

Matrix approximation

COMS 4771 Fall 2019

Overview

- ▶ Compression using k -means
- ▶ Compression using SVD
- ▶ Power iteration
- ▶ Latent semantic analysis
- ▶ Matrix completion

0 / 15

1 / 15

Encoder/decoder interpretation of k -means clustering

- ▶ Given data $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, we learn a lossy encoding/decoding (compression) scheme
- ▶ Scheme is represented by $\mathbf{c}_1, \dots, \mathbf{c}_k \in \mathbb{R}^d$
- ▶ Encoder: $\mathbf{x} \mapsto \arg \min_j \|\mathbf{x} - \mathbf{c}_j\|_2^2$
 - ▶ Need $\log_2 k$ bits to represent encoding of \mathbf{x} !
- ▶ Decoder: $j \mapsto \mathbf{c}_j$

2 / 15

Matrix view

- ▶ Let $\mathbf{A} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \\ \leftarrow & \mathbf{c}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{c}_k^\top & \rightarrow \end{bmatrix} \in \mathbb{R}^{n \times d}$ (forget the $1/\sqrt{n}$ scaling)
- ▶ Let $\mathbf{C} = \begin{bmatrix} \leftarrow & \mathbf{c}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{c}_k^\top & \rightarrow \end{bmatrix} \in \mathbb{R}^{k \times d}$
- ▶ We try to approximate \mathbf{A} with \mathbf{BC} , where $\mathbf{B} \in \{0, 1\}^{n \times k}$ is matrix where $B_{i,j} = 1$ iff \mathbf{x}_i is assigned to group j .
- ▶ The k -means cost is $\|\mathbf{A} - \mathbf{BC}\|_F^2$
 - ▶ Here, $\|\cdot\|_F$ is a matrix norm called Frobenius norm, which treats the $n \times d$ matrix as a vector in nd -dimensional Euclidean space
- ▶ \mathbf{BC} is a particular rank $\leq k$ approximation of the data matrix \mathbf{A}

3 / 15

Matrix factorization

- ▶ Try to approximate A with BC , where $B \in \mathbb{R}^{n \times k}$ and $C \in \mathbb{R}^{k \times d}$, to minimize $\|A - BC\|_F^2$.
 - ▶ Think of B as the encodings of the data in A
 - ▶ At least get “dimension reduction” when $k < d$
 - ▶ Maybe not quite as space-efficient as k -means
- ▶ **Theorem** (Schmidt, 1907; Eckart-Young, 1936): Optimal solution is given by truncating the singular value decomposition (SVD) of A

Singular value decomposition

- ▶ Every matrix $A \in \mathbb{R}^{n \times d}$ —say, with rank r —can be written as

$$A = \sum_{i=1}^r \sigma_i u_i v_i^T$$

where

- ▶ $\sigma_1 \geq \dots \geq \sigma_r > 0$ (singular values)
- ▶ $u_1, \dots, u_n \in \mathbb{R}^n$ (orthonormal left singular vectors)
- ▶ $v_1, \dots, v_n \in \mathbb{R}^d$ (orthonormal right singular vectors)

- ▶ Can also write as

$$A = USV^T$$

where

- ▶ $U = [u_1 | \dots | u_r] \in \mathbb{R}^{n \times r}$, satisfies $U^T U = I$
- ▶ $S = \text{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{r \times r}$
- ▶ $V = [v_1 | \dots | v_r] \in \mathbb{R}^{d \times r}$, satisfies $V^T V = I$
- ▶ This is sometimes called the thin SVD, since the matrices U and V are as thin as possible.

Truncated SVD

- ▶ Let A have SVD $A = \sum_{i=1}^r \sigma_i u_i v_i^T$ (rank of A is r)
- ▶ Truncate at rank k (for any $k \leq r$): rank- k SVD

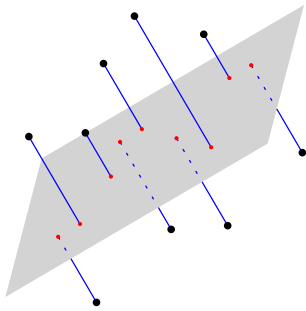
$$A_k := \sum_{i=1}^k \sigma_i u_i v_i^T$$

- ▶ Can write as $A_k := U_k S_k V_k^T$, where
 - ▶ $U_k = [u_1 | \dots | u_k] \in \mathbb{R}^{n \times k}$, satisfies $U^T U = I$
 - ▶ $S_k = \text{diag}(\sigma_1, \dots, \sigma_k) \in \mathbb{R}^{k \times k}$
 - ▶ $V_k = [v_1 | \dots | v_k] \in \mathbb{R}^{d \times k}$, satisfies $V^T V = I$
- ▶ **Theorem** (Schmidt/Eckart-Young):

$$\|A - A_k\|_F^2 = \min_{M: \text{rank}(M)=k} \|A - M\|_F^2 = \sum_{i=k+1}^r \sigma_i^2$$

Encoder/decoder interpretation I

- ▶ Encoder: $x \mapsto V_k^T x \in \mathbb{R}^k$
 - ▶ Encoding rows of A : $AV_k = U_k S_k$
- ▶ Decoder: $z \mapsto V_k z \in \mathbb{R}^d$
 - ▶ Decoding rows of $U_k S_k$: $U_k S_k V_k^T = A_k$
- ▶ Same as k -dimensional PCA mapping!
 - ▶ $A^T A = V S^2 V^T$, so eigenvectors of $A^T A$ are right singular vectors of A , non-zero eigenvalues are squares of the singular values
- ▶ PCA/SVD finds k -dimensional subspace of smallest sum of squared distances to data points.



Encoder/decoder interpretation II

- ▶ Example: OCR data, compare original image to decoding of k -dimensional PCA encoding ($k \in \{1, 10, 50, 200\}$)



Figure 2: OCR digits, original and decoding of k -dimensional PCA encoding

8 / 15

Computation: power iteration

- ▶ How to find v_1, \dots, v_k ?
- ▶ Easier problem: how to find v_1 ?
- ▶ Power iteration:
 - ▶ Start with any unit vector $u \in \mathbb{R}^d$
 - ▶ Repeat: Multiply u by $A^T A$ to get new vector u' ; renormalize back to unit length to replace u

9 / 15

Why does power iteration work?

- ▶ Ignore renormalization for now
- ▶ Then after t iterations, resulting vector is

$$\underbrace{(A^T A) \cdots (A^T A)}_{t \text{ times}} u = \underbrace{(V \Lambda V^T) \cdots (V \Lambda V^T)}_{t \text{ times}} u = V \Lambda^t V^T u$$

- ▶ Look at Λ^t :

$$\Lambda^t = \lambda_1^t \begin{bmatrix} 1 & & & \\ & \left(\frac{\lambda_2}{\lambda_1}\right)^t & & \\ & & \ddots & \\ & & & \left(\frac{\lambda_r}{\lambda_1}\right)^t \end{bmatrix} \approx \lambda_1^t \begin{bmatrix} 1 & & & \\ & 0 & & \\ & & \ddots & \\ & & & 0 \end{bmatrix}$$

whenever $\lambda_2 < \lambda_1$ and t is large.

- ▶ Therefore

$$V \Lambda^t V^T \approx \lambda_1^t v_1 v_1^T$$

and hence

$$V \Lambda^t V^T u \approx \lambda_1^t v_1 v_1^T u = (\lambda_1^t v_1^T u) v_1$$

which renormalizes back to v_1 (assuming $v_1^T u \neq 0$).

10 / 15

Beyond the top eigenvector

- ▶ To get v_2, \dots, v_k :
 - ▶ Can deflate $A^T A$: replace $A^T A$ by $A^T A - \lambda_1 v_1 v_1^T$; now v_2 is top eigenvector.
 - ▶ This only works well for small k , since errors accumulate
- ▶ Alternative: to get some $U \in \mathbb{R}^{d \times k}$ with same (similar) range as V_k , use subspace iteration.

11 / 15

Application: Topic modeling

- ▶ Start with n documents, represent using “bag-of-words” count vectors
- ▶ Arrange in matrix $A \in \mathbb{R}^{n \times d}$, where d is vocabulary size

	aardvark	abacus	abalone	...
doc 1	3	0	0	...
doc 2	7	0	4	...
doc 3	2	4	0	...
\vdots	\vdots	\vdots	\vdots	

- ▶ Rank k SVD provides an approximate factorization

$$A \approx BC$$

- where $B \in \mathbb{R}^{n \times k}$ and $C \in \mathbb{R}^{k \times d}$
- ▶ This use of SVD is called [Latent Semantic Analysis \(LSA\)](#)
 - ▶ Interpret rows of C as “topics”
 - ▶ $B_{i,t}$ is “weight” of document i on topic t
 - ▶ If rows of C were probability distributions, could interpret as $C_{t,w}$ as probability that word w appears in topic t

12 / 15

Application: Matrix completion

- ▶ Start with ratings of movies given by users
- ▶ Arrange in a matrix $A \in \mathbb{R}^{n \times d}$, where $A_{i,j}$ is rating given by user i for movie j .
 - ▶ Netflix: $n = 480000$, $d = 18000$; on avg, each user rates 200 movies
 - ▶ Most entries of A are unknown
- ▶ Idea: Approximate A with low-rank matrix, i.e., find

$$U = \begin{bmatrix} \leftarrow & \mathbf{u}_1^T & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{u}_n^T & \rightarrow \end{bmatrix}, \quad M = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{m}_1 & \cdots & \mathbf{m}_d \\ \downarrow & & \downarrow \end{bmatrix}$$

- with goal of minimizing $\|A - UM\|_F^2$
- ▶ Note: If all entries of A were observed, we could do this with truncated SVD.
 - ▶ Need to find a low-rank approximation without all of A :
[\(low-rank\) matrix completion](#)

13 / 15

SGD for matrix completion

- ▶ Instead, try to minimize

$$\sum_{(i,j) \in \Omega} (A_{i,j} - \mathbf{u}_i^T \mathbf{m}_j)^2$$

- where Ω indexes the collection of observed entries of A .
- ▶ Not convex in (U, M) ; still use (stochastic) gradient descent anyway
 - ▶ For epoch $p = 1, 2, \dots$:
 - ▶ For each $(i, j) \in \Omega$ in random order:

$$\begin{aligned} \mathbf{u}_i &:= \mathbf{u}_i - 2\eta(\mathbf{u}_i^T \mathbf{m}_j - A_{i,j})\mathbf{m}_j \\ \mathbf{m}_j &:= \mathbf{m}_j - 2\eta(\mathbf{u}_i^T \mathbf{m}_j - A_{i,j})\mathbf{u}_i \end{aligned}$$

- (Could also switch or randomize order of updates)
- ▶ Can also use regularization and other objectives / optimization algorithms

14 / 15

Feature representations from matrix completion

- ▶ Small data set ($n = 6040$ users, $d = 3952$ movies, $|\Omega| = 800000$ ratings)
- ▶ Fit U and M by optimizing (regularized) least squares objective
- ▶ Are $\mathbf{m}_1, \dots, \mathbf{m}_d \in \mathbb{R}^k$ useful feature vectors for movies?
- ▶ Some nearest-neighbor pairs ($\mathbf{m}_j, \text{NN}(\mathbf{m}_j)$):
 - ▶ Toy Story (1995), Toy Story 2 (1999)
 - ▶ Sense and Sensibility (1995), Emma (1996)
 - ▶ Heat (1995), Carlito's Way (1993)
 - ▶ The Crow (1994), Blade (1998)
 - ▶ Forrest Gump (1994), Dances with Wolves (1990)
 - ▶ Mrs. Doubtfire (1993), The Bodyguard (1992)

15 / 15