

Singular Value Decomposition

- Data reduction
- Data-driven generalization of Fourier transform (FFT)
- Tailored to specific problem.

- Solve $Ax=b$ for non square A
→ regression
- Basis PCA
→ correlation

Simple/interpretable linear algebra.

Scalable.

(1)

$$\bar{X} = \begin{bmatrix} | & | & \dots & | \\ x_1 & x_2 & \dots & x_m \\ | & | & \dots & | \end{bmatrix} = U \Sigma V^T = \begin{bmatrix} | & | & \dots & | \\ u_1 & u_2 & \dots & u_n \\ | & | & \dots & | \end{bmatrix} \begin{bmatrix} \sigma_1 & \sigma_2 & \dots & \sigma_m \\ & & & 0 \end{bmatrix} \begin{bmatrix} | & | & \dots & | \\ v_1 & v_2 & \dots & v_m \\ | & | & \dots & | \end{bmatrix}^T = U \Sigma V^T$$

x_1, x_2, \dots, x_m are $x_k \in \mathbb{R}^n$ (eigen faces).
 U, V are orthogonal/unitary.
 Σ is diagonal with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m \geq 0$, ordered by importance.
 $U^T U = U U^T = I_{n \times n}$
 $V V^T = V^T V = I_{m \times m}$
 Σ is diagonal.

singular values.
 left singular vectors
 right singular vectors
 eigen time series.
 first column
 mixture of u 's to make x_1

$\bar{X} = \begin{bmatrix} \rightarrow t \\ x_1 & x_2 & \dots & x_m \\ | & | & \dots & | \end{bmatrix}$

- ★ Guaranteed to exist
- ★ Unique

(2)

$$\bar{X} = \begin{bmatrix} | & | & \dots & | \\ x_1 & x_2 & \dots & x_m \\ | & | & \dots & | \end{bmatrix} = U \Sigma V^T = \begin{bmatrix} | & | & \dots & | \\ u_1 & u_2 & \dots & u_n \\ | & | & \dots & | \end{bmatrix} \begin{bmatrix} \sigma_1 & \sigma_2 & \dots & \sigma_m \\ & & & 0 \end{bmatrix} \begin{bmatrix} -v_1^T \\ -v_2^T \\ \vdots \\ -v_m^T \end{bmatrix} = U \hat{\Sigma} V^T$$

$\hat{\Sigma}$ is $m \times m$ matrix.
 m columns of u .
 Economy SVD

Matrix Approximation (SVD)

(3)

$$\bar{X} = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \dots + \sigma_m u_m v_m^T + 0$$

$$= \underbrace{\sigma_1 \begin{bmatrix} | \\ u_1 \\ | \end{bmatrix} \begin{bmatrix} \hline v_1^T \hline \end{bmatrix}}_{\text{rank-1 matrix}} + \underbrace{\sigma_2 \begin{bmatrix} | \\ u_2 \\ | \end{bmatrix} \begin{bmatrix} \hline v_2^T \hline \end{bmatrix}}_{\text{rank-1 matrix}} + \dots + \sigma_m \begin{bmatrix} | \\ u_m \\ | \end{bmatrix} \begin{bmatrix} \hline v_m^T \hline \end{bmatrix}$$

(truncate at rank r) $\approx \tilde{U} \tilde{\Sigma} \tilde{V}^T$: rank r approximation.

Ex: Best rank-2 matrix

Eckard-Young thm [1936]

$$\arg \min_{\tilde{X} \text{ s.t. } \text{rank}(\tilde{X})=r} \|\bar{X} - \tilde{X}\|_F = \tilde{U} \tilde{\Sigma} \tilde{V}^T$$

After truncating

$$\tilde{U}^T \tilde{U} = I_{r \times r}$$

$$\tilde{U} \tilde{U}^T \neq I$$

Dominant Correlations (SVD)

(4)

$$\underbrace{\bar{X}^T \bar{X}}_{m \times m} = \underbrace{\begin{bmatrix} -x_1^T \\ -x_2^T \\ \vdots \\ -x_m^T \end{bmatrix}}_{m \times n} \underbrace{\begin{bmatrix} | & & | \\ x_1 & x_2 & \dots & x_m \\ | & & | \end{bmatrix}}_{n \times m} = \underbrace{\begin{bmatrix} x_1^T x_1 & x_1^T x_2 & \dots & x_1^T x_m \\ x_2^T x_1 & x_2^T x_2 & \dots & x_2^T x_m \\ \vdots & \vdots & \ddots & \vdots \\ x_m^T x_1 & x_m^T x_2 & \dots & x_m^T x_m \end{bmatrix}}_{\text{Correlation Matrix (column-wise)}} \rightarrow \text{PSD}$$

$n \geq m$

$$\text{If } \bar{X} = \hat{U} \hat{\Sigma} \hat{V}^T \\ \bar{X}^T = \hat{V} \hat{\Sigma} \hat{U}^T$$

$$\bar{X}^T \bar{X} = \hat{V} \hat{\Sigma} \underbrace{\hat{U}^T \hat{U}}_I \hat{\Sigma} \hat{V}^T = \hat{V} \hat{\Sigma}^2 \hat{V}^T$$

$$\underbrace{\bar{X}^T \bar{X}}_{\substack{\downarrow \\ \text{eigenvectors} \\ \text{of} \\ \text{correlation} \\ \text{matrix}}} \hat{V} = \hat{V} \hat{\Sigma}^2 \quad \text{(column-wise correlation matrix)} \\ (m \times n)(n \times m) = m \times m$$

\downarrow eigenvalues of correlation matrix.

Similarly,

$$\bar{X} \bar{X}^T = \hat{U} \hat{\Sigma} \underbrace{\hat{V}^T \hat{V}}_I \hat{\Sigma} \hat{U}^T = \hat{U} \hat{\Sigma}^2 \hat{U}^T$$

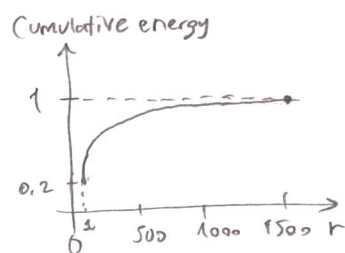
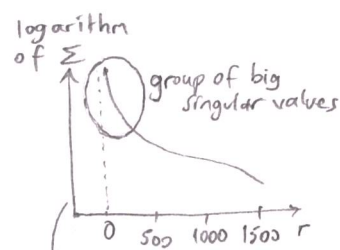
$$\underbrace{\bar{X} \bar{X}^T}_{\substack{\downarrow \\ \text{eigenvectors} \\ \text{of the correlation matrix}}} \hat{U} = \hat{U} \hat{\Sigma}^2 \quad \text{(row-wise correlation matrix)} \\ (n \times m)(m \times n) = n \times n$$

\downarrow eigenvalues of the correlation matrix

bigger than column-wise correlation matrix.

(5)

$\left(\frac{\text{Cumulative sum of first } r \text{ singular values}}{\text{sum of all singular values}} \right) \Rightarrow$ how much energy/information is in the first r modes compared to all modes.



Semilogy

The Frobenius Norm of Matrices

(7)

$$\arg \min_{\tilde{X} \text{ s.t. } \text{rank}(\tilde{X})=r} \|\bar{X} - \tilde{X}\|_F = \hat{U} \hat{\Sigma} \hat{V}^T$$

$$\|A\|_F = \sqrt{\sum_{i,j} |a_{ij}|^2} \quad (A) \text{ matrix to a big vector.}$$

\downarrow frobenius norm

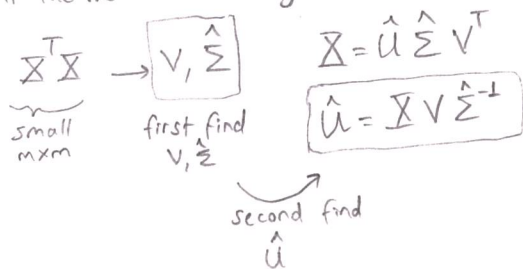
$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$

SVD Method of Snapshots

(8)

Sirovich [1987]

If the matrix is too large



It is not recommended to compute SVD using correlation matrices but if the matrix is too big that you cannot load it to memory:

Load first column and take dot product of itself.
Load first column and second column and dot product.

Only loading two vectors, the result $m \times m$ small correlation matrix can be obtained.

$$\begin{bmatrix} X_1^T X_1 & X_1^T X_2 & \dots \\ \vdots & & \end{bmatrix}$$

Correlation matrix

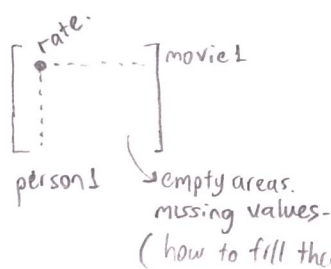
$$X_i^T X_j = \langle X_i, X_j \rangle \text{ dot product.}$$

Matrix Completion and the Netflix Price

(9)

$$X = \begin{bmatrix} | & | & \dots & | \\ X_1 & X_2 & \dots & X_m \\ | & | & \dots & | \\ \text{person 1} & \text{person 2} & \dots & \text{person } m \end{bmatrix}$$

movie 1
movie 2
...



RPCA
Robust Principal Component Analysis.

Unitary Transformations Geometry

(10)

$$X = U \Sigma V^T = \hat{U} \hat{\Sigma} V^T$$

SVD econ SVD

$$U U^T = U^T U = I$$

$$V V^T = V^T V = I$$

(U, V unitary matrices)
Unitary transformations preserve angles and lengths of vectors.

Just rotate vectors.

$$\langle x, y \rangle = \langle Ux, Uy \rangle \quad \forall x, y \in \mathbb{R}^n$$

over the real number field
 U, V are called orthogonal

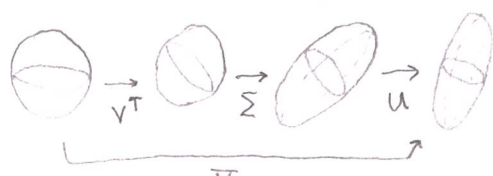
$$Q^T Q = Q Q^T = I$$

$$Q^T = Q^{-1}$$

If X complex $\in \mathbb{C}^{n \times m}$

then X^* is complex conjugate transpose.

SVD: Rotate + Scale + Rotate

$$\mathbb{R}^m \xrightarrow{\quad} \mathbb{R}^n \xrightarrow{\quad} \mathbb{R}^n \xrightarrow{\quad} \mathbb{R}^m$$


(11) (12)

$Ax=b$ linear system of equations
 \downarrow solve \downarrow
 known known

SVD allow us to generalize to non-square A .

economy SVD.
 $A = U \Sigma V^T \Rightarrow A^\dagger$ pseudo-inverse

$Ax=b$ $A^\dagger = V \Sigma^{-1} U^T$ * moore-penrose left pseudo inverse

$U \Sigma V^T = b$

$V \Sigma^{-1} U^T U \Sigma V^T x = V \Sigma^{-1} U^T b$
 $\underbrace{V \Sigma^{-1} U^T U \Sigma V^T}_I = \underbrace{V \Sigma^{-1} U^T}_{A^\dagger} b$

$\tilde{x} = V \Sigma^{-1} U^T b$

$\tilde{x} = A^\dagger b$ dagger

best x

• Underdetermined

$\min \|\tilde{x}\|_2$ s.t. $A\tilde{x}=b$

(minimum norm solution)

↳ cheapest way of multiplying A with x to get b .

minimum 2-norm x .

let's say x is energy.

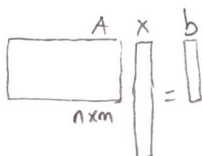
• Overdetermined.

$\min \|A\tilde{x} - b\|_2$ \leftarrow 2 norm error

(least square solution)

• Underdetermined, $n < m$ (short fat A)

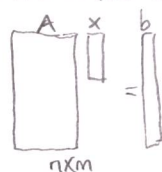
∞ many solution x given b .



(infinite solutions) *

• Overdetermined, $n > m$ (tall skinny A)

zero solutions x for generic b .



(no solution) *

Least Square Regression and the SVD

$A\tilde{x}=b$
 \downarrow
 $U \Sigma V^T V \Sigma^{-1} U^T b$

$\left\{ \begin{array}{l} A\tilde{x} \\ U \Sigma V^T V \Sigma^{-1} U^T b \\ \underbrace{U \Sigma V^T V \Sigma^{-1} U^T}_I \\ \underbrace{U \Sigma V^T V \Sigma^{-1} U^T}_I \end{array} \right\}$

$\hat{U} \hat{U}^T \neq I$ for econ SVD.

$A\tilde{x}$
 $\hat{U} \hat{U}^T b$
 projection of b onto $\text{span}(\hat{U})$
 $= \text{span}(A)$

Linear system of equations, $Ax=b$

* Solution of $Ax=b$ only exists if b is in column space of A .

$A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ $b = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$
 no solution even though $Ax=b$ underdetermined ($n < m$)

$A = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 3 & 4 \end{bmatrix}$ $b = \begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix}$ $x = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$
 there might be a solution even though $Ax=b$ overdetermined ($n > m$)

• $\text{col}(A) = \text{col}(\hat{U})$ range

• $\text{Ker}(A^T)$ orthogonal complement
 \downarrow
 kernel

• $\text{row}(A) = \text{col}(V)$

• $\text{Ker}(A)$ null space

Set of all vec. x_{null} s.t. $Ax_{\text{null}}=0$

* Also, $\dim(\text{Ker}(A)) \neq 0$
 (there are some vectors that map to zero)
 then ∞ many solutions.

$A(\underbrace{x + x_{\text{null}}}_{\text{solution}}) = b$

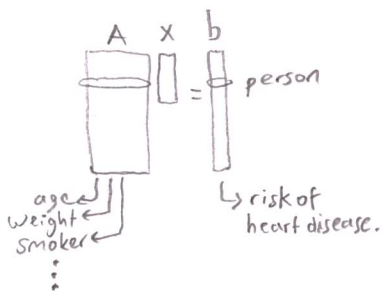
Linear Regression

(16)

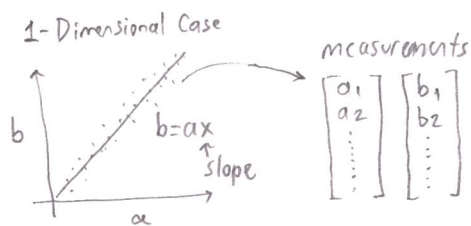
$Ax=b$ linear system of equations $\Rightarrow \tilde{x} = A^+b$

$A = \hat{U}\hat{\Sigma}\hat{V}^T \Rightarrow A^+ = \hat{V}\hat{\Sigma}^{-1}\hat{U}^T$

Linear regression model from data.



best x model



$U = \frac{a}{\|a\|_2}$

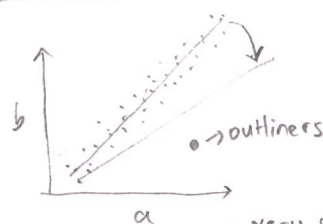
$\Sigma = \|a\|_2$

$V = 1$

best fit slope

$\tilde{x} = \frac{a^T b}{\|a\|_2^2}$

* Add L_1 penalty term to discount outliers.



* very sensitive to the outliers.

If the data is clean or only has gaussian distribution, it is optimal to fit data (SVD)

Linear Regression

$x=3$

$a = [-2, 2] \Delta a = 0.1$

$b = ax + \text{noise}$

both a & b have noisy data.

Example of Data.

$Ax=b$

$A = U\Sigma V^T$

$\tilde{x} = V\Sigma^{-1}U^T b$

$SVD \Rightarrow [U, \Sigma, V] = \text{svd}(a, \text{"econ"})$

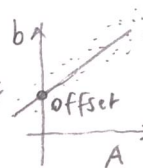
$\tilde{x} = \underbrace{V\Sigma^{-1}U^T}_\text{pseudo-inverse of matrix A} b$

* Note: If there is offset

Add ones for the offset

$\begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} A \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} b \end{bmatrix}$

Add $a \cdot x$ for the offset.



$ax + \text{offset} = b$

$\tilde{x} = V\Sigma^{-1}U^T b$

(17)

$A = \begin{bmatrix} 7 & 26 & 6 & 60 \\ 10 & 59 & 7 & 14 \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$

$b = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$

↑ ↑ ↑ ↑ 13x4

heat

4 ingredients

(overdetermined)

Ex:

- Use 10 data to create model
 - Use 3 data to validate/how accurate
- avoid overfitting.

$A \rightarrow$ parameters

- area
- close to school

$b \rightarrow$ house value.

500 data \rightarrow 250 data: model

\rightarrow 250 data: validate

maybe data collected by one neighbourhood and another neighbourhood.

Magnitude.



* Instead of taking first half, shuffle the data.

Dominant parameter for higher house price.

Dominant parameter for lower house price.

(18)

19, 20, 21
python

Principal Component Analysis (PCA) (dimensionality reduction)

(22)

Statistical Interpretation of the SVD.

Hierarchical coordinate system (based on data)

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix}$$

measurements from a single experiment.
(row instead of column)

1. Compute mean row

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j$$

$$\bar{X} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} [\bar{x}]$$

2. Subtract mean $B = X - \bar{X}$

3. Covariance matrix of the rows of B.

$$C = B^T B$$

4. Compute eigs of C.

$$V_1^T B^T B V_1$$

↳ biggest eigenvector

mean subtracted data
↑

$$C V = V D$$

↳ eigenvalues.
↓
eigenvectors.

$$\lambda = \sigma^2$$

$$T = B V$$

↓
principal components

↳ loadings.

$$B = U \Sigma V^T \Rightarrow T = U \Sigma$$

$$\frac{\sum_{k=1}^r \lambda_k}{\sum_{k=1}^n \lambda_k}$$

Principal Component Analysis

(23)

1. mean $\Rightarrow \bar{x} = \frac{1}{n} \sum_{j=1}^n x_j$

2. $B = X - \bar{X}$ mean centered matrix B.

3. Compute SVD of normalized B
 $B / \sqrt{\text{number of points}}$

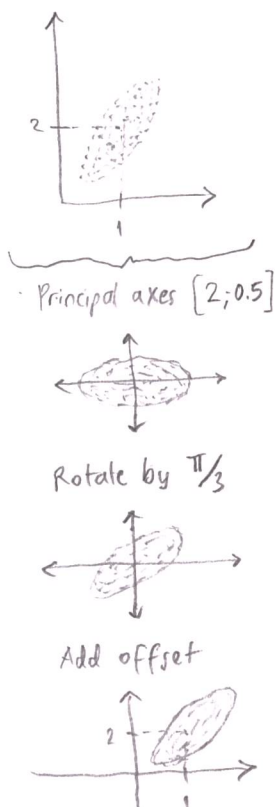
↓
Principal axes and rotation captured from data.

$$\Sigma \text{ should capture } \begin{bmatrix} 2 & 0 \\ 0 & 0.5 \end{bmatrix}$$

$$U \text{ should capture } R \text{ by } \frac{\pi}{3}$$

↓
first column of U \rightarrow first direction of maximal variance

second column of U \rightarrow second direction of maximal variance



★ SVD can be used to compute the principal component analysis which can give you some information about the distribution/statistics how data is distributed, what are the principal directions of variance and the directions that have the least variance, the most variance and you can visualize very high dimensional data.