

CSCI B609: “Foundations of Data Science”

Lecture 8/9: Faster Power Method and Applications of SVD

Slides at <http://grigory.us/data-science-class.html>

Grigory Yaroslavtsev

<http://grigory.us>

Faster Power Method

- PM drawback: $A^T A$ is dense even for sparse A
- Pick random Gaussian \mathbf{x} and compute $B^k \mathbf{x}$
- $\mathbf{x} = \sum_{i=1}^d c_i \mathbf{v}_i$ (augment \mathbf{v}_i 's to o.n.b. if $r < d$)
- $B^k \mathbf{x} \approx (\sigma_1^{2k} \mathbf{v}_1 \mathbf{v}_1^T) (\sum_{i=1}^d c_i \mathbf{v}_i) = \sigma_1^{2k} c_1 \mathbf{v}_1$
 $B^k \mathbf{x} = (A^T A)(A^T A) \dots (A^T A) \mathbf{x}$

- **Theorem:** If \mathbf{x} is unit \mathbb{R}^d -vector, $|\mathbf{x}^T \mathbf{v}_1| \geq \delta$:
 - V = subspace spanned by \mathbf{v}_i 's for $\sigma_j \geq (1 - \epsilon)\sigma_1$
 - \mathbf{w} = unit vector after $k = \frac{1}{2\epsilon} \ln \left(\frac{1}{\epsilon\delta} \right)$ iterations of PM

$\Rightarrow \mathbf{w}$ has a component at most ϵ orthogonal to V

Faster Power Method: Analysis

- $A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ and $\mathbf{x} = \sum_{i=1}^d c_i \mathbf{v}_i$
- $B^k \mathbf{x} = \sum_{i=1}^d \sigma_i^{2k} \mathbf{v}_i \mathbf{v}_i^T \sum_{j=1}^d c_j \mathbf{v}_j = \sum_{i=1}^d \sigma_i^{2k} c_i \mathbf{v}_i$

$$\|B^k \mathbf{x}\|_2^2 = \left\| \sum_{i=1}^d \sigma_i^{2k} c_i \mathbf{v}_i \right\|_2^2 = \sum_{i=1}^d \sigma_i^{4k} c_i^2 \geq \sigma_1^{4k} c_1^2 \geq \sigma_1^{4k} \delta^2$$

- (Squared) component orthogonal to V is

$$\sum_{i=m+1}^d \sigma_i^{4k} c_i^2 \leq (1 - \epsilon)^{4k} \sigma_1^{4k} \sum_{i=m+1}^d c_i^2 \leq (1 - \epsilon)^{4k} \sigma_1^{4k}$$

- Component of $\mathbf{w} \perp V \leq (1 - \epsilon)^{2k} / \delta \leq \epsilon$

Choice of \mathbf{x}

- \mathbf{y} random spherical Gaussian with unit variance

- $\mathbf{x} = \frac{\mathbf{y}}{\|\mathbf{y}\|_2}$:

$$\Pr \left[|\mathbf{x}^T \mathbf{v}| \leq \frac{1}{20\sqrt{d}} \right] \leq \frac{1}{10} + 3e^{-d/64}$$

- $\Pr \left[\|\mathbf{y}\|_2 \geq 2\sqrt{d} \right] \leq 3e^{-d/64}$ (Gaussian Annulus)

- $\mathbf{y}^T \mathbf{v} \sim N(0,1) \Rightarrow \Pr \left[\left| \mathbf{y}^T \mathbf{v} \right| \leq \frac{1}{10} \right] \leq \frac{1}{10}$

- Can set $\delta = \frac{1}{20\sqrt{d}}$ in the “faster power method”

Singular Vectors and Eigenvectors

- Right singular vectors are eigenvectors of $A^T A$
- σ_i^2 are eigenvalues of $A^T A$
- Left singular vectors are eigenvectors of AA^T
- $A^T A$ satisfies $\forall \mathbf{x}: \mathbf{x}^T B \mathbf{x} \geq 0$
 - $B = \sum_i \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T$
 - $\forall \mathbf{x}: \mathbf{x}^T \mathbf{v}_i \mathbf{v}_i^T \mathbf{x} = (\mathbf{x}^T \mathbf{v}_i)^2 \geq 0$
 - Such matrices are called positive semi-definite
- Any p.s.d matrix can be decomposed as $A^T A$

Application of SVD: Centering Data

- Minimize sum of squared distances from A_i to S_k
- **SVD**: best fitting S_k if data is centered
- What if not?
- **Thm.** S_k that minimizes squared distance goes through centroid of the point set:

$$\frac{1}{n} \sum A_i$$

- Will only prove for $k = 1$, analogous proof for arbitrary k (see textbook)

Application of SVD: Centering Data

- **Thm.** Line that minimizes squared distance goes through the centroid
- Line: $\ell = \mathbf{a} + \lambda \mathbf{v}$; distance $dist(\mathbf{A}_i, \ell)$
- $\|\mathbf{A}_i - \mathbf{a}\|_2^2 = dist(\mathbf{A}_i, \ell)^2 + \langle \mathbf{v}, \mathbf{A}_i \rangle^2$
- Center so that $\sum_{i=1}^n \mathbf{A}_i = \mathbf{0}$ by subtracting the centroid

$$\begin{aligned} \sum_i^n dist(\mathbf{A}_i, \ell)^2 &= \sum_{i=1}^n (\|\mathbf{A}_i - \mathbf{a}\|_2^2 - \langle \mathbf{v}, \mathbf{A}_i \rangle^2) \\ &= \sum_{i=1}^n (\|\mathbf{A}_i\|_2^2 + \|\mathbf{a}\|_2^2 - 2\langle \mathbf{A}_i, \mathbf{a} \rangle - \langle \mathbf{v}, \mathbf{A}_i \rangle^2) \\ &= \sum_{i=1}^n \|\mathbf{A}_i\|_2^2 + n\|\mathbf{a}\|_2^2 - 2\langle \sum_{i=1}^n \mathbf{A}_i, \mathbf{a} \rangle - \sum_{i=1}^n \langle \mathbf{v}, \mathbf{A}_i \rangle^2 \\ &= \sum_{i=1}^n \|\mathbf{A}_i\|_2^2 + n\|\mathbf{a}\|_2^2 - \sum_{i=1}^n \langle \mathbf{v}, \mathbf{A}_i \rangle^2 \end{aligned}$$

- Minimized when $\mathbf{a} = \mathbf{0}$

Principal Component Analysis

- $n \times d$ matrix: customers \times movies preference
- $n = \#$ customers, $d = \#$ movies
- $A_{ij} =$ how much customer i likes movie j
- Assumption: A_{ij} can be described with k factors
 - Customers and movies: vectors in \mathbf{u}_i and $\mathbf{v}_i \in \mathbb{R}^k$

- $A_{ij} = \langle \mathbf{u}_i, \mathbf{v}_j \rangle$

- Solution: A_k

$$\begin{matrix} & & & \text{factors} \\ & & & \left(\begin{array}{c} \\ \\ \\ \end{array} \right) \\ \text{customers} & \left(\begin{array}{c} \\ \\ \\ \end{array} \right) & = & \left(\begin{array}{c} \\ \\ \\ \end{array} \right) \left(\begin{array}{c} \text{movies} \\ \\ \\ \end{array} \right) \\ & A & = & U \quad V \end{matrix}$$

Class Project

- **Survey of 3-5 research papers**
 - Closely related to the topics of the class
 - Algorithms for high-dimensional data
 - Fast algorithms for numerical linear algebra
 - Algorithms for machine learning and/or clustering
 - Algorithms for streaming and massive data
 - Office hours if you need suggestions
 - Individual (not a group) project
 - **1-page Proposal Due: October 31, 2016 at 23:59 EST**
 - **Final Deadline: December 09, 2016 at 23:59 EST**
- Submission by e-mail to **Lisul Islam (IU id: islammdl)**
 - Submission Email Title: Project + Space + “Your Name”
 - Submission format: **PDF from LaTeX**

Separating mixture of k Gaussians

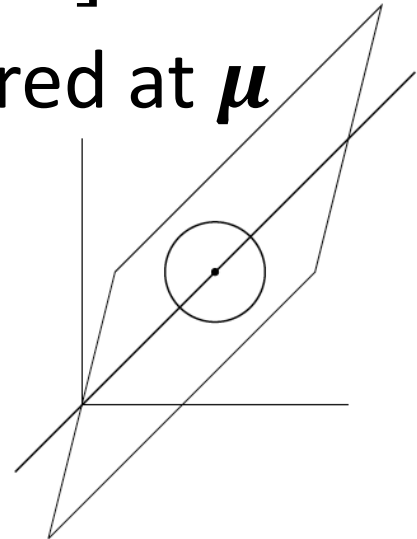
- **Sample origin problem:**
 - Given samples from k **well-separated** spherical Gaussians
 - **Q:** Did they come from the same Gaussian?
- δ = distance between centers
- For two Gaussians naïve separation requires
$$\delta > \omega(d^{1/4})$$
- **Thm.** $\delta = \Omega(k^{1/4})$ suffices
- **Idea:**
 - Project on a k -dimensional subspace through centers
 - **Key fact:** This subspace can be found via SVD
 - Apply naïve algorithm

Separating mixture of k Gaussians

- **Easy fact:** Projection preserves the property of being a unit-variance spherical Gaussian
- **Def.** If p is a probability distribution, **best fit line** $\{c\mathbf{v}, c \in \mathbb{R}\}$ is:

$$\mathbf{v} = \operatorname{argmax}_{|\mathbf{v}|=1} \mathbb{E}_{\mathbf{x} \sim p} \left[(\mathbf{v}^T \mathbf{x})^2 \right]$$

- **Thm:** Best fit line for a Gaussian centered at $\boldsymbol{\mu}$ passes through $\boldsymbol{\mu}$ and the origin



Best fit line for a Gaussian

- **Thm:** Best fit line for a Gaussian centered at $\boldsymbol{\mu}$ passes through $\boldsymbol{\mu}$ and the origin

$$\begin{aligned}\mathbb{E}_{\mathbf{x} \sim p} \left[(\mathbf{v}^T \mathbf{x})^2 \right] &= \mathbb{E}_{\mathbf{x} \sim p} \left[(\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu}) + \mathbf{v}^T \boldsymbol{\mu})^2 \right] \\ &= \mathbb{E}_{\mathbf{x} \sim p} \left[\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu})^2 + 2(\mathbf{v}^T \boldsymbol{\mu}) \mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu}) + (\mathbf{v}^T \boldsymbol{\mu})^2 \right] \\ &= \mathbb{E}_{\mathbf{x} \sim p} \left[\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu})^2 \right] + 2(\mathbf{v}^T \boldsymbol{\mu}) \mathbb{E}_{\mathbf{x} \sim p} \left[\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu}) \right] + (\mathbf{v}^T \boldsymbol{\mu})^2 \\ &= \mathbb{E}_{\mathbf{x} \sim p} \left[\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu})^2 \right] + (\mathbf{v}^T \boldsymbol{\mu})^2 \\ &= \sigma^2 + (\mathbf{v}^T \boldsymbol{\mu})^2\end{aligned}$$

- Where we used:

- $\mathbb{E}_{\mathbf{x} \sim p} [\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu})] = \mathbf{0}$

- $\mathbb{E}_{\mathbf{x} \sim p} [\mathbf{v}^T (\mathbf{x} - \boldsymbol{\mu})^2] = \sigma^2$

- Best fit line maximizes $(\mathbf{v}^T \boldsymbol{\mu})^2$

Best fit subspace for one Gaussian

- Best fit k -dimensional subspace V_k :

$$V_k = \underset{V: \dim(V)=k}{\operatorname{argmax}} \mathbb{E}_{\mathbf{x} \sim p} \left[\|\operatorname{proj}(\mathbf{x}, V)\|_2^2 \right]$$

- For a spherical Gaussian V is a best-fit k -dimensional subspace **iff** it contains $\boldsymbol{\mu}$
- If $\boldsymbol{\mu} = 0$ then any k -dim. subspace is best fit
- If $\boldsymbol{\mu} \neq 0$ then best fit line \mathbf{v} goes through $\boldsymbol{\mu}$
 - Same greedy process as SVD projects on \mathbf{v}
 - After projection we have Gaussian with $\boldsymbol{\mu} = 0$
 - Any $(k - 1)$ -dimensional subspace would do

Best fit subspace for k Gaussians

- **Thm.** p is a mixture of k spherical Gaussians \Rightarrow best fit k -dim. subspace contains their centers
- $p = w_1 \mathbf{p}_1 + w_2 \mathbf{p}_2 + \dots + w_k \mathbf{p}_k$
- Let V be a subspace of dimension $\leq k$

$$\mathbb{E}_{x \sim p} \left[\left\| \text{proj}(x, V) \right\|_2^2 \right] = \sum_{i=1}^k w_i \mathbb{E}_{x \sim p_i} \left[\left\| \text{proj}(x, V) \right\|_2^2 \right]$$

- Each term is maximized if V contains all μ'_i s
- If we only have a finite number of samples then accuracy has to be analyzed carefully

HITS Algorithm for Hubs and Authorities

- Document ranking: project on 1st singular vector
- WWW: directed graph with links = edges
- n Authorities: pages containing original info
- d Hubs: collections of links to authorities
 - Authority depends on importance of pointing hubs
 - Hub quality depends on how authoritative links are
- Authority vector: $\mathbf{v}_j, j = 1, \dots, n: \mathbf{v}_j \sim \sum_{i=1}^d \mathbf{u}_i \mathbf{A}_{ij}$
- Hub vector: $\mathbf{u}_i, i = 1, \dots, d: \mathbf{u}_i \sim \sum_{j=1}^n \mathbf{v}_j \mathbf{A}_{ij}$
- Use power method: $\mathbf{u} = \mathbf{A}\mathbf{v}, \mathbf{v} = \mathbf{A}^T \mathbf{u}$
- Converges to first left/right singular vectors

Exercises

- Ex. 1: A is $n \times n$ matrix with orthonormal rows
 - Show that it has orthonormal columns
- Ex. 2: Interpret the left and right singular vectors of the document x term matrix
- Ex. 3. Use power method to compute singular values of the matrix:

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$