

New topic: Low-Rank Approximation via Randomized Algorithms

- We've learned a lot of methods to handle matrices:
 - When A has no structure: LU / QR / SVD (general but expensive)
 - When A is sparse (or Ax easy to evaluate): Krylov (sparse A usually arises in PDE problems)
- One of the most significant shifts in numerical analysis / applied math in recent years is the need to handle massive volume of data.

Challenge: 1) massive high-dim data sets / matrices

2) The structure is less explicit in many cases

3) Presence of noise and corruption in matrix entries

- How do we deal with high-dim data?

Observation: high-dim data can often be approximated with low-rank matrices

$$\begin{matrix} A & \approx & B & \cdot & C \\ m \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix}_n & & m \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix}_k & & \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix}_k \end{matrix} \quad \leftarrow \begin{matrix} \text{cheaper to} \\ \text{store and operate} \end{matrix}$$

$k < \min\{m, n\}$

Finding such B and C is not a new math problem.

We can formalize it as follows.

Goal: Given $A \in \mathbb{R}^{m \times n}$, $k < n$ (assume $m \geq n$)

$$\text{Find } \min_{\text{rank}(\hat{A}) \leq k} \|A - \hat{A}\| \quad (*)$$

Here we take $\|\cdot\|$ to be 2-norm or Frobenius norm

Solution to (*) is given by the truncated SVD of A

Let $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$ be singular value of A

$\mathbb{R}^{m \times m} \ni U = [\vec{u}_1, \vec{u}_2, \dots, \vec{u}_m]$ be left-singular vectors of A

$\mathbb{R}^{n \times n} \ni V = [\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n]$ be right-singular vectors of A

$$A = U \Sigma V^T, \quad \Sigma = \text{diag}(\sigma_1, \dots, \sigma_n) \in \mathbb{R}^{m \times n}$$
$$= \sum_{i=1}^n \sigma_i u_i v_i^T \quad \leftarrow \text{rank}(A) = \# \text{ nonzero } \sigma_i \text{'s}$$

Now we take the truncated SVD of A

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

$$= \begin{bmatrix} \vec{u}_1 & \dots & \vec{u}_k \end{bmatrix} \begin{bmatrix} \sigma_1 & \dots & 0 \\ & \ddots & \\ 0 & & \sigma_k \end{bmatrix} \begin{bmatrix} -\vec{v}_1^T \\ \vdots \\ -\vec{v}_k^T \end{bmatrix} \in \mathbb{R}^{m \times n}$$
$$= \begin{bmatrix} \vec{u}_1 & \dots & \vec{u}_m \end{bmatrix} \underbrace{\begin{bmatrix} \sigma_1 & \dots & 0 \\ & \ddots & \\ 0 & & \sigma_k & \dots & 0 \end{bmatrix}}_{=: \Sigma_k} \begin{bmatrix} -\vec{v}_1^T \\ \vdots \\ -\vec{v}_n \end{bmatrix}$$

clearly $\text{rank}(A_k) \leq k$,

$$\|A - A_k\|_2 = \|U(\Sigma - \Sigma_k)V^T\|_2 = \|\Sigma - \Sigma_k\|_2 = \sigma_{k+1}$$

$$\|A - A_k\|_F = \|\Sigma - \Sigma_k\|_F = \left(\sum_{j=k+1}^n \sigma_j^2 \right)^{1/2}$$

Thm (Eckart - Young)

$$\min_{\text{rank}(\hat{A}) \leq k} \|A - \hat{A}\|_2 = \sigma_{k+1}$$

$$\min_{\text{rank}(\hat{A}) \leq k} \|A - \hat{A}\|_F = \left(\sum_{j=k+1}^n \sigma_j^2 \right)^{1/2}$$

Pf: We prove the 2-norm case only.

It suffices to show that $\|A - \hat{A}\|_2 \geq \sigma_{k+1}$

for any $\text{rank}(\hat{A}) \leq k$

It suffices to show that $\exists x \in \mathbb{R}^n$, s.t. $\frac{\|(A - \hat{A})x\|_2}{\|x\|_2} \geq \sigma_{k+1}$

I want to find $x \in \mathbb{R}^n$ such that $\hat{A}x = 0$

and $x \in \text{span}\{v_1, \dots, v_{k+1}\}$ ($x = \sum_{i=1}^{k+1} \alpha_i v_i$, $\sum_{i=1}^{k+1} \alpha_i^2 = 1$)

$$\Rightarrow \frac{\|(A - \hat{A})x\|_2}{\|x\|_2} = \frac{\|Ax\|_2}{\|x\|_2} = \left\| \sum_{i=1}^{k+1} \sigma_i u_i \alpha_i \right\|_2 \geq \sigma_{k+1}$$

Such x always exists:

Since $\text{rank}(\hat{A}) \leq k$, we know $\dim \text{Null}(\hat{A}) \geq n - k$

but $\dim \text{span}\{v_1, \dots, v_{k+1}\} = k + 1$

$$\Rightarrow \text{Null}(\hat{A}) \cap \text{span}\{v_1, \dots, v_{k+1}\} \neq \emptyset$$



∴ The best rank k approximation of A is given by

k -truncated SVD of A

-
- An equivalent formulation of low-rank approximation is the Principal Component Analysis (PCA).

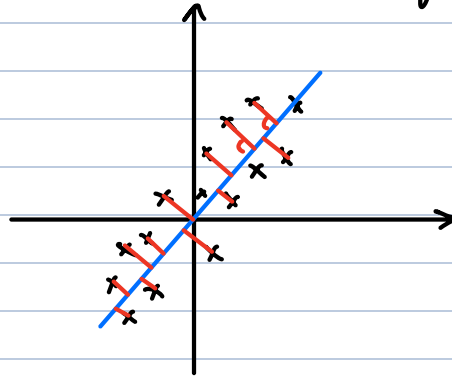
Let data points be m -dim vectors, stored in n columns of A .

PCA aims to find k vectors whose span best contains the data points in A . We can assume these vectors are orthonormal basis and form the following problem:

$$\min_{\substack{Q \in \mathbb{R}^{m \times k} \\ Q^T Q = I_m}} \|A - \underbrace{Q Q^T A}_{=: P_Q}\| \quad (**)$$

projection onto span of Q

that is, $A \approx Q Q^T A$, Q is an approximate basis
for the range of A



PCA is essentially low-rank approximation:

- On one hand, $\min_{\text{rank}(\hat{A}) \leq k} \|A - \hat{A}\| \leq \min_{\substack{Q \in \mathbb{R}^{m \times k} \\ Q^T Q = I_m}} \|A - \underbrace{Q Q^T A}_{\text{rank}(Q Q^T A) \leq k}\|$

- On the other hand, we take $Q = [u_1 \dots u_k] \in \mathbb{R}^{m \times k}$

$$\text{then } \underbrace{Q Q^T A}_{=: C} = [u_1 \dots u_k] \begin{bmatrix} -u_1^T \\ \vdots \\ -u_k^T \end{bmatrix} [u_1 \dots u_m] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \begin{bmatrix} -v_1^T \\ \vdots \\ -v_n^T \end{bmatrix}$$

$$= [u_1 \dots u_k] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_k \end{bmatrix} \begin{bmatrix} -v_1^T \\ \vdots \\ -v_k^T \end{bmatrix} \leftarrow \text{truncated SVD of } A$$

• How to efficiently compute low-rank approximation of A ?

- Apply SVD to A then truncate \leftarrow expensive! $O(mn^2)$

- Workaround: two stage algorithm from PCA

Stage A: compute orthonormal basis Q whose span
approximates $\text{span}(A)$

$$\Rightarrow A \approx Q Q^T A \text{ and } Q^T Q = I$$

Stage B: Compute $C = Q^T A \in \mathbb{R}^{k \times n}$ $\leftarrow O(kmn)$

then compute SVD of C (using whatever method):

$$B = \tilde{U} \Sigma V^T \leftarrow O(k^2 n)$$

$$\text{then } A \approx (Q \tilde{U}) \Sigma V^T \leftarrow O(k^2 m)$$

• How to find Q ? ("Randomized sketch")

Idea: Approximate $\text{span}\{u_1, \dots, u_k\}$, the top k singular vectors of A , with a single power iteration

Intuition 1: Suppose A has exactly rank k , so the best rank k approximation of A is A itself. Suppose $\sigma_k > 0$

$$\text{Compute } Y = A \Omega, \quad \Omega \in \mathbb{R}^{n \times k}$$

$$[\dot{y}_1 \dots \dot{y}_k] = [\dot{u}_1 \dots \dot{u}_k] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_k \end{bmatrix} \begin{bmatrix} -v_1^T \\ \vdots \\ -v_k^T \end{bmatrix} [\dot{w}_1 \dots \dot{w}_k]$$

\uparrow

$\dot{w}_i \in \mathbb{R}^n$ random vectors
 Y has linearly independent columns as long as $\text{rank}(V^T \Omega) = k$

For iid normal entries of Ω , this happens almost surely.

To get orthonormal basis, we compute $Y = QR$

$$\text{Clearly } \|A - Q Q^T A\|_2 = 0 = \sigma_{k+1} \text{ a.s.}$$

$V^T \Omega$ still has iid normal entries

$$\text{Intuition 2: Now let } A = \underbrace{\sum_{i=1}^k \sigma_i u_i v_i^T}_{=: \hat{A}_k} + \underbrace{\sum_{j=k+1}^n \sigma_j u_j v_j^T}_{=: E}$$

($\|E\|_2 = \sigma_{k+1}$ small)

then $\tilde{Y}_i = A w_i = \hat{A}_k w_i + E w_i$

$$\tilde{Y} = \underset{\substack{\uparrow \\ \text{basis for} \\ \text{span}(\hat{A}_k)}}{Y} + E \underset{\substack{\nwarrow \\ \text{small perturbation}}}{\Omega}$$

The span of \tilde{Y} is a good approximation to Y with high probability as long as we oversample, i.e. take $\Omega \in \mathbb{R}^{n \times (k+p)}$ instead (p is oversampling parameter)

Algorithm: (Randomized SVD)

Stage A: 1) Generate iid Gaussian random matrix $\Omega \in \mathbb{R}^{n \times (k+p)}$

2) Compute $Y = A \Omega \in \mathbb{R}^{m \times (k+p)}$

3) Compute QR fact. $Y = Q R$, $Q \in \mathbb{R}^{m \times (k+p)}$

Stage B: 1) Compute $B = Q^T A \in \mathbb{R}^{(k+p) \times n}$

2) Compute SVD: $B = \tilde{U} \Sigma V^T$
 $\tilde{U} \in \mathbb{R}^{(k+p) \times (k+p)}$

$$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_{k+p}) \in \mathbb{R}^{(k+p)^2}$$

$$V \in \mathbb{R}^{n \times (k+p)}$$

3) Compute $U = Q \tilde{U} \in \mathbb{R}^{m \times (k+p)}$

Output: $A \approx \sum_{i=1}^k \sigma_i u_i v_i^T$