# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2022.
Lecture 18

- Problem Set 3 is due Monday at 11:59pm.
- No quiz due Monday.
- I will hold additional office hours on Thursday from 11:30am-12:40pm.

## Summary

### Last Class

- The Singular Value Decomposition (SVD) and its connection to eigendecomposition of $X^TX$ and $XX^T$, and low-rank approximation.

### This Class: Application of Low-Rank Approximation Beyond Compression

- Low-rank matrix completion (predicting missing measurements using low-rank structure).
- Entity embeddings (e.g., word embeddings, node embeddings).
- Low-rank approximation for non-linear dimensionality reduction.
- Eigendecomposition to partition graphs into clusters.

## SVD Review

- Every $X \in \mathbb{R}^{n \times d}$ can be written in its SVD as $U\Sigma V^T$.

$X =$

4

## SVD Review



- Every $X \in \mathbb{R}^{n \times d}$ can be written in its SVD as $U\Sigma V^T$.
  $$\text{rank}(X) = r$$
- $U \in \mathbb{R}^{n \times r}$ (orthonormal) contains the eigenvectors of $XX^T$.
  $V \in \mathbb{R}^{d \times r}$ (orthonormal) contains the eigenvectors of $X^TX$.
  $\Sigma \in \mathbb{R}^{r \times r}$ (diagonal) contains their eigenvalues.

4

## SVD Review

$$n \left[ \begin{array}{c} \overset{k}{\overbrace{U_k \Sigma_k}} \end{array} \right]$$

- Every $X \in \mathbb{R}^{n \times d}$ can be written in its SVD as $U\Sigma V^T$.

- $U \in \mathbb{R}^{n \times r}$ (orthonormal) contains the eigenvectors of $XX^T$.
  $V \in \mathbb{R}^{d \times r}$ (orthonormal) contains the eigenvectors of $X^TX$.
  $\Sigma \in \mathbb{R}^{r \times r}$ (diagonal) contains their eigenvalues.

- $U_k U_k^T X = X V_k V_k^T = U_k \Sigma_k V_k^T = \underset{\text{B s.t. rank(B)} \leq k}{\arg \min} \|X - B\|_F$.

row
proj.        $d \times d$
proj

$$n \left[ \begin{array}{c} \overset{k}{U} \end{array} \right] \left[ \overset{k}{\Sigma_k} \right]^k \left[ \overset{l}{V_k^T} \right] k$$
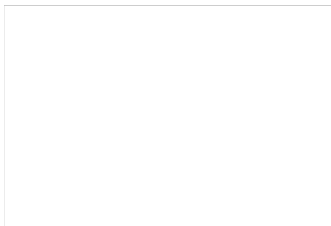
4

## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix).

## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem.
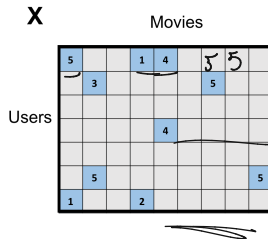
## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem.
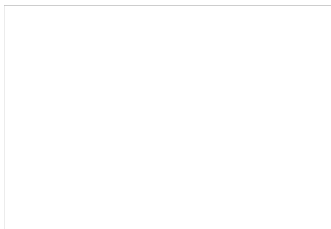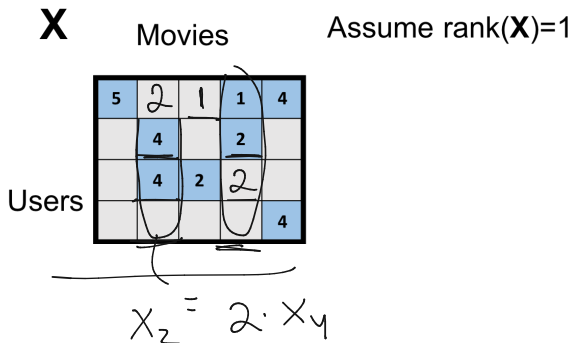
## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem.
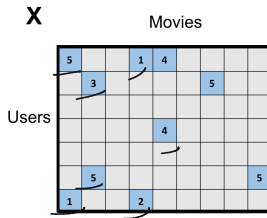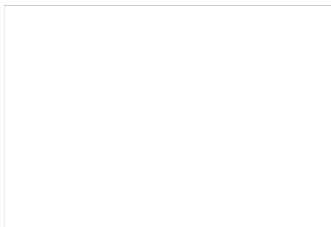


$$X_2 = 2 \cdot X_4$$

## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem.



**Solve:** $Y = \underset{B \text{ s.t. } \underline{\text{rank}(B)} \leq k}{\arg\min} \sum_{\text{observed } (j,k)} \left[ X_{j,k} - B_{j,k} \right]^2$

$$\underset{\substack{B: \\ \text{rank}(B) \leq k}}{\min} \| X - B \|_F^2$$

## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem.
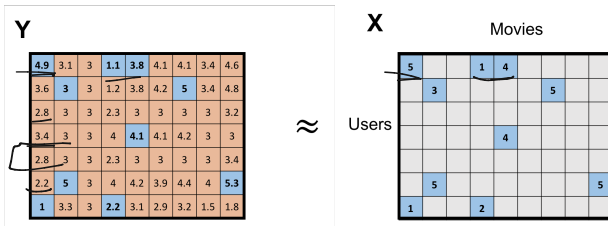


Solve: $Y = \underset{B \text{ s.t. } \mathrm{rank}(B) \leq k}{\arg\min} \sum_{\text{observed } (j,k)} \left[ X_{j,k} - B_{j,k} \right]^2$
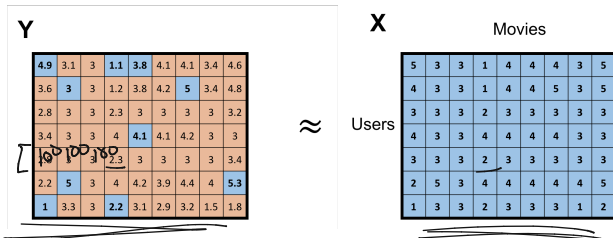
## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem.



**Solve:** $Y = \underset{B \text{ s.t. } \mathrm{rank}(B) \leq k}{\arg \min} \sum_{\text{observed } (j,k)} \left[ X_{j,k} - B_{j,k} \right]^2$
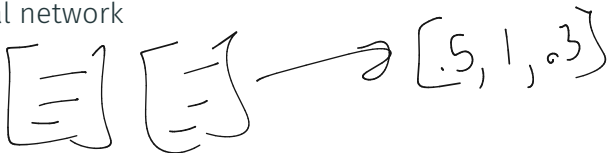
Under certain assumptions, can show that Y well approximates X on both the observed and (most importantly) unobserved entries.

Dimensionality reduction embeds *d*-dimensional vectors into *k* dimensions. But what about when you want to embed objects other than vectors?

- Documents (for topic-based search and classification)
- Words (to identify synonyms, translations, etc.)
- Nodes in a social network

## Entity Embeddings

Dimensionality reduction embeds $d$-dimensional vectors into $k$ dimensions. But what about when you want to embed objects other than vectors?

- Documents (for topic-based search and classification)
- Words (to identify synonyms, translations, etc.)
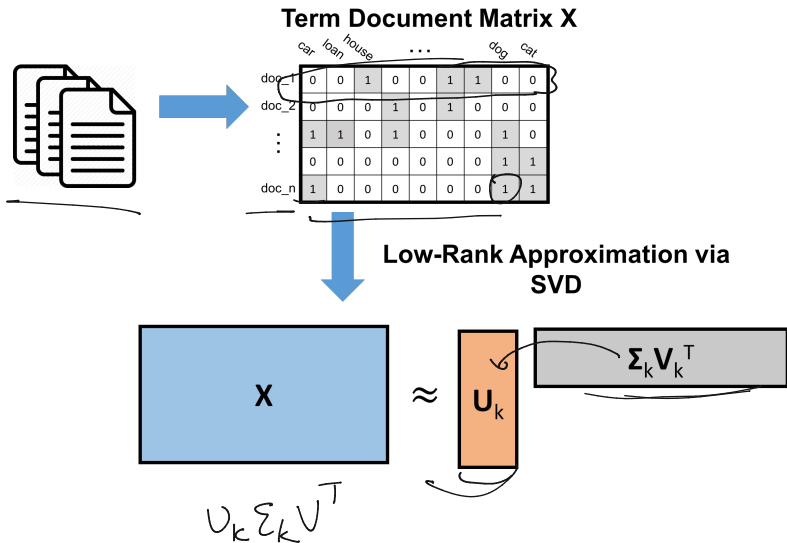- Nodes in a social network

**Classic Approach:** Convert each item into a (very) high-dimensional feature vector and then apply low-rank approximation.

**Term Document Matrix X**

**Low-Rank Approximation via SVD**

$$U_k \Sigma_k V^T$$

**Term Document Matrix X**

**Low-Rank Approximation via SVD**

$X \approx Y Z^T$

$UU^T X$

$|docs| \times |words|$

**X**

$|docs|$ $k$

$k$ **Z**$^T$ $(words)$

represents a word

$n \times d$

$XVV^T$

$XV_{n \times k}$

$k$ represents a doc

7

# Example: Latent Semantic Analysis



**Term Document Matrix X**

|         | car | loan | house | ... |   |   | dog | cat |   |
|---------|-----|------|-------|-----|---|---|-----|-----|---|
| doc_1   | 0   | 0    | 1     | 0   | 0 | 1 | 1   | 0   | 0 |
| doc_2   | 0   | 0    | 0     | 1   | 0 | 1 | 0   | 0   | 0 |
|   ⋮     | 1   | 1    | 0     | 1   | 0 | 0 | 0   | 1   | 0 |
|         | 0   | 0    | 0     | 0   | 0 | 0 | 0   | 1   | 1 |
| doc_n   | 1   | 0    | 0     | 0   | 0 | 0 | 0   | 1   | 1 |

**Low-Rank Approximation via SVD**
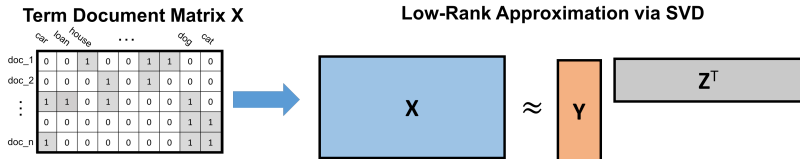
$$X \approx Y \; Z^{\mathsf{T}}$$

**Term Document Matrix X**

**Low-Rank Approximation via SVD**

- If the error $\|X - YZ^T\|_F$ is small, then on average,

$$X_{i,a} \approx (YZ^T)_{i,a} = \langle \vec{y}_i, \vec{z}_a \rangle.$$

**Term Document Matrix X**



**Low-Rank Approximation via SVD**

$$X \approx Y \, Z^T$$

- If the error $\|X - YZ^T\|_F$ is small, then on average,

$$X_{i,a} \approx (YZ^T)_{i,a} = \langle \vec{y}_i, \vec{z}_a \rangle.$$

- I.e., $\langle \vec{y}_i, \vec{z}_a \rangle \approx 1$ when $doc_i$ contains $word_a$.

# Example: Latent Semantic Analysis



**Term Document Matrix X**

| | car | loan | house | | ... | | dog | cat | |
|---|---|---|---|---|---|---|---|---|---|
| doc_1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 |
| doc_2 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |
| ⋮ | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| doc_n | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |

**Low-Rank Approximation via SVD**
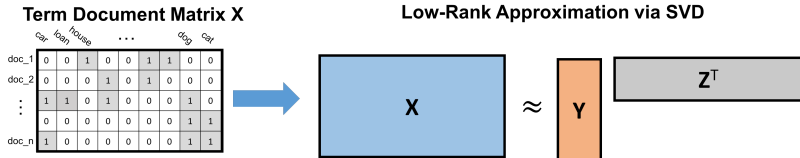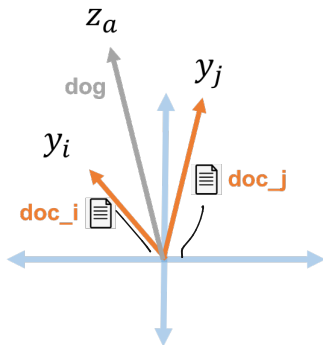
$$X \approx Y \, Z^T$$

- If the error $\|X - YZ^T\|_F$ is small, then on average,

$$X_{i,a} \approx (YZ^T)_{i,a} = \langle \vec{y}_i, \vec{z}_a \rangle.$$

- I.e., $\langle \vec{y}_i, \vec{z}_a \rangle \approx 1$ when $doc_i$ contains $word_a$.

- If $doc_i$ and $doc_j$ both contain $word_a$, $\langle \vec{y}_i, \vec{z}_a \rangle \approx \langle \vec{y}_j, \vec{z}_a \rangle \approx 1$.
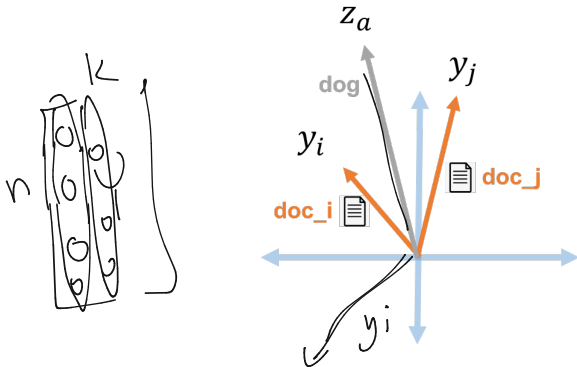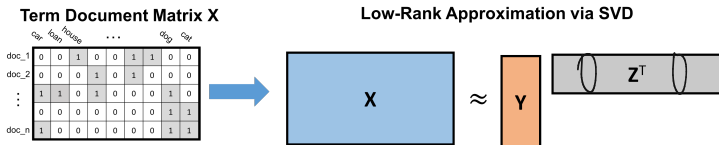
## Example: Latent Semantic Analysis

If $doc_i$ and $doc_j$ both contain $word_a$, $\langle \vec{y}_i, \vec{z}_a \rangle \approx \langle \vec{y}_j, \vec{z}_a \rangle \approx 1$

If $doc_i$ and $doc_j$ both contain $word_a$, $\langle \vec{y}_i, \vec{z}_a \rangle \approx \langle \vec{y}_j, \vec{z}_a \rangle \approx 1$
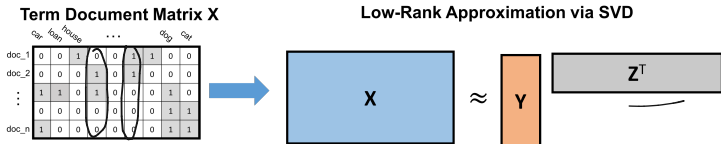


**Another View:** Each column of $Y$ represents a 'topic'. $\vec{y}_i(j)$ indicates how much $doc_i$ belongs to topic $j$. $\vec{z}_a(j)$ indicates how much $word_a$ associates with that topic.

# Example: Latent Semantic Analysis



**Term Document Matrix X**      **Low-Rank Approximation via SVD**

$$X \approx Y \, Z^T$$

- Just like with documents, $\vec{z}_a$ and $\vec{z}_b$ will tend to have high dot product if $word_a$ and $word_b$ appear in many of the same documents.

**Term Document Matrix X**

**Low-Rank Approximation via SVD**

- Just like with documents, $\vec{z}_a$ and $\vec{z}_b$ will tend to have high dot product if $word_a$ and $word_b$ appear in many of the same documents.

- In an SVD decomposition we set $Z^T = \Sigma_k V_K^T$.

- The columns of $V_k$ are equivalently: the top $k$ eigenvectors of $X^T X$.

**Term Document Matrix X**
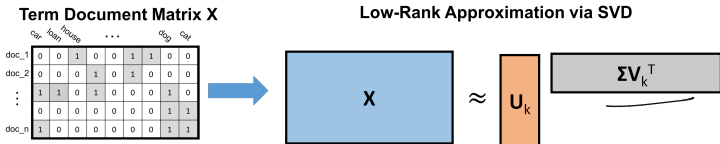
**Low-Rank Approximation via SVD**

- Just like with documents, $\vec{z}_a$ and $\vec{z}_b$ will tend to have high dot product if $word_a$ and $word_b$ appear in many of the same documents.

- In an SVD decomposition we set $Z^T = \Sigma_k V_K^T$.

- The columns of $V_k$ are equivalently: the top $k$ eigenvectors of $X^TX$.

- **Claim:** $ZZ^T$ is the best rank-$k$ approximation of $X^TX$. I.e.,
  $$\arg\min_{\text{rank}-k \; B} \|X^TX - B\|_F$$

SVD $X^TX = V\Sigma U^T U \Sigma V^T = V\Sigma^2 V^T$

$V_k \Sigma_k^2 V_k^T$

## Example: Word Embedding

LSA gives a way of embedding words into $k$-dimensional space.

- Embedding is via low-rank approximation of $X^T X$: where $(X^T X)_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.

## Example: Word Embedding

LSA gives a way of embedding words into $k$-dimensional space.

- Embedding is via low-rank approximation of $X^T X$: where $(X^T X)_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.

- Think about $X^T X$ as a similarity matrix (gram matrix, kernel matrix) with entry $(a, b)$ being the similarity between $word_a$ and $word_b$.

## Example: Word Embedding

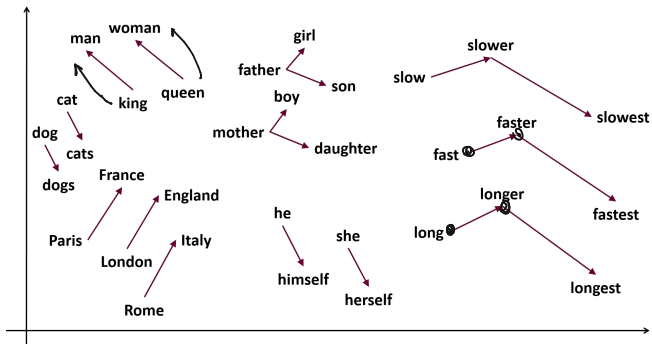LSA gives a way of embedding words into $k$-dimensional space.

- Embedding is via low-rank approximation of $X^T X$: where $(X^T X)_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.
- Think about $X^T X$ as a similarity matrix (gram matrix, kernel matrix) with entry $(a, b)$ being the similarity between $word_a$ and $word_b$.
- Many ways to measure similarity: number of sentences both occur in, number of times both appear in the same window of $w$ words, in similar positions of documents in different languages, etc.
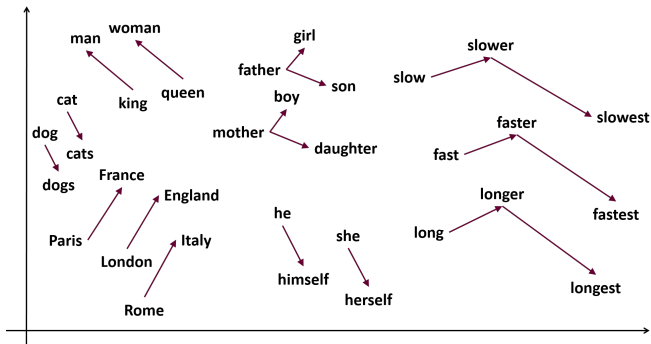
## Example: Word Embedding

LSA gives a way of embedding words into $k$-dimensional space.

- Embedding is via low-rank approximation of $X^T X$: where $(X^T X)_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.
- Think about $X^T X$ as a similarity matrix (gram matrix, kernel matrix) with entry $(a, b)$ being the similarity between $word_a$ and $word_b$.
- Many ways to measure similarity: number of sentences both occur in, number of times both appear in the same window of $w$ words, in similar positions of documents in different languages, etc.
- Replacing $X^T X$ with these different metrics (sometimes appropriately transformed) leads to popular word embedding algorithms: word2vec, GloVe, fastText, etc.
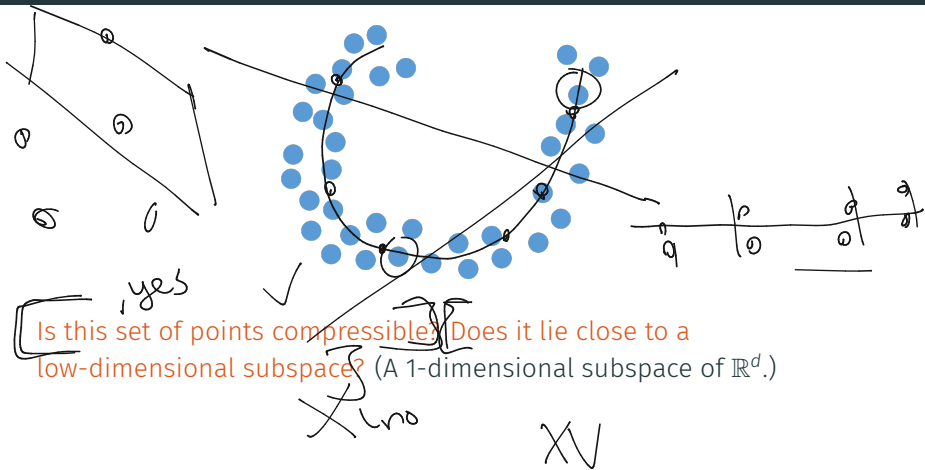
## Example: Word Embedding

## Example: Word Embedding



**Note:** word2vec is typically described as a neural-network method, but can be viewed as just a low-rank approximation of a specific similarity matrix. *Neural word embedding as implicit matrix factorization,* Levy and Goldberg.

# Non-Linear Dimensionality Reduction



Is this set of points compressible? Does it lie close to a low-dimensional subspace? (A 1-dimensional subspace of $\mathbb{R}^d$.)
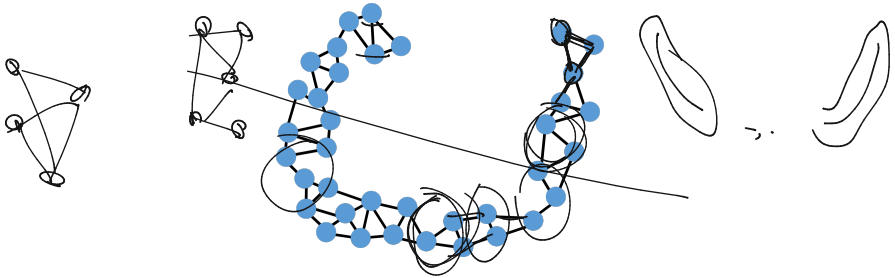
# Non-Linear Dimensionality Reduction



Is this set of points compressible? Does it lie close to a low-dimensional subspace? (A 1-dimensional subspace of $\mathbb{R}^d$.)

# Non-Linear Dimensionality Reduction



Is this set of points compressible? Does it lie close to a low-dimensional subspace? (A 1-dimensional subspace of $\mathbb{R}^d$.)

A common way of automatically identifying this non-linear structure is to connect data points in a graph. E.g., a $r$-nearest neighbor graph.

- Connect items to similar items, possibly with higher weight edges when they are more similar.
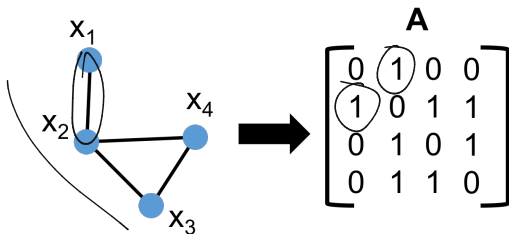
## Linear Algebraic Representation of a Graph

Once we have connected $n$ data points $x_1, \ldots, x_n$ into a graph, we can represent that graph by its (weighted) adjacency matrix.

$A \in \mathbb{R}^{n \times n}$ with $A_{i,j} = $ edge weight between nodes $i$ and $j$

Once we have connected $n$ data points $x_1, \ldots, x_n$ into a graph, we can represent that graph by its (weighted) adjacency matrix.

$\mathbf{A} \in \underline{\mathbb{R}^{n \times n} \text{ with } \mathbf{A}_{i,j} =}$ edge weight between nodes $i$ and $j$

## Linear Algebraic Representation of a Graph

Once we have connected $n$ data points $x_1, \ldots, x_n$ into a graph, we can represent that graph by its (weighted) adjacency matrix.
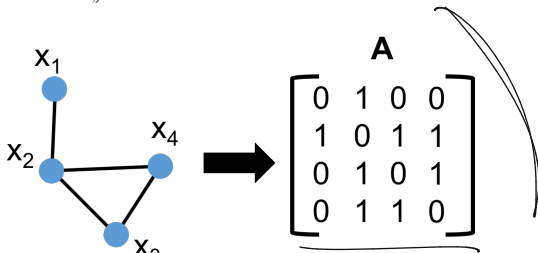
$A \in \mathbb{R}^{n \times n}$ with $A_{i,j} =$ edge weight between nodes $i$ and $j$



$$\left(X^T X\right)_{ab} = \text{it does contain both } a, b$$

In LSA example, when $X$ is the term-document matrix, $X^T X$ is like an adjacency matrix, where $word_a$ and $word_b$ are connected if they appear in at least 1 document together (edge weight is $\#$ documents they appear in together).

How do we compute an optimal low-rank approximation of $\mathbf{A}$?

- Project onto the top $k$ eigenvectors of $\mathbf{A}^T\mathbf{A} = \mathbf{A}^2$. These are just the eigenvectors of $\mathbf{A}$.

$$Ax = \lambda x$$

$$A^2 x = A(Ax) = A(\lambda x)$$
$$= \lambda(\lambda x)$$
$$= \lambda^2 x$$

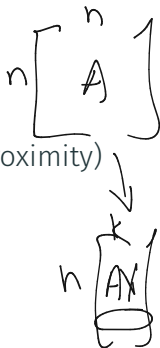$$A(\lambda x) = \lambda Ax = \lambda \cdot \lambda x$$
$$= \lambda^2 x$$

# Adjacency Matrix Eigenvectors
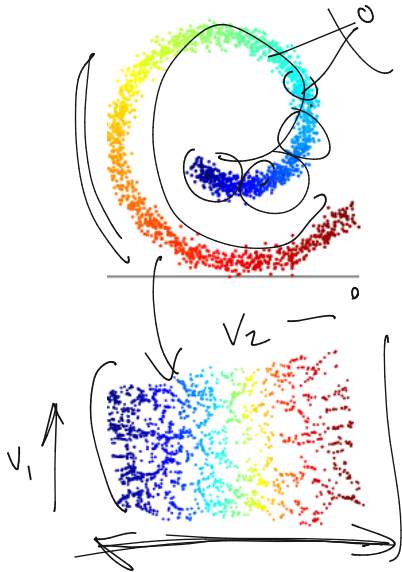
$$A^T A = A \cdot A = A^2 \qquad A^T = A$$

How do we compute an optimal low-rank approximation of $A$?

- Project onto the top $k$ eigenvectors of $A^T A = A^2$. These are just the eigenvectors of $A$.

- $A \approx AVV^T$. The rows of $AV$ can be thought of as 'embeddings' for the vertices.

- Similar vertices (close with regards to graph proximity) should have similar embeddings.
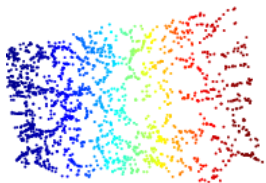
# Spectral Embedding



**Step 1:** Produce a nearest neighbor graph based on your input data in $\mathbb{R}^d$.

**Step 2:** Apply low-rank approximation to the graph adjacency matrix to produce embeddings in $\mathbb{R}^k$.

**Step 3:** Work with the data in the embedded space. Where distances represent distances in your original 'non-linear space.'

# Spectral Embedding



What other methods do you know for embedding or representing data points with non-linear structure?

Questions?