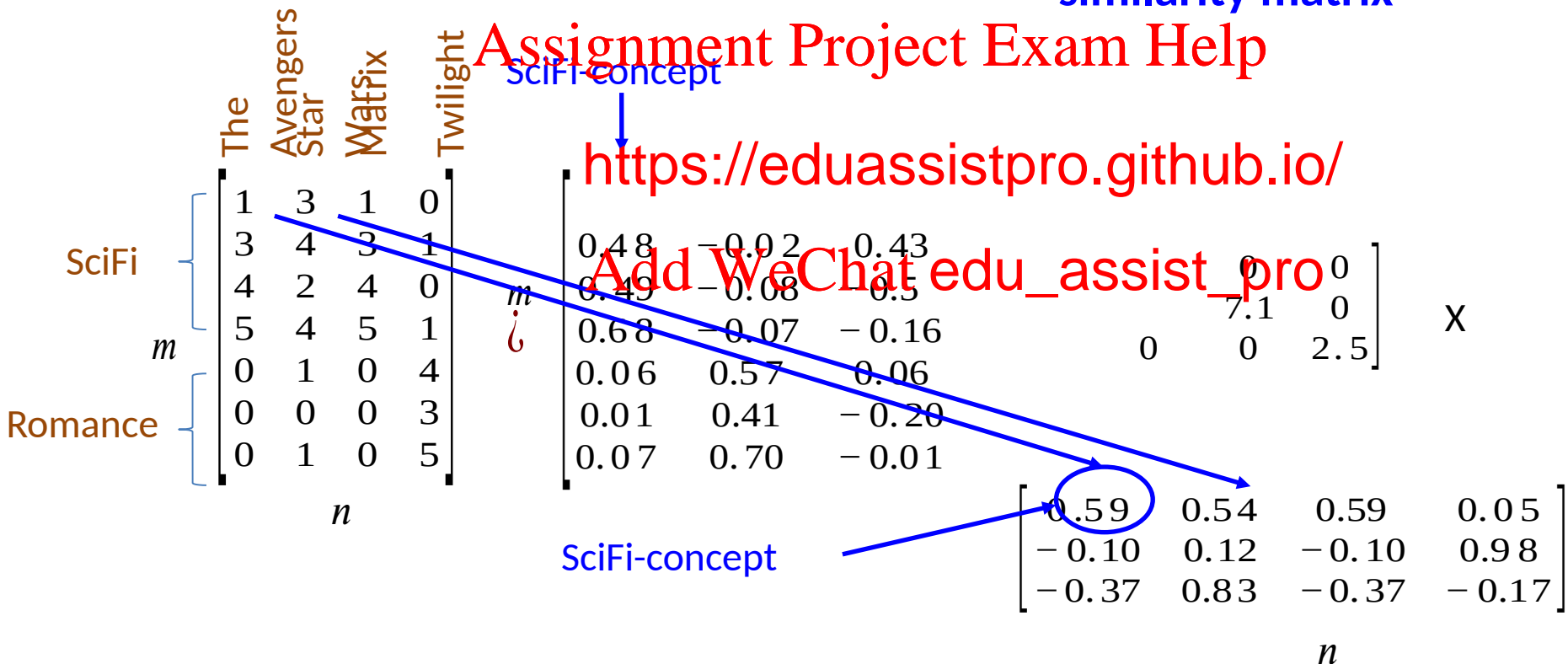
$$\begin{bmatrix}
 0.48 & -0.02 & 0.43 & 0 \\
 0.49 & -0.08 & 0.5 & 0 \\
 0.68 & -0.07 & -0.16 & 7.1 \\
 0.06 & 0.57 & 0.06 & 0 \\
 0.01 & 0.41 & -0.20 & 0 \\
 0.07 & 0.70 & -0.01 & 2.5
 \end{bmatrix}$$

$$\begin{bmatrix}
 0.59 & 0.54 & 0.59 & 0.05 \\
 -0.10 & 0.12 & -0.10 & 0.98 \\
 -0.37 & 0.83 & -0.37 & -0.17
 \end{bmatrix}$$

SVD: Example – Users-to-Movies

- example

V is “movie-to-concept”
similarity matrix



Q: Does the movie “Twilight” relate to concept “Romance”?

SVD: Interpretation #1

- “movies”, “users” and “concepts”
 - : user-to-concept similarity matrix
 - : movie-to-concept similarity matrix
 - : its diagonal
 - ‘strength’ of each concept

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

SVD: Interpretations #2

- SVD gives 'best' axis to project on
– 'best' = minimum squares of p errors
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro
- In other words,
minimum reconstruction error

SVD: Interpretation #2

- example

- U : user-to-concept matrix

- V : movie-to-c

Assignment Project Exam Help

<https://eduassistpro.github.io/>

$$\begin{bmatrix} 1 & 3 & 1 & 0 \\ 3 & 4 & 3 & 1 \\ 4 & 2 & 4 & 0 \\ 5 & 4 & 5 & 1 \\ 0 & 1 & 0 & 4 \\ 0 & 0 & 0 & 3 \\ 0 & 1 & 0 & 5 \end{bmatrix} \cdot \begin{bmatrix} 0.24 & 0.02 & 0.69 \\ 0.48 & -0.02 & 0.43 \\ 0.49 & -0.08 & -0.52 \\ 0.68 & -0.07 & -0.16 \\ 0.06 & 0.57 & 0.06 \\ 0.01 & 0.41 & -0.20 \\ 0.07 & 0.70 & -0.01 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 0 & 2.5 \end{bmatrix}$$

$$\begin{bmatrix} 0.59 & 0.54 & 0.59 & 0.05 \\ -0.10 & 0.12 & -0.10 & 0.98 \\ -0.37 & 0.83 & -0.37 & -0.17 \end{bmatrix}$$

SVD: Interpretation #2

- example

Assignment Project Exam Help

variance ("spread")

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

$$\begin{bmatrix} 1 & 3 & 1 & 0 \\ 3 & 4 & 3 & 1 \\ 4 & 2 & 4 & 0 \\ 5 & 4 & 5 & 1 \\ 0 & 1 & 0 & 4 \\ 0 & 0 & 0 & 3 \\ 0 & 1 & 0 & 5 \end{bmatrix}$$

;

$$\begin{bmatrix} 0.24 & 0.02 & 0.69 \\ 0.48 & -0.02 & 0.43 \\ 0.49 & -0.08 & -0.52 \\ 0.68 & -0.07 & -0.16 \\ 0.06 & 0.57 & 0.06 \\ 0.01 & 0.41 & -0.20 \\ 0.07 & 0.70 & -0.01 \end{bmatrix}$$

$$\begin{bmatrix} X & 0 & 0 & 0 \\ 0 & 0 & 0 & 2.5 \end{bmatrix} X$$

$$\begin{bmatrix} 0.59 & 0.54 & 0.59 & 0.05 \\ -0.10 & 0.12 & -0.10 & 0.98 \\ -0.37 & 0.83 & -0.37 & -0.17 \end{bmatrix}$$

SVD: Interpretation #2

- example

- : the coordinates of the points in the projection axis

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Projection of users on

the "Sci-Fi" axis

Add WeChat edu_assist_pro

$$\begin{bmatrix} 1 & 3 & 1 & 0 \\ 3 & 4 & 3 & 1 \\ 4 & 2 & 4 & 0 \\ 5 & 4 & 5 & 1 \\ 0 & 1 & 0 & 4 \\ 0 & 0 & 0 & 3 \\ 0 & 1 & 0 & 5 \end{bmatrix}$$

$$\begin{bmatrix} 2.86 & 0.24 & 8.21 \\ 5.71 & -0.24 & 5.12 \\ 5.83 & -0.95 & -6.19 \\ 8.09 & -0.83 & -1.90 \\ 0.71 & 6.78 & 0.71 \\ 0.12 & 4.88 & -2.38 \\ 0.83 & 8.33 & -0.12 \end{bmatrix}$$

SVD: Interpretation #2

- Q: how exactly is dimension reduction done?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

$$\begin{bmatrix} 1 & 3 & 1 & 0 \\ 3 & 4 & 3 & 1 \\ 4 & 2 & 4 & 0 \\ 5 & 4 & 5 & 1 \\ 0 & 1 & 0 & 4 \\ 0 & 0 & 0 & 3 \\ 0 & 1 & 0 & 5 \end{bmatrix} = \begin{bmatrix} 0.24 & 0.02 & 0.69 \\ 0.48 & -0.02 & 0.43 \\ 0.49 & -0.08 & -0.52 \\ 0.68 & -0.07 & -0.16 \\ 0.06 & 0.57 & 0.06 \\ 0.01 & 0.41 & -0.20 \\ 0.07 & 0.70 & -0.01 \end{bmatrix} \begin{bmatrix} X & 0 & 0 & 0 \\ 0 & 0 & 0 & 2.5 \end{bmatrix} \begin{bmatrix} X & 0 & 0 & 0 \\ 0 & 0 & 0 & 2.5 \end{bmatrix}$$

Add WeChat edu_assist_pro

$$\begin{bmatrix} 0.59 & 0.54 & 0.59 & 0.05 \\ -0.10 & 0.12 & -0.10 & 0.98 \\ -0.37 & 0.83 & -0.37 & -0.17 \end{bmatrix}$$

SVD: Interpretation #2

- Q: how exactly is dimension reduction done?

- A: Set smallest singular values to zero

Assignment Project Exam Help

<https://eduassistpro.github.io/>

$$\begin{bmatrix} 1 & 3 & 1 & 0 \\ 3 & 4 & 3 & 1 \\ 4 & 2 & 4 & 0 \\ 5 & 4 & 5 & 1 \\ 0 & 1 & 0 & 4 \\ 0 & 0 & 0 & 3 \\ 0 & 1 & 0 & 5 \end{bmatrix}$$

;

$$\begin{bmatrix} 0.24 & 0.02 & 0.69 \\ 0.48 & -0.02 & 0.43 \\ 0.49 & -0.08 & -0.52 \\ 0.68 & -0.07 & -0.16 \\ 0.06 & 0.57 & 0.06 \\ 0.01 & 0.41 & -0.20 \\ 0.07 & 0.70 & -0.01 \end{bmatrix}$$

$$\begin{bmatrix} X & 0 & 0 & 0 \\ 0 & 0 & 0 & 2.5 \end{bmatrix}$$

$$\begin{bmatrix} 0.59 & 0.54 & 0.59 & 0.05 \\ -0.10 & 0.12 & -0.10 & 0.98 \\ -0.37 & 0.83 & -0.37 & -0.17 \end{bmatrix}$$

SVD: Interpretation #2

- Q: how exactly is dimension reduction done?
- A: Set smallest singular values to zero
 - Approximate -rank matrices

$$\begin{bmatrix} 1 & 3 & 1 & 0 \\ 3 & 4 & 3 & 1 \\ 4 & 2 & 4 & 0 \\ 5 & 4 & 5 & 1 \\ 0 & 1 & 0 & 4 \\ 0 & 0 & 0 & 3 \\ 0 & 1 & 0 & 5 \end{bmatrix} \approx \begin{bmatrix} 0.24 & 0.02 & 0.69 & 0.43 \\ 0.48 & -0.02 & 0.43 & 0.52 \\ 0.49 & -0.08 & -0.52 & 0.51 \\ 0.68 & -0.07 & -0.16 & 0.51 \\ 0.06 & 0.57 & 0.06 & 0.51 \\ 0.01 & 0.41 & -0.20 & 0.51 \\ 0.07 & 0.70 & -0.01 & 0.51 \end{bmatrix} \begin{bmatrix} X & 0 & 0 & 2.5 \\ 0 & 0 & 0 & 2.5 \\ 0 & 0 & 0 & 2.5 \\ 0 & 0 & 0 & 2.5 \end{bmatrix} \begin{bmatrix} X & 0 & 0 & 2.5 \\ 0 & 0 & 0 & 2.5 \\ 0 & 0 & 0 & 2.5 \\ 0 & 0 & 0 & 2.5 \end{bmatrix}$$

<https://eduassistpro.github.io/>
 Add WeChat edu_assist_pro

$$\begin{bmatrix} 0.59 & 0.54 & 0.59 & 0.05 \\ -0.10 & 0.12 & -0.10 & 0.98 \\ -0.37 & 0.83 & -0.37 & -0.17 \end{bmatrix}$$

SVD: Interpretation #2

- Q: how exactly is dimension reduction done?
- A: Set smallest singular values to zero
 - Approximate -rank matrices

$$\begin{bmatrix} 1 & 3 & 1 & 0 \\ 3 & 4 & 3 & 1 \\ 4 & 2 & 4 & 0 \\ 5 & 4 & 5 & 1 \\ 0 & 1 & 0 & 4 \\ 0 & 0 & 0 & 3 \\ 0 & 1 & 0 & 5 \end{bmatrix}$$

\approx

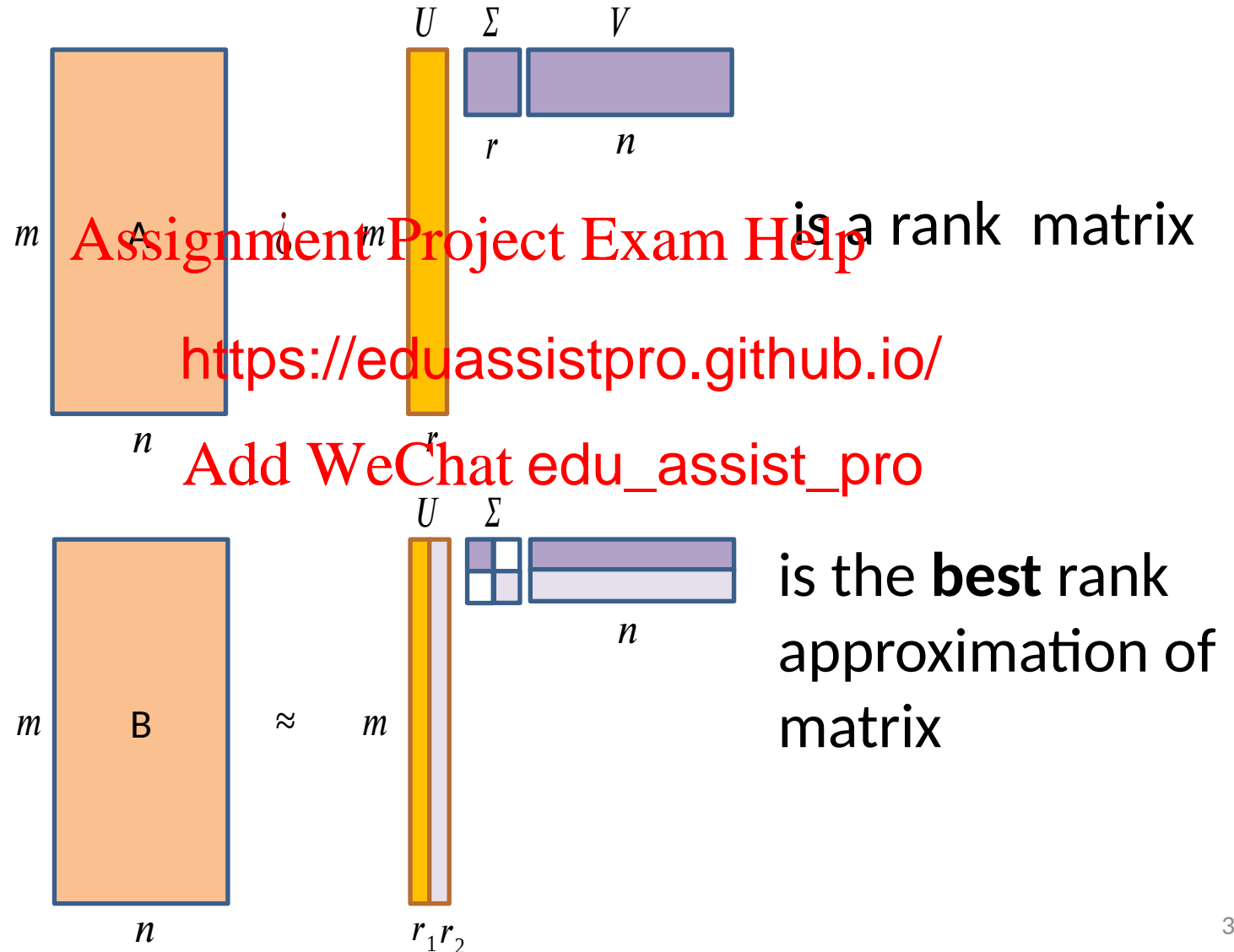
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

X

X

SVD: Best Low Rank Approximation



SVD: Best Low Rank Approximation

- **Theorem**: Let U , Σ , and V

- Σ = diagonal matrix where σ_i (and σ_{i+1})

- or equivalently, Σ is the best rank- k approximation to A

- or equivalently,

- Intuition (spectra

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- Why setting small σ_i to 0 is the right thing to do?

- Vectors u_i and v_i are unit length, so σ_i scales them.

- Therefore, zeroing small σ_i introduces less error.

SVD: Interpretation #2

- Q: How many σ_i to keep?

- A: Rule-of-a thumb

Keep 80~90

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

$$\begin{matrix} m & \begin{bmatrix} 1 & 3 & 1 & 0 \\ 3 & 4 & 3 & 1 \\ 4 & 2 & 4 & 0 \\ 5 & 4 & 5 & 1 \\ 0 & 1 & 0 & 4 \\ 0 & 0 & 0 & 3 \\ 0 & 1 & 0 & 5 \end{bmatrix} \\ & n \end{matrix}$$

$$\sigma_1 u_1 u_1^T + \sigma_2 u_2 u_2^T + \cdots$$

Assume: $\sigma_1 \geq \sigma_2 \geq \cdots$

SVD: Complexity

- SVD for full matrix

Assignment Project Exam Help

- But

- Less work, if we only want first k singular values
- or if we only want first k singular values (thin-svd).
- or if the matrix is sparse (sparse svd).

- Stable implementations

- LINPACK, Matlab, Splus, Mathematica...
- Available in most common languages

SVD: Conclusions so far

- SVD: : unique
 - user-to-concept similarities
 - movie-to-co
 - : strength to <https://eduassistpro.github.io/>
- Dimensionality reduction
 - Keep the few largest singular values (80-90% of “energy”)
 - SVD: picks up linear correlations

SVD: Relationship to Eigen-decomposition

- SVD gives us

Assignment Project Exam Help

- Eigen-decomp

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- is symmetric
- are orthonormal (
- are diagonal

SVD: Relationship to Eigen-decomposition

- Eigen-decomposition of Σ and Λ

Assignment Project Exam Help

<https://eduassistpro.github.io/>

– So, U , Σ , and V

– That is, U is the matrix of eigenvectors of Σ and V is the matrix of eigenvectors of Λ

– This shows how to use eigen-decomposition to compute SVD

– The singular values of Σ are the square roots of the corresponding eigenvalues of Λ

- Note: Σ and Λ are the dataset covariance matrices

A Brief Review of Eigen-Decomposition

- Eigenvalues and eigenvectors

- matrix.
 - eigenvalue of , :
- <https://eduassistpro.github.io/>

- Simple computation

- Solve the equation
- Example

- Then
- Then
- Solve , we get

A Brief Review of Eigen-Decomposition

- Example (continued)

Assignment Project Exam Help

- solve , we get <https://eduassistpro.github.io/>
- now we compute an eigenvector [Add WeChat edu_assist_pro](#)
 - for eigenvalue we need to find
 - solve
 - We get Since needs to be a unit vector, therefore
- Similarly, we can compute

Computing Eigenvalues: Power Method

- Power method
 - choose an arbitrary
 - .
 - Theorem: sequence converges to the principal eigenvector (i.e., the eigenvector corresponding to the largest eigenvalue)
- Normalized power method
 - choose an arbitrary
 - Theorem: sequence converges to the principal eigenvector.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

In-class Practice

- Go to [practice](#)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

SVD Case Study: How to Query?

- Q: Find users that like “The Avengers”
- A: Map query into a “concept space” – how?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

	The Avengers	Star Matrix	Twilight
SciFi	1	3	1
	3	4	3
	4	2	4
	5	4	5
Romance	0	1	0
	0	0	0
	0	1	0

$$m \begin{bmatrix} 0.24 & 0.02 & 0.69 \\ 0.48 & -0.02 & 0.43 \\ 0.49 & -0.08 & -0.52 \\ 0.68 & -0.07 & -0.16 \\ 0.06 & 0.57 & 0.06 \\ 0.01 & 0.41 & -0.20 \\ 0.07 & 0.70 & -0.01 \end{bmatrix}$$

$$X \begin{bmatrix} 11.9 & 0 & 0 \\ 0 & 7.1 & 0 \\ 0 & 0 & 2.5 \end{bmatrix} X$$

$$\begin{bmatrix} 0.59 & 0.54 & 0.59 & 0.05 \\ -0.10 & 0.12 & -0.10 & 0.98 \\ -0.37 & 0.83 & -0.37 & -0.17 \end{bmatrix}$$

Case Study: How to Query?

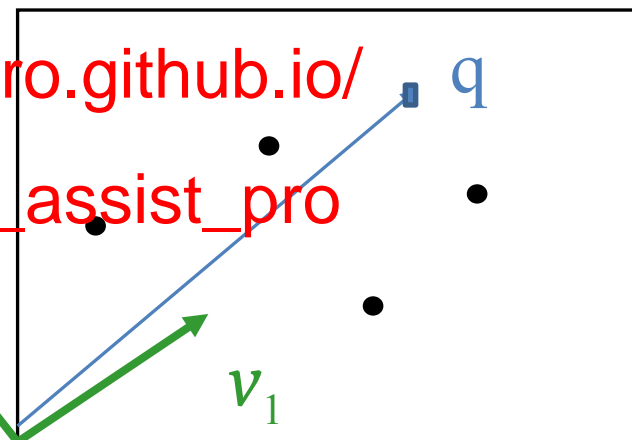
- Q: Find users that like “The Avengers”
- A: Map query into a “concept space” – how?

Assignment Project Exam Help

	The Avengers	Star Wars	Matrix	Twilight
q	[5	0	0	0]

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



Project into concept space:

Inner product with each concept
vector v_i

Case Study: How to Query?

- Q: Find users that like “The Avengers”
- A: Map query into a “concept space” – how?

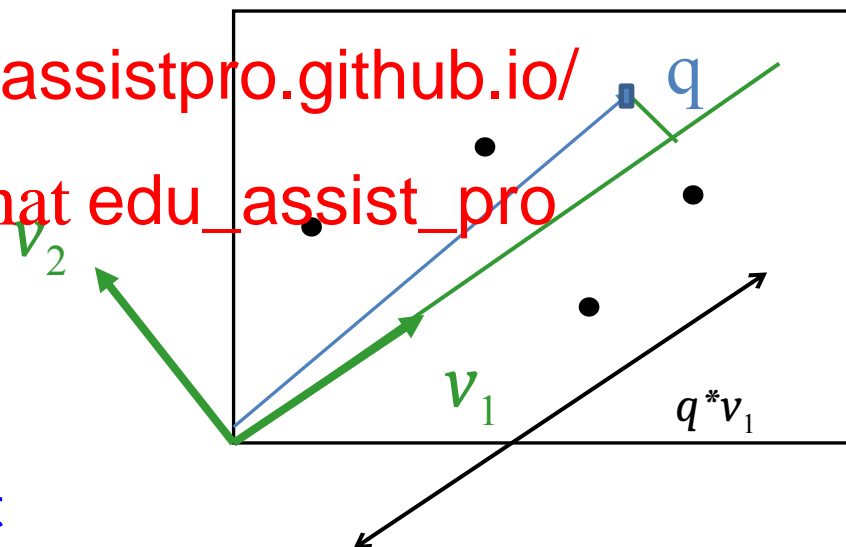
Assignment Project Exam Help

	The Avengers	Star Wars	Matrix	Twilight
q	[5	0	0	0]

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Project into concept space:
Inner product with each concept
vector v_i



Case Study: How to Query?

- Compactly, we have

$$-q_c = qV$$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro^{ept}

$$q \begin{matrix} & \text{The} & \text{Avengers} & \text{Matrix} & \text{Twilight} \\ & \text{5} & \text{0} & \text{0} & \text{0} \end{matrix} \times \begin{bmatrix} 0.59 & -0.10 \\ 0.54 & 0.12 \\ 0.59 & -0.10 \\ 0.05 & 0.98 \end{bmatrix} = \begin{bmatrix} 2.95 & -0.50 \end{bmatrix}$$

movie-to-concept
similarities ()

Case Study: How to Query?

- How would the user d that rated ('Star Wars', 'Matrix') be handled?

– $d_c = d V$ Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat [edu_assist_pro](#)^{ept}

$$\begin{array}{c}
 \text{The} \\
 \text{Avengers} \\
 \text{Star} \\
 \text{Wars} \\
 \text{Matrix} \\
 \text{Twilight}
 \end{array}
 \begin{array}{c}
 d \\
 [0 \quad 4 \quad 5 \quad 0]
 \end{array}
 \times
 \begin{bmatrix}
 0.59 & -0.10 \\
 0.54 & 0.12 \\
 0.59 & -0.10 \\
 0.05 & 0.98
 \end{bmatrix}
 =
 \begin{bmatrix}
 5.11 & -0.02
 \end{bmatrix}$$

movie-to-concept
similarities ()

Case Study: How to Query?

- Observation

- User d that rated ('Star Wars') will be similar to user q that rate ('The Avengers'), although d and q have zero ratings in common

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

ept



	The Avengers Star	War Matrix	Twilight
d	[0 4 5 0]		
q	[5 0 0 0]		

----->

[5.11 -0.02]

----->

[2.95 -0.50]

Zero ratings in common

Similarity $\neq 0$

Cosine similarity: 0.99

SVD: Drawbacks

+ Optimal low-rank approximation

in terms of Euclidean norm

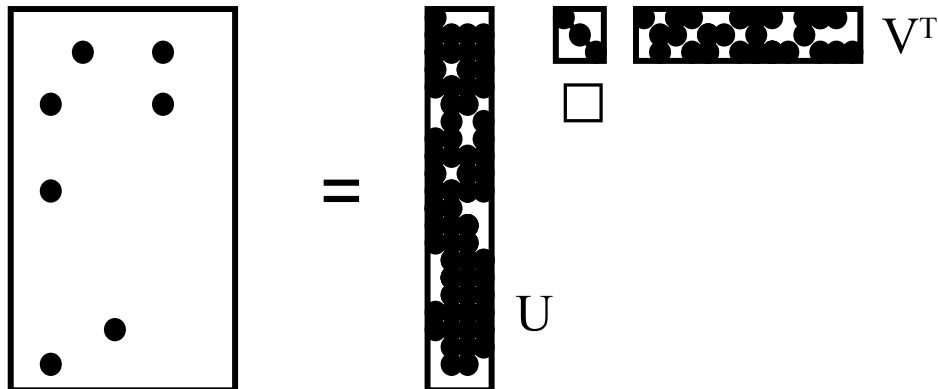
– Interpretability problem:

- A singular vector specifies a linear combination

<https://eduassistpro.github.io/>

– Lack of sparsity:

- Singular vectors are **dense**!



CUR Decomposition

- Goal: express A as a product of matrices
 - Minimize
- Constraints o

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

CUR Decomposition

- Goal: express A as a product of matrices
 - Minimize
- Constraints o

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

CUR: Good Approximation to SVD

- Let A_k be the best rank k approximation of A (obtain by SVD)

Assignment Project Exam Help

- Theorem

- CUR algorithm <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- with probability at least $1 - \epsilon$, by picking
 - columns and
 - rows
 - (in practice, choose k columns/rows)

CUR: How it Works

- Sampling columns (similarly for rows):

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

CUR: Computing U

- Let C be the “intersection” of sampled columns C and rows R
 - Let SVD of C be $U_C \Sigma_C V_C^T$
- Then: $U \approx U_C U_C^T U$, where
 - U_C is the “Moore-Penrose ps”
 - $U_C^T U_C = I$, if C has full column rank.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

CUR: Pros & Cons

+ easy interpretation

– the basis vectors are actual columns and rows

- duplicate columns

– columns of A are repeated many times

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

CUR: Duplicate Columns

- If we want to get rid of the duplicates
 - Throw them away
 - Scale the columns by the square root of the number of duplicates

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

SVD vs CUR

Question: Large or small? Dense or sparse?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

SVD vs CUR

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

SVD & CUR: Simple Experiments

- DBLP data
 - author-to-conference matrix
 - very sparse
 - : number of pap at conference .
 - 428k authors (r
 - 3659 conferences (column)
- Dimensionality reduction
 - Running time?
 - Space?
 - Reconstruction error?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Results: DBLP

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

courtesy: Sun, Faloutsos: *Less is more, Compact Matrix Decomposition for Large Sparse Graph*, SDM'07

- accuracy: 1-relative sum squared errors
- space ratio: $\# \text{ output non-zero matrix entries} / \# \text{ input non-zero matrix entries}$

The Linearity Assumption

- SVD is limited to linear projections

- Data lies on a low-dimensional linear space

Assignment Project Exam Help

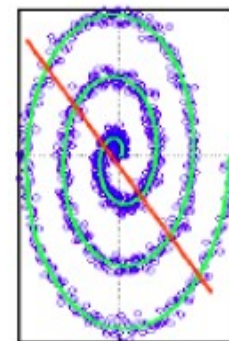
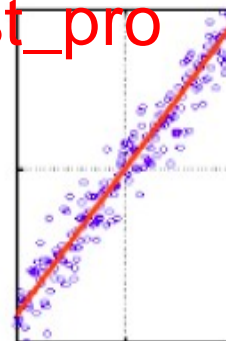
- Non-linear me

- Data lies on a <https://eduassistpro.github.io/>

- Non-linear Add WeChat edu_assist_pro

- How?

- Build adjacency graph
 - SVD the graph adjacency matrix
 - Further reading: wikipage of Isomap



PCA: An Application of SVD

- PCA = Principle Component Analysis

Assignment Project Exam Help

- Motivation

- Visualization

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

PCA: Data Visualization

- Example:
 - Given 53 blood samples (features) from 65 people (data item or instance)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- How can we visualize the samples

PCA: Data Visualization

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

How can we visualize the other variables???

... difficult to see in 4 or higher dimensional spaces ...

PCA: Data Visualization

- Is there a representation better than the coordinate axes?
- Is it really necessary to show all the 53 dimensions?

Assignment Project Exam Help

- What if there are too many features?
<https://eduassistpro.github.io/>
- How could we find the smallest subspace of the 53-D space that keeps the most information of the original data?
Add WeChat edu_assist_pro

- A solution: Principal Component Analysis
 - An application of SVD.

PCA: Definition and Algorithms

- PCA

- Orthogonal projection of the data onto a lower-dimensional line

that

- Maximize variance of projected data (purple line)
 - Minimize mean squared distance between

- Data point
 - Projection (sum of blue lines)

- Look data from a literally different angle.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

PCA: Idea

- Given data points in a d -dimensional space, project them into a lower dimensional space while preserving as much information as possible.
 - Find best plane for 3-D data
 - Find best 12-D approximation for 100-D data
- In particular, choose projection that minimizes squared error in reconstructing the original data.
 - Implement through SVD

PCA

- **PCA Vectors** originate from the center of mass.
- Principal component #1: points in the direction of the **largest variance**
- Each subsequent component
 - is **orthogonal** to the previous ones
 - points in the directions of the **largest variance of the residual subspace**

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

PCA: 2D Gaussian dataset

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

1st PCA axis

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

2nd PCA axis

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

PCA: Algorithm

- Given centered data , compute principle vectors

- 1st principle vector

Assignment Project Exam Help

- maximize the <https://eduassistpro.github.io/> of

Add WeChat edu_assist_pro

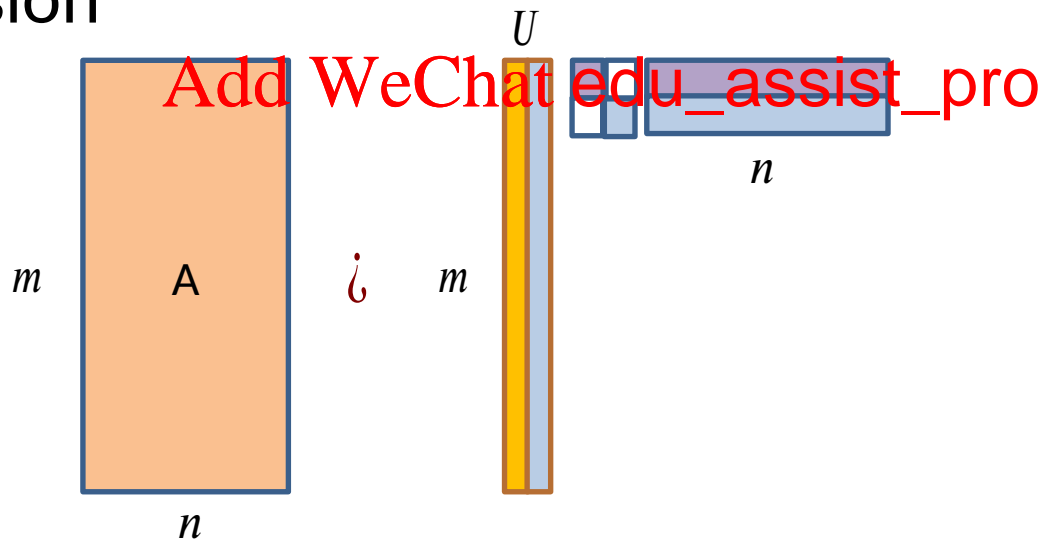
PCA: Algorithm by SVD

- SVD of the centered data matrix

Assignment Project Exam Help

– : number of i

– : dimension <https://eduassistpro.github.io/>



PCA: Algorithm by SVD

- Columns of
 - is exactly the principal vectors.
 - orthogonal a
- Matrix
 - Diagonal
 - Strength of each eigenvector
- Columns of
 - Coefficients for reconstructing the samples.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Application: Face Recognition

- Want to identify specific person, based on facial image
- Can't just use the given 256 x 256 pixels

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Applying PCA

- **Method A:** Build a PCA subspace for each person and check which subspace can reconstruct the test image the best
- **Method B:** Build a PCA subspace for the whole dataset and then classify based on the weights.
- Example data set: Images of faces
- Each face is ...
 - values

Principal Components (Method B)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Reconstructing ... (Method B)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Happiness Subspace (Method A)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Disgust Subspace (Method A)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Image Compression

Assignment Project Exam Help

- Divide the original 372x492 image into patches:
 - Each patch is an instance that contains 12x12 pixels on a grid
 - View each as a 144-D vector

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

PCA Compression: 144D => 60D

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

PCA Compression: 144D => 16D

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

PCA Compression: 144D => 6D

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

6 Most Important Eigenvectors

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

PCA Compression: 144D => 3D

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

3 Most Important Eigenvectors

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Noisy Image

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Denoised Image using 15 PCA Components

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

PCA: Shortcomings

- PCA cannot capture non-linear structure
 - Similar with SVD

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

PCA: Shortcomings

- PCA does not know labels

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

PCA: Conclusions

- PCA
 - find orthonormal basis for data
 - sort dimensions in order of “strength”
 - discard low significance dimensions
- Uses <https://eduassistpro.github.io/>
 - Get compact description
 - Ignore noise
 - Improve classification (hopefully)
- Not magic:
 - Doesn't know class labels
 - Can only capture linear variations
 - One of many tricks to reduce dimensionality!

Extra: Compute PCA Using Eigen-Decomposition

- Given centered data compute covariance matrix

Assignment Project Exam Help

- Top PCA components are the principal vectors of Σ .
 - Equivalence position and SVD
 - SVD decomposition of Σ ,
 - SVD-based algorithm for PCA
 - Eigen-decomposition of Σ .
 - Eigen-based algorithm for PCA
 - The equivalence gives $\mathbf{U} = \mathbf{V}$.

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

One-slide Takeaway

- Dimensionality reduction
 - compress/reduce dimension
 - reconstruct the original matrix by two or more smaller matrices
- Singular value dec
 - decompose a matrix $A = U \Sigma V^T$
 - U : column-orthonormal. Σ : diagonal matrix.
- CUR decomposition
 - set of C columns of A . set of R rows of A .
- Principle component analysis (PCA)
 - reconstruct data matrix by a smaller number of eigenvectors
 - view the data from a *literally* different angle.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

In-class Practice

- **1.** Describe briefly (informally or formally) the relationship between singular value decomposition and eigenvalue decomposition.
- **2.1** Compute the eigenvalues and eigenvectors of matrix
- **2.2** Let $A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$. It is easy to check that $A^2 = A + I$. What are the singular values of A ?
- **2.3** Obtain SVD for A where $A = U \Sigma V^T$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro