

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Spring 2020.

Lecture 17

- Problem Set 2 was released this weekend. Due Monday 4/13.
- See Piazza (and email from college) for clarification on P/F policy.

Last Few Classes: Low-Rank Approximation and PCA

- Compress data that lies close to a k -dimensional subspace.
- Equivalent to finding a low-rank approximation of the data matrix \mathbf{X} : $\mathbf{X} \approx \mathbf{X}\mathbf{V}\mathbf{V}^T$ for orthonormal $\mathbf{V} \in \mathbb{R}^{d \times k}$.
- Optimal solution via PCA (eigendecomposition of $\mathbf{X}^T\mathbf{X}$ or equivalently, SVD of \mathbf{X}).
- Singular vectors of \mathbf{X} are the eigenvectors of $\mathbf{X}\mathbf{X}^T$ and $\mathbf{X}^T\mathbf{X}$. Singular values squared are the eigenvalues.

Last Few Classes: Low-Rank Approximation and PCA

- Compress data that lies close to a k -dimensional subspace.
- Equivalent to finding a low-rank approximation of the data matrix \mathbf{X} : $\mathbf{X} \approx \mathbf{X}\mathbf{V}\mathbf{V}^T$ for orthonormal $\mathbf{V} \in \mathbb{R}^{d \times k}$.
- Optimal solution via PCA (eigendecomposition of $\mathbf{X}^T\mathbf{X}$ or equivalently, SVD of \mathbf{X}).
- Singular vectors of \mathbf{X} are the eigenvectors of $\mathbf{X}\mathbf{X}^T$ and $\mathbf{X}^T\mathbf{X}$. Singular values squared are the eigenvalues.

This Class: Applications of low-rank approx. beyond compression.

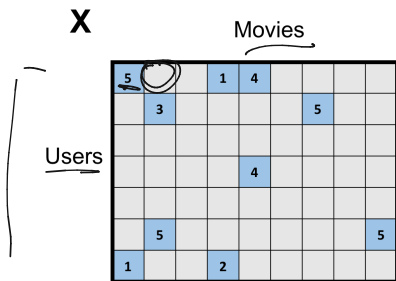
- Matrix completion and collaborative filtering
- Entity embeddings (word embeddings, node embeddings, etc.)
- Low-rank approximation for **non-linear dimensionality reduction**.
- Spectral graph theory, spectral clustering.

MATRIX COMPLETION

Consider a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ 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 $\mathbf{X} \in \mathbb{R}^{n \times d}$ 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

Movies

5		1	4		
	3				5
			4		
	5				5
1		2			

Users

observed all of X
↓
eigen bump
 $X^T X$

$$\min \|X - B\|_F^2$$

Solve: $Y = \arg \min_{\text{rank } -k \text{ } B} \sum_{\text{observed } (j,k)} [X_{j,k} - B_{j,k}]^2$

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.

Y

Movies

4.9	3.1	3	1.1	3.8	4.1	4.1	3.4	4.6
3.6	3	3	1.2	3.8	4.2	5	3.4	4.8
2.8	3	3	2.3	3	3	3	3	3.2
3.4	3	3	4	4.1	4.1	4.2	3	3
2.8	3	3	2.3	3	3	3	3	3.4
2.2	5	3	4	4.2	3.9	4.4	4	5.3
1	3.3	3	2.2	3.1	2.9	3.2	1.5	1.8

Users

regression
using low rank
approx.

Solve: $Y = \arg \min_{\text{rank}-k B} \sum_{\text{observed } (j,k)} [X_{j,k} - B_{j,k}]^2$

$B = [U][W]$

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 $\mathcal{X} = d'$ dimensions. But what about when you want to embed objects other than vectors?

Dimensionality reduction embeds d -dimensional vectors into d' 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

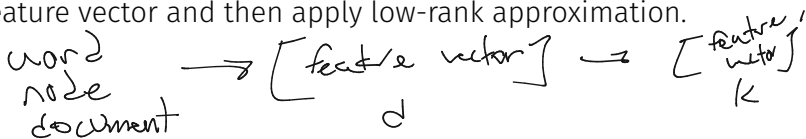
$$X \in \mathbb{R}^{n \times d} \approx \underbrace{XW^T}_{\text{compression}} \in \mathbb{R}^{n \times d}$$

$\|X - XWV^T\|_F^2$ $XV \in \mathbb{R}^{n \times k}$

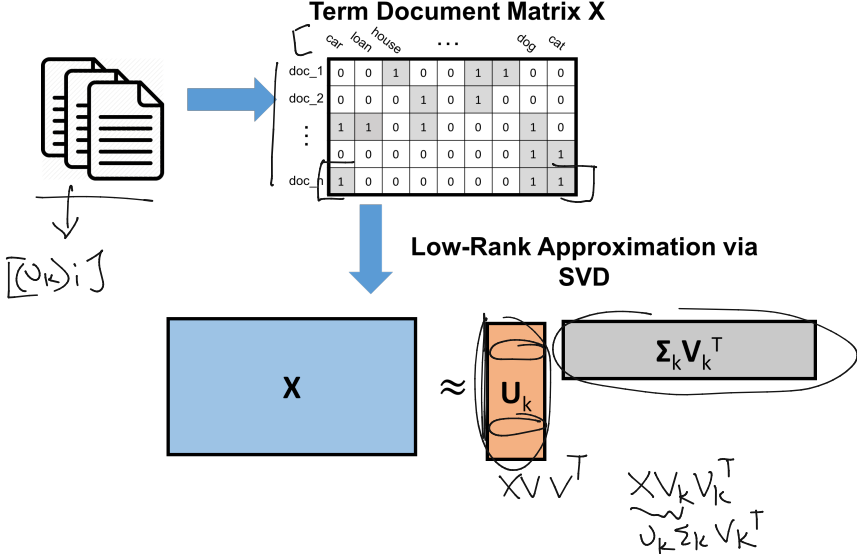
Dimensionality reduction embeds d -dimensional vectors into d' dimensions. But what about when you want to embed objects other than vectors?

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

Usual Approach: Convert each item into a high-dimensional feature vector and then apply low-rank approximation.



EXAMPLE: LATENT SEMANTIC ANALYSIS



EXAMPLE: LATENT SEMANTIC ANALYSIS

$$G_i(x) = x_i (x^T x) = \lambda_i (x x^T)$$

$$X = [U] \begin{bmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \ddots \end{bmatrix} [V]^T$$

Term Document Matrix X

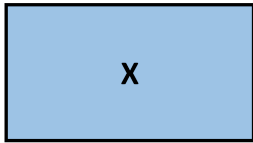


	car	loan	house	...	dog	cat			
doc_1	0	0	1	0	0	1	1	0	0
doc_2	0	0	0	1	0	1	0	0	0
⋮									
doc_n	1	0	0	0	0	0	0	1	1

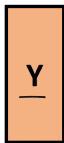
$$U \Sigma_k V_k^T$$

$$[U_k] [\Sigma_k] [V_k^T] = X_k$$

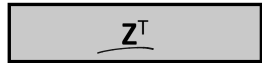
Low-Rank Approximation via SVD



≈



U_k



$\Sigma_k V_k^T$



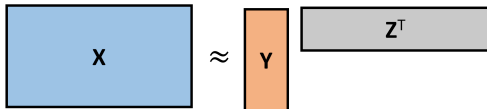
EXAMPLE: LATENT SEMANTIC ANALYSIS

Term Document Matrix X

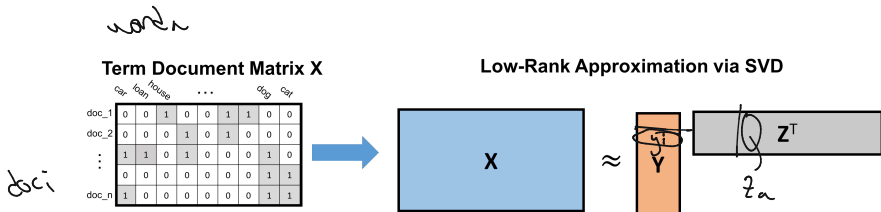
	car	loan	house	...	dog	cat			
doc_1	0	0	1	0	0	1	1	0	0
doc_2	0	0	0	1	0	1	0	0	0
⋮	1	1	0	1	0	0	0	1	0
⋮	0	0	0	0	0	0	0	1	1
doc_n	1	0	0	0	0	0	0	1	1



Low-Rank Approximation via SVD



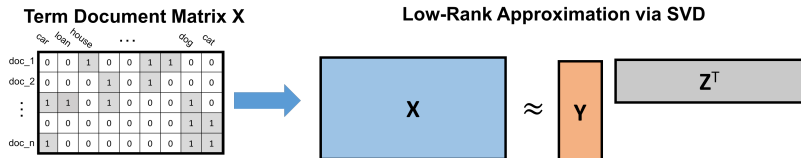
EXAMPLE: LATENT SEMANTIC ANALYSIS



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

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

EXAMPLE: LATENT SEMANTIC ANALYSIS

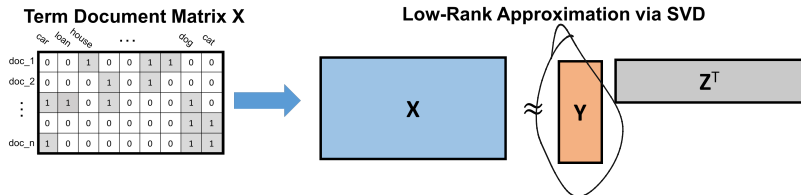


- If the error $\|X - YZ^T\|_F$ is small, then on average,

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

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

EXAMPLE: LATENT SEMANTIC ANALYSIS



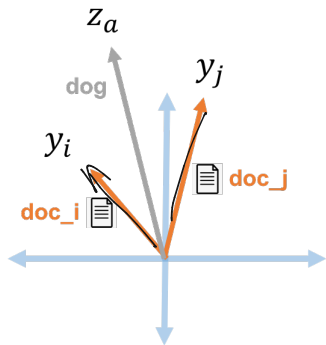
- If the error $\|X - YZ^T\|_F$ is small, then on average,

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

- 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

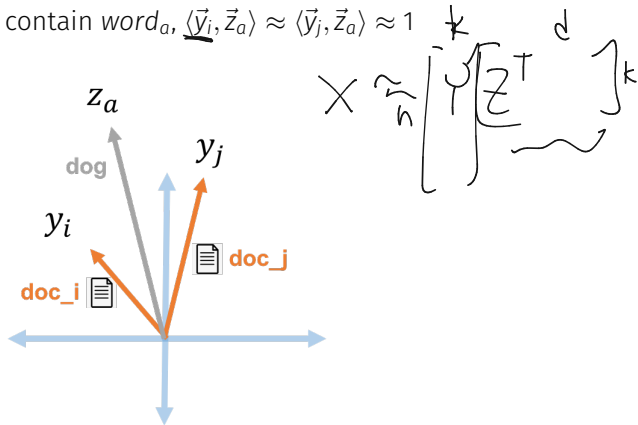
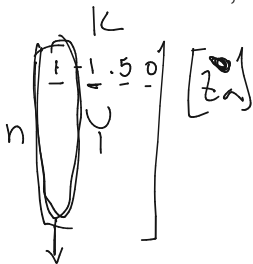
If doc_i and doc_j both contain $word_a$, $\langle \vec{y}_i, \vec{z}_a \rangle \approx \langle \vec{y}_j, \vec{z}_a \rangle \approx 1$



$$\|X - YZ^T\|_F^2$$

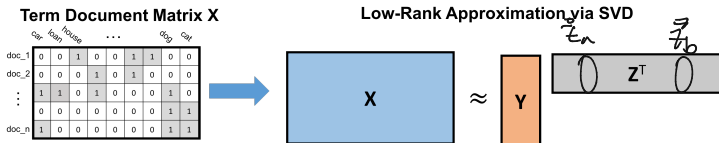
EXAMPLE: LATENT SEMANTIC ANALYSIS

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



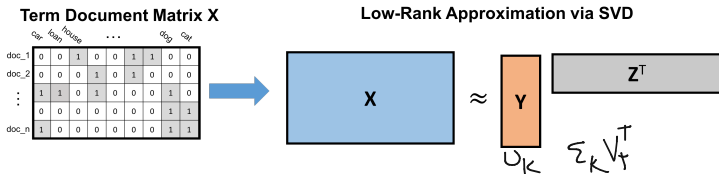
Another View: Each column of Y represents a 'topic'. $\vec{y}_i(j)$ indicates how much doc_i belongs to topic j . $\vec{z}_a(j)$ indicates how much $word_a$ associates with that topic.

EXAMPLE: LATENT SEMANTIC ANALYSIS



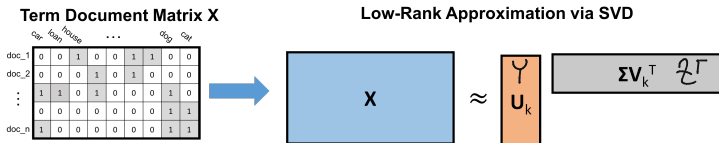
- Just like with documents, \vec{z}_a and \vec{z}_b will tend to have high dot product if $word_a$ and $word_b$ appear in many of the same documents.

EXAMPLE: LATENT SEMANTIC ANALYSIS



- 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 $Z = \sum_k \mathbf{v}_k^T$.
- The columns of \mathbf{v}_k are equivalently: the top k eigenvectors of $\underline{X^T X}$.

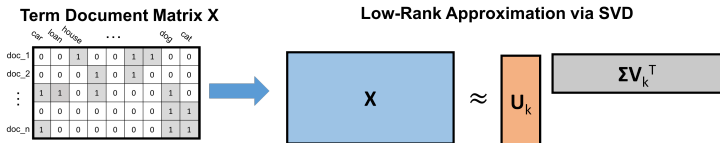
EXAMPLE: LATENT SEMANTIC ANALYSIS



- 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 $X^T X$.
The eigendecomposition of $X^T X$ is $X^T X = V \Sigma^2 V^T$.

$$X = U \Sigma V^T \quad X^T X = V \Sigma U^T U \Sigma V^T = \underline{V \Sigma^2 V^T}$$

EXAMPLE: LATENT SEMANTIC ANALYSIS



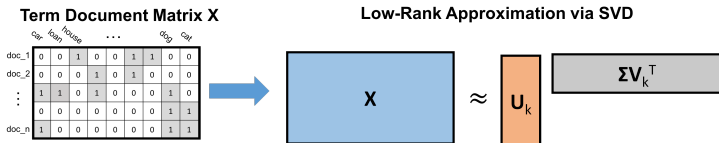
- Just like with documents, \vec{z}_a and \vec{z}_b will tend to have high dot product if $word_a$ and $word_b$ appear in many of the same documents.
- In an SVD decomposition we set $\underline{Z^T} = \underline{\Sigma_k V_k^T}$.
- The columns of V_k are equivalently: the top k eigenvectors of $X^T X$. The eigendecomposition of $X^T X$ is $X^T X = \underline{V \Sigma^2 V^T}$.
- What is the best rank- k approximation of $X^T X$? i.e.

$$\arg \min_{\text{rank} = k} \|X^T X - B\|_F$$

$Z Z^T$ is best LRA of $X^T X$

$$\begin{bmatrix} V_k \\ \Sigma_k \\ V_k^T \end{bmatrix} = \begin{bmatrix} Z Z^T \end{bmatrix}$$

EXAMPLE: LATENT SEMANTIC ANALYSIS



- 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 $Z = \mathbf{\Sigma}_k \mathbf{V}_k^T$.
- The columns of \mathbf{V}_k are equivalently: the top k eigenvectors of $\mathbf{X}^T \mathbf{X}$. The eigendecomposition of $\mathbf{X}^T \mathbf{X}$ is $\mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}^T$.
- **What is the best rank- k approximation of $\mathbf{X}^T \mathbf{X}$?** I.e.
$$\arg \min_{\text{rank} = k} \mathbf{B} \|\mathbf{X}^T \mathbf{X} - \mathbf{B}\|_F$$
- $\mathbf{X} \mathbf{X}^T = \mathbf{V}_k \mathbf{\Sigma}_k^2 \mathbf{V}_k^T = \mathbf{Z} \mathbf{Z}^T$.

EXAMPLE: WORD EMBEDDING

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

[2]

- 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.

doc $\begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix}$ words

word $[1 \ 0 \ 1 \ 0]$ $\begin{bmatrix} \cdot \\ 1 \\ \cdot \\ \cdot \end{bmatrix} = 2$

word $\begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} X^T$

$X^T X$

EXAMPLE: WORD EMBEDDING

$$X \approx [U]C[Z^T]$$

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.

$$ZZ^T = \operatorname{argmin} \|X^T X - B\|_F^2$$

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

- Think about $\mathbf{X}^T\mathbf{X}$ 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.

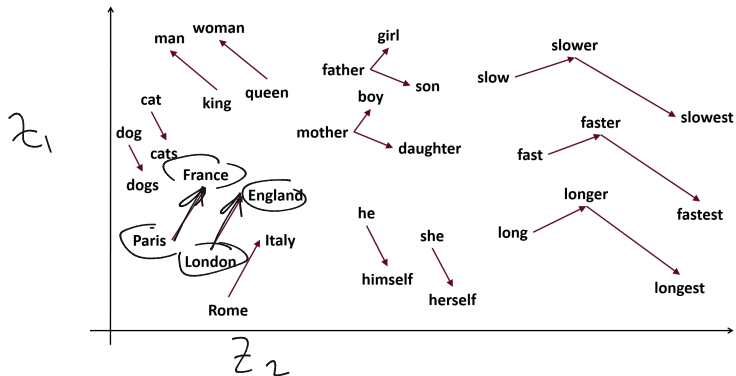
$$\begin{matrix} \uparrow \\ \left[\begin{matrix} S \\ \end{matrix} \right] \end{matrix} \begin{matrix} \downarrow \\ \left(\begin{matrix} d \\ \end{matrix} \right) \end{matrix}$$

EXAMPLE: WORD EMBEDDING

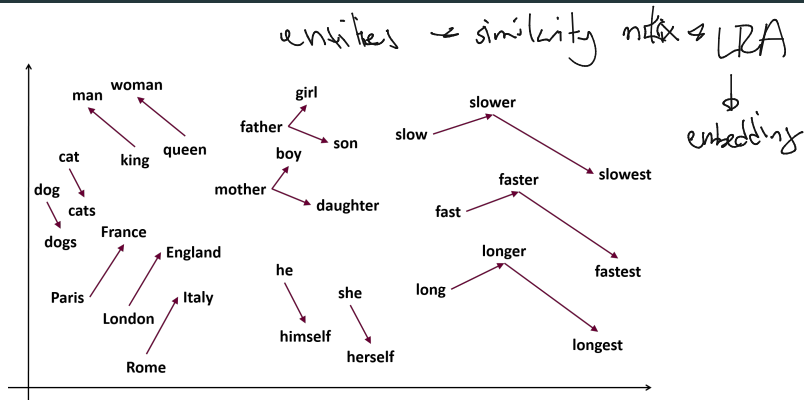
LSA gives a way of embedding words into k -dimensional space.

- Embedding is via low-rank approximation of $\mathbf{X}^T\mathbf{X}$: where $(\mathbf{X}^T\mathbf{X})_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.
- Think about $\mathbf{X}^T\mathbf{X}$ 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 $\mathbf{X}^T\mathbf{X}$ 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 it is really just low-rank approximation of a specific similarity matrix. (Neural word embedding as implicit matrix factorization, Levy and Goldberg)

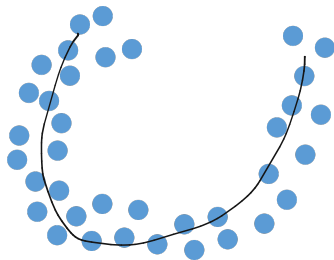
A common way of encoding similarity is via a graph. E.g., a k -nearest neighbor graph.

- Connect items to similar items, possibly with higher weight edges when they are more similar.

SIMILARITY VIA GRAPHS

A common way of encoding similarity is via a graph. E.g., a k -nearest neighbor graph.

- Connect items to similar items, possibly with higher weight edges when they are more similar.

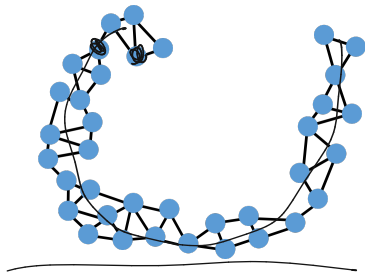


SIMILARITY VIA GRAPHS

A common way of encoding similarity is via a graph. E.g., a k -nearest neighbor graph.

- Connect items to similar items, possibly with higher weight edges when they are more similar.

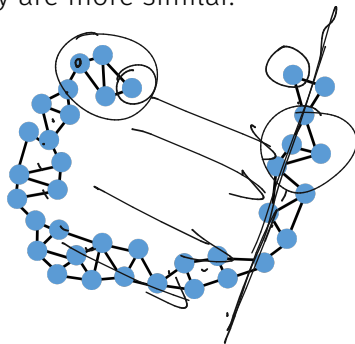
Non-linear
compression



SIMILARITY VIA GRAPHS

A common way of encoding similarity is via a graph. E.g., a k -nearest neighbor graph.

- Connect items to similar items, possibly with higher weight edges when they are more similar.

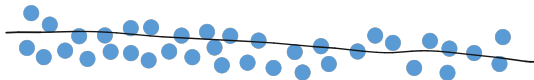


Is this set of points compressible? Does it lie close to a low-dimensional subspace?

SIMILARITY VIA GRAPHS

A common way of encoding similarity is via a graph. E.g., a k -nearest neighbor graph.

- Connect items to similar items, possibly with higher weight edges when they are more similar.



Is this set of points compressible? Does it lie close to a low-dimensional subspace?

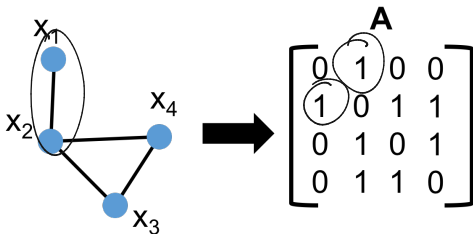
Once we have connected n data points x_1, \dots, x_n into a graph, we can represent that graph by its (weighted) adjacency matrix.

$\mathbf{A} \in \mathbb{R}^{n \times n}$ 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, \dots, x_n into a graph, we can represent that graph by its (weighted) adjacency matrix.

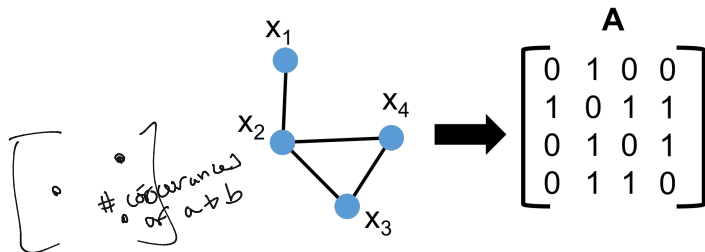
$A \in \mathbb{R}^{n \times n}$ with $A_{i,j}$ = edge weight between nodes i and j



LINEAR ALGEBRAIC REPRESENTATION OF A GRAPH

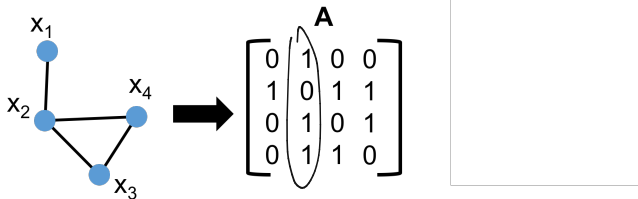
Once we have connected n data points x_1, \dots, x_n into a graph, we can represent that graph by its (weighted) adjacency matrix.

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



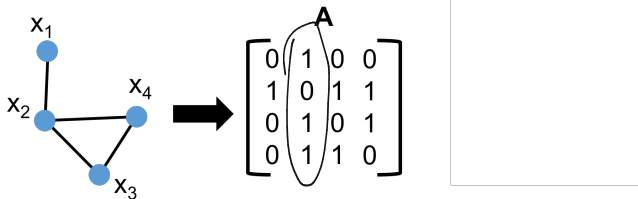
In LSA example, when X is the term-document matrix, $X^T X$ is like an adjacency matrix, where $word_a$ and $word_b$ are connected if they appear in at least 1 document together (edge weight is # documents they appear in together).

NORMALIZED ADJACENCY MATRIX



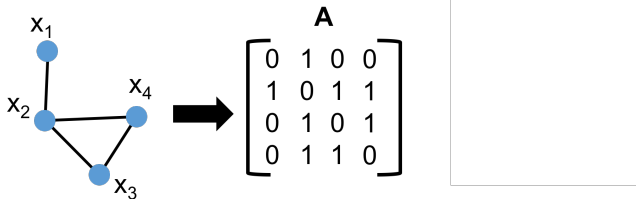
What is the sum of entries in the i^{th} column of A ?

NORMALIZED ADJACENCY MATRIX



What is the sum of entries in the i^{th} column of A ? The (weighted) degree of vertex i .

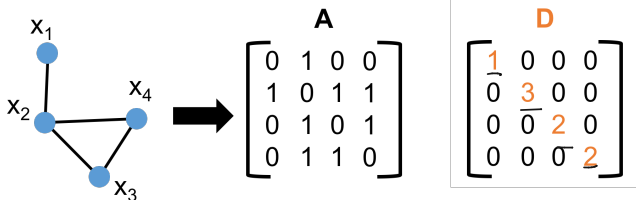
NORMALIZED ADJACENCY MATRIX



What is the sum of entries in the i^{th} column of A ? The (weighted) degree of vertex i .

Often, A is normalized as $\bar{A} = \underline{D^{-1/2}AD^{-1/2}}$ where D is the degree matrix.

NORMALIZED ADJACENCY MATRIX



What is the sum of entries in the i^{th} column of A ? The (weighted) degree of vertex i .

Often, A is normalized as $\bar{A} = D^{-1/2}AD^{-1/2}$ where D is the degree matrix.

NORMALIZED ADJACENCY MATRIX

$$\bar{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$$

$$\begin{bmatrix} 0 & \frac{1}{\sqrt{3}} & 0 & 0 \\ \frac{1}{\sqrt{3}} & 0 & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{6}} \\ 0 & \frac{1}{\sqrt{6}} & 0 & \frac{1}{2} \\ 0 & \frac{1}{\sqrt{3}} & \frac{1}{2} & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & \frac{1}{3} & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & \frac{1}{3} & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & \frac{1}{2} \end{bmatrix}$$

What is the sum of entries in the i^{th} column of A ? The (weighted) degree of vertex i .

Often, \mathbf{A} is normalized as $\bar{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ where \mathbf{D} is the degree matrix.

NORMALIZED ADJACENCY MATRIX

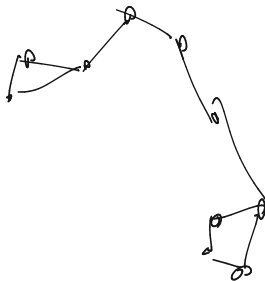
$$\bar{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$$
$$\begin{bmatrix} 0 & \frac{1}{\sqrt{3}} & 0 & 0 \\ \frac{1}{\sqrt{3}} & 0 & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{6}} \\ 0 & \frac{1}{\sqrt{6}} & 0 & \frac{1}{2} \\ 0 & \frac{1}{\sqrt{3}} & \frac{1}{2} & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & \frac{1}{3} & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & \frac{1}{3} & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & \frac{1}{2} \end{bmatrix}$$

What is the sum of entries in the i^{th} column of \mathbf{A} ? The (weighted) degree of vertex i .

Often, \mathbf{A} is normalized as $\bar{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ where \mathbf{D} is the degree matrix.

Spectral graph theory is the field of representing graphs as matrices and applying linear algebraic techniques.

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



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

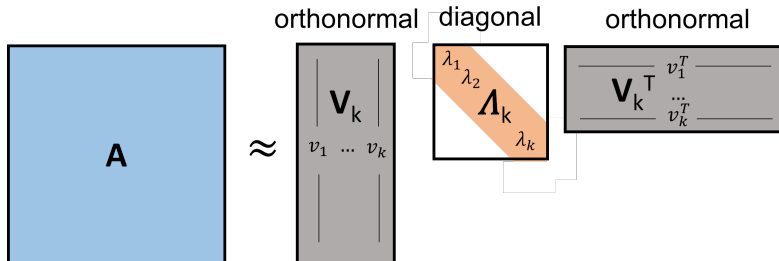
- Project onto the top k eigenvectors of $\mathbf{A}^T\mathbf{A} = \mathbf{A}^2$.

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

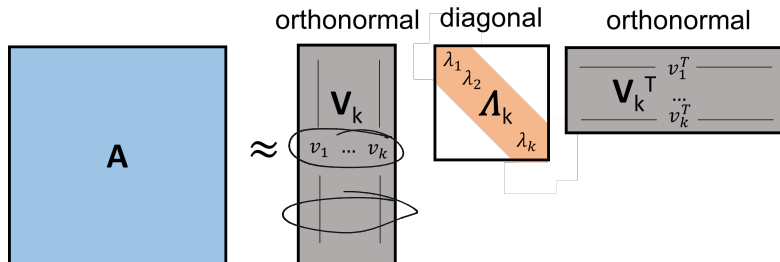
- Project onto the top k eigenvectors of $\mathbf{A}^T\mathbf{A} = \mathbf{A}^2$. These are just the eigenvectors of \mathbf{A} .

$$\begin{aligned} \mathbf{A} &= \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T \\ \mathbf{A}^T\mathbf{A} &= \mathbf{A}^2 = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T\mathbf{V}\mathbf{\Lambda}\mathbf{V}^T \\ &= \mathbf{V}\mathbf{\Lambda}^2\mathbf{V}^T \end{aligned}$$

ADJACENCY MATRIX EIGENVECTORS



ADJACENCY MATRIX EIGENVECTORS



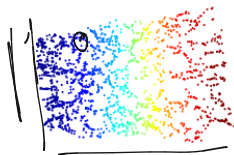
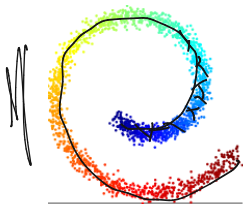
- Similar vertices (close with regards to graph proximity) should have similar embeddings. I.e., $\mathbf{V}_k(i)$ should be similar to $\mathbf{V}_k(j)$.

"jelly roll"

$$\begin{bmatrix} A \end{bmatrix}$$



$$\begin{bmatrix} z \end{bmatrix}_2 \begin{bmatrix} z^T \end{bmatrix}_2$$



Summary:

LRA used beyond
compression of vectors

- matrix completion
- entity embedding
- nonlinear dimension