

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Fall 2021.

Lecture 18

- Problem Set 3 is due Monday at 11:59pm.
- No quiz due Monday.

Last Class

- The Singular Value Decomposition (SVD) and its connection to eigendecomposition of $\mathbf{X}^T\mathbf{X}$ and $\mathbf{X}\mathbf{X}^T$, and low-rank approximation.
- Low-rank matrix completion (predicting missing measurements using low-rank structure).

This Class: More applications of Low-Rank Approximation

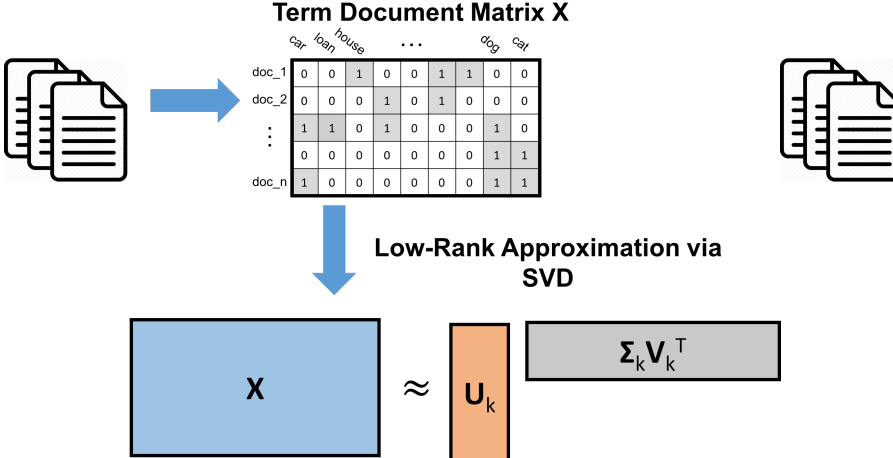
- Entity embeddings.
- Low-rank approximation for non-linear dimensionality reduction.
- Eigendecomposition to partition graphs into clusters.

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

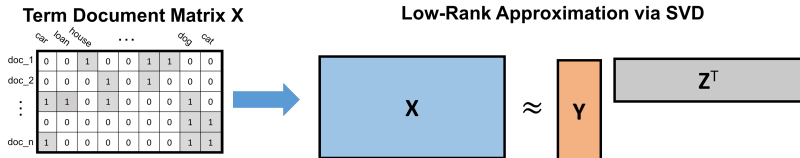
- Documents (for topic-based search and classification)
- Words (to identify synonyms, translations, etc.)
- Nodes in a social network

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

EXAMPLE: LATENT SEMANTIC ANALYSIS



EXAMPLE: LATENT SEMANTIC ANALYSIS



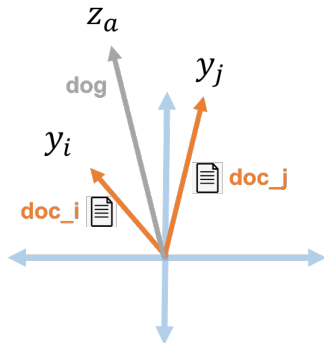
- If the error $\|X - YZ^T\|_F$ is small, then on average,

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

- 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 $word_a$, $\langle \vec{y}_i, \vec{z}_a \rangle \approx \langle \vec{y}_j, \vec{z}_a \rangle \approx 1$.

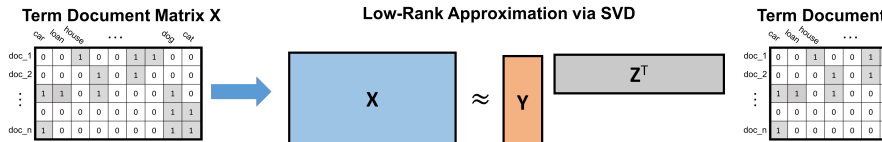
EXAMPLE: LATENT SEMANTIC ANALYSIS

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



Another View: Each column of \mathbf{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



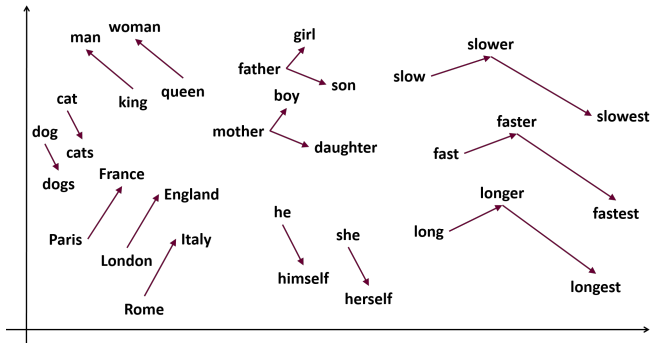
- Just like with documents, \vec{z}_a and \vec{z}_b will tend to have high dot product if $word_a$ and $word_b$ appear in many of the same documents.
- In an SVD decomposition we set $Z^T = \sum_k V_k^T$.
- The columns of V_k are equivalently: the top k eigenvectors of $X^T X$.
- **Claim:** ZZ^T is the best rank- k approximation of $X^T X$. I.e.,
$$\arg \min_{\text{rank} = k, B} \|X^T X - B\|_F$$

EXAMPLE: WORD EMBEDDING

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

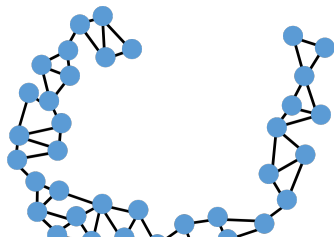
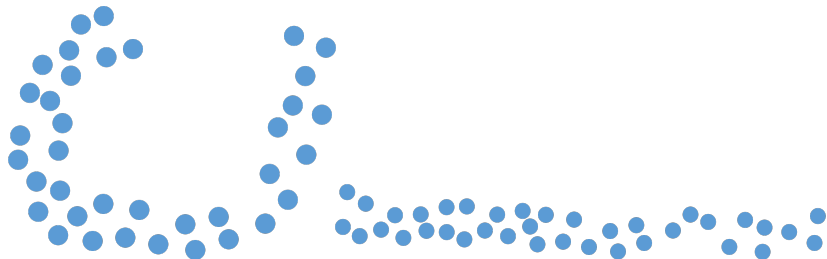
- Embedding is via low-rank approximation of $\mathbf{X}^T\mathbf{X}$: where $(\mathbf{X}^T\mathbf{X})_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.
- Think about $\mathbf{X}^T\mathbf{X}$ as a **similarity matrix** (gram matrix, kernel matrix) with entry (a, b) being the 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



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.

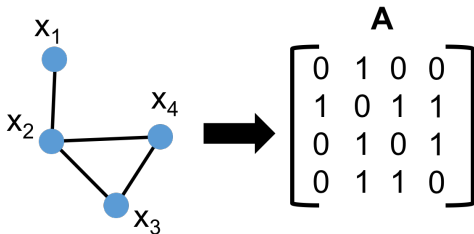
NON-LINEAR DIMENSIONALITY REDUCTION



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.

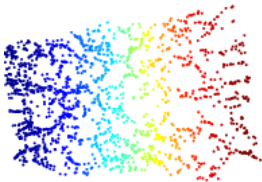
$\mathbf{A} \in \mathbb{R}^{n \times n}$ with $\mathbf{A}_{i,j}$ = edge weight between nodes i and j



In LSA example, when \mathbf{X} is the term-document matrix, $\mathbf{X}^T\mathbf{X}$ 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} ?

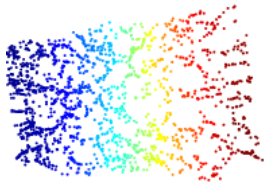
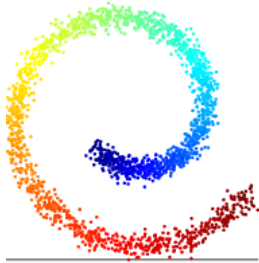
- Project onto the top k eigenvectors of $\mathbf{A}^T\mathbf{A} = \mathbf{A}^2$. These are just the eigenvectors of \mathbf{A} .
- $\mathbf{A} \approx \mathbf{A}\mathbf{V}\mathbf{V}^T$. The rows of $\mathbf{A}\mathbf{V}$ can be thought of as ‘embeddings’ for the vertices.
- Similar vertices (close with regards to graph proximity) should have similar embeddings.



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.'



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

Questions?