

SVD and Low-rank Approximation

Lecture 19

Nov 1, 2022

Matrix Rank

Given $m \times n$ matrix A the *column rank* of A is the maximum number of linearly independent columns of A . The *row rank* is the maximum number of linearly independent rows of A .

Non-obvious fact: column rank = to row rank = rank(A)

Fact: A has rank r iff A can be written as sum of k rank 1 matrices

$$A = \sum_{i=1}^r y_i z_i^T = YZ^T$$

where Y is $m \times r$ matrix and Z is $r \times n$ matrix.

Singular Value Decomposition (SVD)

Let A be a $m \times n$ real-valued matrix

- a_i denotes vector corresponding to row i
- m rows. think of each row as a data point in \mathbb{R}^n
- Data applications: $m \gg n$
- Other notation: A is a $n \times d$ matrix.

Singular Value Decomposition (SVD)

Let A be a $m \times n$ real-valued matrix

- a_i denotes vector corresponding to row i
- m rows. think of each row as a data point in \mathbb{R}^n
- Data applications: $m \gg n$
- Other notation: A is a $n \times d$ matrix.

SVD theorem: A can be written as UDV^T where

- V is a $n \times n$ orthonormal matrix
- D is a $m \times n$ diagonal matrix with $\leq \min\{m, n\}$ non-zeroes called the singular values of A
- U is a $m \times m$ orthonormal matrix

SVD

Let $d = \min\{m, n\}$.

- u_1, u_2, \dots, u_m columns of U , left singular vectors of A
- v_1, v_2, \dots, v_n columns of V (rows of V^T) right singular vectors of A
- $\sigma_1 \geq \sigma_2 \geq \dots, \geq \sigma_d \geq 0$ are non-negative singular values where $d = \min\{m, n\}$. And $\sigma_i = D_{i,i}$

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

SVD

Let $d = \min\{m, n\}$.

- u_1, u_2, \dots, u_m columns of U , left singular vectors of A
- v_1, v_2, \dots, v_n columns of V (rows of V^T) right singular vectors of A
- $\sigma_1 \geq \sigma_2 \geq \dots, \geq \sigma_d \geq 0$ are non-negative singular values where $d = \min\{m, n\}$. And $\sigma_i = D_{i,i}$

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

We can in fact restrict attention to r the rank of A .

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

SVD

Interpreting A as a linear operator $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$

- Columns of V is an orthonormal basis and hence $V^T x$ for $x \in \mathbb{R}^n$ expresses x in the V basis. Note that $V^T x$ is a rigid transformation (does not change length of x).
- Let $y = V^T x$. D is a diagonal matrix which only stretches y along the coordinate axes. Also adjusts dimension to go from n to m with right number of zeroes.
- Let $z = Dy$. Then Uz is a rigid transformation that expresses z in the basis corresponding to rows of U .

Thus any linear operator can be split into a sequence of three simpler/basic type of transformations

Low rank approximation property of SVD

Question: Given $A \in \mathbb{R}^{m \times n}$ and integer k find a matrix B of rank at most k such that $\|A - B\|$ is minimized

Low rank approximation property of SVD

Question: Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and integer k find a matrix \mathbf{B} of rank at most k such that $\|\mathbf{A} - \mathbf{B}\|$ is minimized

Fact: For Frobenius norm *and* spectral norm optimum for all k is captured by SVD.

That is, $\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ is the best rank k approximation to \mathbf{A}

$$\|\mathbf{A} - \mathbf{A}_k\|_F = \min_{\mathbf{B}: \text{rank}(\mathbf{B}) \leq k} \|\mathbf{A} - \mathbf{B}\|_F = \sqrt{\sum_{i>k} \sigma_i^2}$$

$$\|\mathbf{A} - \mathbf{A}_k\|_2 = \min_{\mathbf{B}: \text{rank}(\mathbf{B}) \leq k} \|\mathbf{A} - \mathbf{B}\|_2 = \sigma_{k+1}$$

Low rank approximation property of SVD

Question: Given $A \in \mathbb{R}^{m \times n}$ and integer k find a matrix B of rank at most k such that $\|A - B\|$ is minimized

Fact: For Frobenius norm *and* spectral norm optimum for all k is captured by SVD.

That is, $A_k = \sum_{i=1}^k \sigma_i u_i v_i^T$ is the best rank k approximation to A

$$\|A - A_k\|_F = \min_{B: \text{rank}(B) \leq k} \|A - B\|_F = \sqrt{\sum_{i>k} \sigma_i^2}$$

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

Why this magic? Frobenius norm and basic properties of vector

Geometric meaning

What is the best rank 1 matrix B that minimizes $\|A - B\|_F$

Since B is rank 1, $B = uv^T$ where $v \in \mathbb{R}^n$ and $u \in \mathbb{R}^m$

Without loss of generality v is a unit vector

Geometric meaning

What is the best rank 1 matrix B that minimizes $\|A - B\|_F$

Since B is rank 1, $B = uv^T$ where $v \in \mathbb{R}^n$ and $u \in \mathbb{R}^m$

Without loss of generality v is a unit vector

$$\|A - uv^T\|_F^2 = \sum_{i=1}^m \|a_i - u(i)v\|^2$$

Geometric meaning

What is the best rank 1 matrix B that minimizes $\|A - B\|_F$

Since B is rank 1, $B = uv^T$ where $v \in \mathbb{R}^n$ and $u \in \mathbb{R}^m$

Without loss of generality v is a unit vector

$$\|A - uv^T\|_F^2 = \sum_{i=1}^m \|a_i - u(i)v\|^2$$

If we know v then best u to minimize above is determined. Why?

Geometric meaning

What is the best rank 1 matrix B that minimizes $\|A - B\|_F$

Since B is rank 1, $B = uv^T$ where $v \in \mathbb{R}^n$ and $u \in \mathbb{R}^m$

Without loss of generality v is a unit vector

$$\|A - uv^T\|_F^2 = \sum_{i=1}^m \|a_i - u(i)v\|^2$$

If we know v then best u to minimize above is determined. Why?

For fixed v , $u(i) = \langle a_i, v \rangle$

Geometric meaning

What is the best rank 1 matrix B that minimizes $\|A - B\|_F$

Since B is rank 1, $B = uv^T$ where $v \in \mathbb{R}^n$ and $u \in \mathbb{R}^m$

Without loss of generality v is a unit vector

$$\|A - uv^T\|_F^2 = \sum_{i=1}^m \|a_i - u(i)v\|^2$$

If we know v then best u to minimize above is determined. Why?

For fixed v , $u(i) = \langle a_i, v \rangle$

$\|a_i - \langle a_i, v \rangle v\|_2$ is distance of a_i from line described by v .

Geometric meaning

What is the best rank 1 matrix B that minimizes $\|A - B\|_F$
It is to find unit vector/direction \mathbf{v} to minimize

$$\sum_{i=1}^m \|\mathbf{a}_i - \langle \mathbf{a}_i, \mathbf{v} \rangle \mathbf{v}\|^2$$

which is same as finding unit vector \mathbf{v} to maximize

$$\sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2$$

Geometric meaning

What is the best rank 1 matrix B that minimizes $\|A - B\|_F$
It is to find unit vector/direction \mathbf{v} to minimize

$$\sum_{i=1}^m \|\mathbf{a}_i - \langle \mathbf{a}_i, \mathbf{v} \rangle \mathbf{v}\|^2$$

which is same as finding unit vector \mathbf{v} to maximize

$$\sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2$$

Note: Maximum value is $\|A\|_2^2$, the spectral norm square! How to find best \mathbf{v} ? Not obvious: we will come to it a bit later

Best rank two approximation

Consider $k = 2$. What is the best rank 2 matrix B that minimizes $\|A - B\|_F$

Since B has rank 2 we can assume without loss of generality that $B = u_1 v_1^T + u_2 v_2^T$ where v_1, v_2 are orthogonal unit vectors (span a space of dimension 2)

Best rank two approximation

Consider $k = 2$. What is the best rank 2 matrix B that minimizes $\|A - B\|_F$

Since B has rank 2 we can assume without loss of generality that $B = u_1 v_1^T + u_2 v_2^T$ where v_1, v_2 are orthogonal unit vectors (span a space of dimension 2)

Minimizing $\|A - B\|_F^2$ is same as finding orthogonal vectors v_1, v_2 to maximize

$$\sum_{i=1}^m (\langle a_i, v_1 \rangle^2 + \langle a_i, v_2 \rangle^2)$$

in other words the best fit 2-dimensional space

Greedy algorithm

- Find \mathbf{v}_1 as the best rank 1 approximation. That is
$$\mathbf{v}_1 = \arg \max_{\mathbf{v}, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2$$
- For \mathbf{v}_2 solve
$$\arg \max_{\mathbf{v} \perp \mathbf{v}_1, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2.$$

Alternatively: let $\mathbf{a}'_i = \mathbf{a}_i - \langle \mathbf{a}_i, \mathbf{v}_1 \rangle \mathbf{v}_1$. Let
$$\mathbf{v}_2 = \arg \max_{\mathbf{v}, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}'_i, \mathbf{v} \rangle^2$$

Greedy algorithm

- Find \mathbf{v}_1 as the best rank 1 approximation. That is
$$\mathbf{v}_1 = \arg \max_{\mathbf{v}, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2$$
- For \mathbf{v}_2 solve
$$\arg \max_{\mathbf{v} \perp \mathbf{v}_1, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2.$$

Alternatively: let $\mathbf{a}'_i = \mathbf{a}_i - \langle \mathbf{a}_i, \mathbf{v}_1 \rangle \mathbf{v}_1$. Let
$$\mathbf{v}_2 = \arg \max_{\mathbf{v}, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}'_i, \mathbf{v} \rangle^2$$

Greedy algorithm works!

Greedy algorithm correctness

Proof that Greedy works for $k = 2$.

Suppose w_1, w_2 are orthogonal unit vectors that form the best fit 2-d space. Let H be the space spanned by w_1, w_2 .

Claim: Any two orthogonal unit vectors in H will yield same value.

Suffices to prove that

$$\sum_{i=1}^m (\langle a_i, v_1 \rangle^2 + \langle a_i, v_2 \rangle^2) \geq \sum_{i=1}^m (\langle a_i, w_1 \rangle^2 + \langle a_i, w_2 \rangle^2)$$

Greedy algorithm correctness

Case 1: $v_1 \in H$ then done because we can assume wlog that $w_1 = v_1$ and v_2 is at least as good as w_2 .

Case 2: $v_1 \notin H$. Let v'_1 be projection of v_1 onto H and $v''_1 = v_1 - v'_1$ be the component of v_1 orthogonal to H .

Greedy algorithm correctness

Case 1: $\mathbf{v}_1 \in H$ then done because we can assume wlog that $\mathbf{w}_1 = \mathbf{v}_1$ and \mathbf{v}_2 is at least as good as \mathbf{w}_2 .

Case 2: $\mathbf{v}_1 \notin H$. Let \mathbf{v}'_1 be projection of \mathbf{v}_1 onto H and $\mathbf{v}''_1 = \mathbf{v}_1 - \mathbf{v}'_1$ be the component of \mathbf{v}_1 orthogonal to H . Note that $\|\mathbf{v}'_1\|^2 + \|\mathbf{v}''_1\|_2^2 = \|\mathbf{v}_1\|_2^2 = 1$.

Wlog we can assume by rotation that $\mathbf{w}_1 = \frac{1}{\|\mathbf{v}'_1\|_2} \mathbf{v}'_1$ and \mathbf{w}_2 is orthogonal to \mathbf{v}'_1 . Hence \mathbf{w}_2 is orthogonal to \mathbf{v}_1 .

Greedy algorithm correctness

Case 1: $v_1 \in H$ then done because we can assume wlog that $w_1 = v_1$ and v_2 is at least as good as w_2 .

Case 2: $v_1 \notin H$. Let v'_1 be projection of v_1 onto H and $v''_1 = v_1 - v'_1$ be the component of v_1 orthogonal to H . Note that $\|v'_1\|^2 + \|v''_1\|_2^2 = \|v_1\|_2^2 = 1$.

Wlog we can assume by rotation that $w_1 = \frac{1}{\|v'_1\|_2} v'_1$ and w_2 is orthogonal to v'_1 . Hence w_2 is orthogonal to v_1 .

Therefore v_2 is at least as good as w_2 , and v_1 is at least as good as w_1 which implies the desired claim.

Greedy algorithm for general k

- Find \mathbf{v}_1 as the best rank 1 approximation. That is $\mathbf{v}_1 = \arg \max_{\mathbf{v}, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2$
- For \mathbf{v}_k solve $\arg \max_{\mathbf{v} \perp \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k-1}, \|\mathbf{v}\|_2=1} \sum_{i=1}^k \langle \mathbf{a}_i, \mathbf{v} \rangle^2$ which is same as solving $k = 1$ with vectors $\mathbf{a}'_1, \mathbf{a}'_2, \dots, \mathbf{a}'_m$ that are residuals. That is $\mathbf{a}'_i = \mathbf{a}_i - \sum_{j=1}^{k-1} \langle \mathbf{a}_i, \mathbf{v}_j \rangle \mathbf{v}_j$

Proof of correctness is via induction and is a straight forward generalization of the proof for $k = 2$

Summarizing

$$\sigma_j^2 = \sum_{i=1}^m \langle a_i, v_j \rangle^2$$

By greedy construction $\sigma_1 \geq \sigma_2 \geq \dots$,

Let r be the (row) rank of A . v_1, v_2, \dots, v_r span the row space of A and $\sigma_j = 0$ for $j > r$. Can choose v_{r+1}, \dots, v_n to ensure orthonormal basis of R^n

u_1 determined by v_1 and u_2 determined by v_1, v_2 and so on. Can show that they are orthogonal.

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

Power method

Thus SVD relies on being able to solve $k = 1$ case

Given m vectors $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m \in \mathbb{R}^n$ solve

$$\max_{\mathbf{v} \in \mathbb{R}^n, \|\mathbf{v}\|_2=1} \langle \mathbf{a}_i, \mathbf{v} \rangle^2$$

How do we solve the above problem?

Let $\mathbf{B} = \mathbf{A}^T \mathbf{A}$ Then

$$\begin{aligned} \mathbf{B} &= \left(\sum_{i=1}^m \sigma_i \mathbf{v}_i \mathbf{u}_i^T \right) \left(\sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T \right) \\ &= \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T \end{aligned}$$

Power method continued

Let $B = A^T A$ Then

$$\begin{aligned} B^2 &= \left(\sum_{i=1}^r \sigma_i^2 v_i v_i^T \right) \left(\sum_{i=1}^r \sigma_i^2 v_i v_i^T \right) \\ &= \sum_{i=1}^r \sigma_i^4 v_i v_i^T. \end{aligned}$$

More generally

$$B^k = \sum_{i=1}^r \sigma_i^{2k} v_i v_i^T$$

Power method continued

Let $B = A^T A$ Then

$$\begin{aligned} B^2 &= \left(\sum_{i=1}^r \sigma_i^2 v_i v_i^T \right) \left(\sum_{i=1}^r \sigma_i^2 v_i v_i^T \right) \\ &= \sum_{i=1}^r \sigma_i^4 v_i v_i^T. \end{aligned}$$

More generally

$$B^k = \sum_{i=1}^r \sigma_i^{2k} v_i v_i^T$$

If $\sigma_1 > \sigma_2$ then B^k converges to $\sigma_1^{2k} v_1 v_1^T$ and we can identify v_1 from B^k . But expensive to compute B^k

Power method continued

Pick a random (unit) vector $\mathbf{x} \in \mathbb{R}^n$. Then $\mathbf{x} = \sum_{i=1}^n \lambda_i \mathbf{v}_i$ since $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ is a basis for \mathbb{R}^n .

$$B^k \mathbf{x} = \left(\sum_{i=1}^r \sigma_i^{2k} \mathbf{v}_i \mathbf{v}_i^T \right) \left(\sum_{i=1}^d \lambda_i \mathbf{v}_i \right) \rightarrow \sigma_1^{2k} \lambda_1 \mathbf{v}_1$$

Can obtain \mathbf{v}_1 by normalizing $B^k \mathbf{x}$ to a unit vector.

Computing $B^k \mathbf{x}$ is easier via a series of matrix vector multiplications

Power method continued

Pick a random (unit) vector $\mathbf{x} \in \mathbb{R}^n$. Then $\mathbf{x} = \sum_{i=1}^n \lambda_i \mathbf{v}_i$ since $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ is a basis for \mathbb{R}^n .

$$B^k \mathbf{x} = \left(\sum_{i=1}^r \sigma_i^{2k} \mathbf{v}_i \mathbf{v}_i^T \right) \left(\sum_{i=1}^d \lambda_i \mathbf{v}_i \right) \rightarrow \sigma_1^{2k} \lambda_1 \mathbf{v}_1$$

Can obtain \mathbf{v}_1 by normalizing $B^k \mathbf{x}$ to a unit vector.

Computing $B^k \mathbf{x}$ is easier via a series of matrix vector multiplications

Why random \mathbf{x} ? So as to ensure $\lambda_1 > 0$ with good probability.

Theorem

Suppose $\sigma_1 > \sigma_2$. Then with probability $(1 - \delta)$, power method converges to a vector \mathbf{v} such that $\langle \mathbf{v}, \mathbf{v}_1 \rangle \geq (1 - \epsilon)$ after

$O\left(\frac{\log n + \log(1/\epsilon) + \log(1/\delta)}{\log(\sigma_1/\sigma_2)}\right)$ iterations.

Power method continued

Pick a random (unit) vector $\mathbf{x} \in \mathbb{R}^n$. Then $\mathbf{x} = \sum_{i=1}^n \lambda_i \mathbf{v}_i$ since $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ is a basis for \mathbb{R}^n .

$$B^k \mathbf{x} = \left(\sum_{i=1}^r \sigma_i^k \mathbf{v}_i \mathbf{v}_i^T \right) \left(\sum_{i=1}^d \lambda_i \mathbf{v}_i \right) \rightarrow \sigma_1^{2k} \lambda_1 \mathbf{v}_1$$

Convergence depends on σ_1/σ_2 . What if $\sigma_1 \simeq \sigma_2$? Power method may not converge to \mathbf{v}_1 but output will be some vector in the space spanned by $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_h$ where σ_h is the largest h such that $\sigma_1 \simeq \sigma_h$. This is good enough in various applications. See references.

Principal Component Analysis

Consider A a $m \times n$ matrix where rows a_1, a_2, \dots, a_m are data points in \mathbb{R}^n

$B = A^T A$ is a symmetric positive definite matrix and has real non-negative eigenvalues

Via SVD $B = (UDV^T)^T(UDV^T) = (VD^T DV^T)$

Can check that v_1, v_2, \dots, v_r are eigen vectors of B with eigen values $\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2$

Principal Component Analysis

Consider A a $m \times n$ matrix where rows a_1, a_2, \dots, a_m are data points in \mathbb{R}^n

- Compute eigenvectors of $B = A^T A$ or singular vectors v_1, v_2, \dots, v_n which are also called the principal directions
- Approximate each a_i by its projection onto the first k singular vectors for some small k . That is $a'_i = \sum_{j=1}^k \langle a_i, v_j \rangle v_j$.
- Thus a'_1, a'_2, \dots, a'_m , a kind of *dimensionality reduction* along first k principal directions. Different from JL and is motivated by different applications (mainly statistical analysis)

PCA and Covariance Matrix

Covariance of two real-valued random variables \mathbf{X} , \mathbf{Y} is defined as

$$\text{Cov}(\mathbf{X}, \mathbf{Y}) := E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{Y} - E[\mathbf{Y}])]$$

Note that $\text{Cov}(\mathbf{X}, \mathbf{X}) = \text{Var}(\mathbf{X})$. If \mathbf{X} , \mathbf{Y} independent then $\text{Cov}(\mathbf{X}, \mathbf{Y}) = 0$ but converse is not necessarily true. There is also a related normalized measure (value in $[-1, 1]$)

$$\text{Correlation}(\mathbf{X}, \mathbf{Y}) = \frac{\text{Cov}(\mathbf{X}, \mathbf{Y})}{\sigma_{\mathbf{X}}\sigma_{\mathbf{Y}}}$$

The sign of $\text{Cov}(\mathbf{X}, \mathbf{Y})$ is an “indication” of positive vs negative correlation. Non-linear relationships between \mathbf{X} , \mathbf{Y} are not necessarily captured by covariance but still useful in many situations.

PCA and Covariance Matrix

Suppose $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)$ is a n -dimensional random variable. Thus, each \mathbf{X}_i is a random variable and may be correlated with the other variables.

Given \mathbf{X} we can define a covariance matrix \mathbf{C} where $C_{i,j} = \text{Cov}(\mathbf{X}_i, \mathbf{X}_j)$. Note that if the \mathbf{X}_i are independent then \mathbf{C} will be a diagonal matrix. Similarly one can also define a correlation matrix where the entries are the correlation coefficients instead of covariances.

PCA of \mathbf{C} reveals useful information if \mathbf{X} is in fact obtained via a linear transformation from another random variable \mathbf{Y} that lives in a lower dimension. Typically \mathbf{X} will be a noisy version of \mathbf{Y} and hence will not be a pure low rank matrix but a low rank approximation gives the important directions.

PCA and Covariance Matrix

Suppose $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)$ is a n -dimensional random variable.

Suppose we have m data points $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m \in \mathbb{R}^n$ drawn independently from the distribution of \mathbf{X} . We create a $m \times n$ matrix \mathbf{A} where \mathbf{a}_i is i 'th row. Given the empirical data matrix \mathbf{A} we would like to estimate the covariance matrix \mathbf{C} of \mathbf{X} .

Assuming we know for each i , $\mu_i = \mathbb{E}[\mathbf{X}_i]$ we can estimate $\text{Cov}(\mathbf{X}_i, \mathbf{X}_j)$ from the m data samples as

$$\frac{1}{m} \sum_{\ell=1}^m (\mathbf{a}_\ell(i) - \mu_i)(\mathbf{a}_\ell(j) - \mu_j).$$

By setting $\mathbf{a}'_\ell = \mathbf{a}_\ell - \boldsymbol{\mu}$ where $\boldsymbol{\mu}$ is the vector of expectations we see that $\mathbf{C} = \frac{1}{m} (\mathbf{A}')^T \mathbf{A}'$ is the desired estimated covariance matrix. Thus PCA on $(\mathbf{A}')^T \mathbf{A}'$ helps identify important features in the underlying distribution \mathbf{X}

PCA and Covariance Matrix

Suppose $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)$ is a n -dimensional random variable. Suppose we have m data points $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m \in \mathbb{R}^n$ drawn independently from the distribution of \mathbf{X} . We create a $m \times n$ matrix \mathbf{A} where \mathbf{a}_i is i 'th row. Given the empirical data matrix \mathbf{A} we would like to estimate the covariance matrix \mathbf{C} of \mathbf{X} .

Suppose we do not know the means $\mu_i = \mathbb{E}[\mathbf{X}_i]$. We can compute an empirical estimate from the data itself as $\frac{1}{m} \sum_{\ell=1}^m \mathbf{a}_\ell(i)$ and then the empirical mean vector it to "center" the data to compute an estimated covariance matrix as in the previous slide. Sometimes data is already assumed to be centered in which case we simply work with $\mathbf{A}^T \mathbf{A}$.

Linear least square/Regression and SVD

Linear least squares: Given $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$ find x to minimize $\|Ax - b\|_2$.

Interesting when $m > n$ the over constrained case when there is no solution to $Ax = b$ and want to find best fit.

Linear least square/Regression and SVD

Linear least squares: Given $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$ find x to minimize $\|Ax - b\|_2$.

Interesting when $m > n$ the over constrained case when there is no solution to $Ax = b$ and want to find best fit.

Geometrically Ax is a linear combination of columns of A . Hence we are asking what is the vector z in the column space of A that is closest to vector b in ℓ_2 norm.

Linear least square/Regression and SVD

Linear least squares: Given $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$ find x to minimize $\|Ax - b\|_2$.

Interesting when $m > n$ the over constrained case when there is no solution to $Ax = b$ and want to find best fit.

Geometrically Ax is a linear combination of columns of A . Hence we are asking what is the vector z in the column space of A that is closest to vector b in ℓ_2 norm.

Closest vector to b is the projection of b into the column space of A so it is “obvious” geometrically. How do we find it?

Linear least square/Regression and SVD

Linear least squares: Given $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$ find x to minimize $\|Ax - b\|_2$.

Interesting when $m > n$ the over constrained case when there is no solution to $Ax = b$ and want to find best fit.

Geometrically Ax is a linear combination of columns of A . Hence we are asking what is the vector z in the column space of A that is closest to vector b in ℓ_2 norm.

Closest vector to b is the projection of b into the column space of A so it is “obvious” geometrically. How do we find it? Find an orthonormal basis z_1, z_2, \dots, z_r for the columns of A . Compute projection b' as $b' = \sum_{j=1}^r \langle b, z_j \rangle z_j$ and output answer as $\|b - b'\|_2$.

Linear least square/Regression and SVD

Linear least squares: Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$ find \mathbf{x} to minimize $\|\mathbf{Ax} - \mathbf{b}\|_2$.

Closest vector to \mathbf{b} is the projection of \mathbf{b} into the column space of \mathbf{A} so it is “obvious” geometrically. Find an orthonormal basis $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_r$ for the columns of \mathbf{A} . Compute projection \mathbf{b}' as $\mathbf{b}' = \sum_{j=1}^r \langle \mathbf{b}, \mathbf{z}_j \rangle \mathbf{z}_j$ and output answer as $\|\mathbf{b} - \mathbf{b}'\|_2$.

Finding the basis is the expensive part. Recall SVD gives $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$ which form a basis for the *row* space of \mathbf{A} but then $\mathbf{u}_1^T, \mathbf{u}_2^T, \dots, \mathbf{u}_m^T$ form a basis for the *column* space of \mathbf{A} . Hence SVD gives us all the information to find \mathbf{b}' . In fact we have

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2^2 = \sum_{i=r+1}^m \langle \mathbf{u}_i^T, \mathbf{b} \rangle^2$$