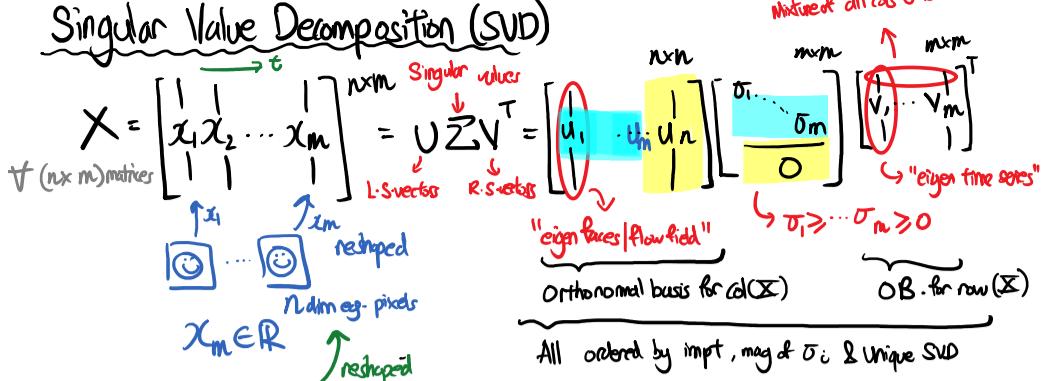


Singular Value Decomposition (SVD)



$X_m \in \mathbb{R}^n$ reshaped

$x_m \in \mathbb{R}^n$ reshaped

$$\text{Outerpate} \quad -\sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \dots + \sigma_m u_m v_m^T + 0$$

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

"Economy SVD" ($n \gg m$)

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

If $\sigma_{m+1} \dots \sigma_m$ all negligibly small, we truncate at rank r

$$X = \begin{bmatrix} U_1 & \dots & U_r & U_m \end{bmatrix}_{n \times n} \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_m & \\ & 0 & \dots & 0 \end{bmatrix}_{r \times r} \begin{bmatrix} V_1^T & & \\ & \ddots & \\ & & V_m^T \end{bmatrix}_{r \times m}$$

rank theorem $\hookrightarrow r \leq m$ and \hookrightarrow assume, $n \gg m \hookrightarrow r \ll n$

Eckard-Young Theorem [1936]

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

the best rank r approx \tilde{X} that minimizes error

Frobenius Norm, $\|A\|_F$

$$\text{eg. } \|A\|_F = \sqrt{\sum_{i,j} (A_{ij})^2}, \text{ Take matrix } A, \text{ reshape into 1 long vector, calc 2 norm of it}$$

Now after truncation, no longer symmetrical U & V .

$$\tilde{U}^T \tilde{U} = I_{rr} \neq \tilde{U} \tilde{U}^T$$

orthovectors \tilde{U}^T_{rxn} \tilde{U}_{nr} $= I_{rr}$

BUT $\tilde{U}^T_{rxn} \tilde{U}_{nr} \neq I_{rxr}$

* U & V is unitary, X isn't

either \tilde{X} or $\tilde{U} \tilde{\Sigma} \tilde{V}^T$

unit vectors, v

Linear transform $\tilde{X} \rightarrow R^n$

(matrix multip.)

Take, $\tilde{X} = U \tilde{\Sigma} V^T$ identity matrix here

• Orientation depends on U & \tilde{X}

• Axis length transform do \tilde{X} $\tilde{\Sigma}$

Axis are the PC = Vol $\frac{\tilde{\Sigma}^T \tilde{X}_c}{\tilde{n}!}$, if \tilde{X} centered already,

unit vectors, u

Compare to a transform \tilde{X} some mapping

$$Av = R_x R_y R_z \begin{bmatrix} \text{scale}_x & 0 & 0 \\ 0 & \text{scale}_y & 0 \\ 0 & 0 & \text{scale}_z \end{bmatrix} I_{\tilde{V}}$$

SAME!

Eigen decomposition of Square Cross Pdt Matrices

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

$$= V \tilde{\Sigma}^2 V^T \hookrightarrow I$$

$$\Rightarrow (\tilde{X}^T \tilde{X}) V = V \tilde{\Sigma}^2$$

eigenvalues = \sum diagonal singular values 2

Similarly $(\tilde{X} \tilde{X}^T) U = U \tilde{\Sigma}^2$

some eigenvalues for the nonzero.

- Manually do SVD of $\tilde{X}^T \tilde{X}$ or directly SVD to get $U, \tilde{\Sigma}, V$
- BUT let's say if X is so large to store as a memory. (Very rare)

Linear System, Pseudo A

$Ax = b$ linear system, solve x but now A^{-1} doesn't exist ($n \times m$)

Using SVD: $Ax = b$

$$\hat{U} \hat{\Sigma} \hat{V}^T x = b \quad (\text{Assume "E-SVD")}$$

$$\hat{V} \hat{\Sigma} \hat{U}^T \hat{U} \hat{\Sigma} \hat{V}^T x = \hat{V} \hat{\Sigma} \hat{U}^T b$$

\approx approx $\tilde{x} = \hat{A}^+ b$

A^+ := Moore-Penrose (left) pseudo inverse

So now $A \tilde{x} = b$ is true for underdetermined, but overdetermined it is not: $A \tilde{x} = \hat{U} \hat{\Sigma} \hat{V}^T \hat{V} \hat{\Sigma}^{-1} \hat{U}^T b$

$\neq b$ if not $A \tilde{x} = b$

$\hat{U} \hat{\Sigma} \hat{U}^T b$:= projection of b on the span(\hat{U}), \therefore cols \hat{U} span cols A , \therefore project of b on the span(A)

Application: Linear Regression

A $\xrightarrow{\text{age... weight}}$ x $\xrightarrow{\text{person i}}$ b $\xrightarrow{\text{Risk of heart failure}}$

To find the best x st we can predict for any person their heart failure

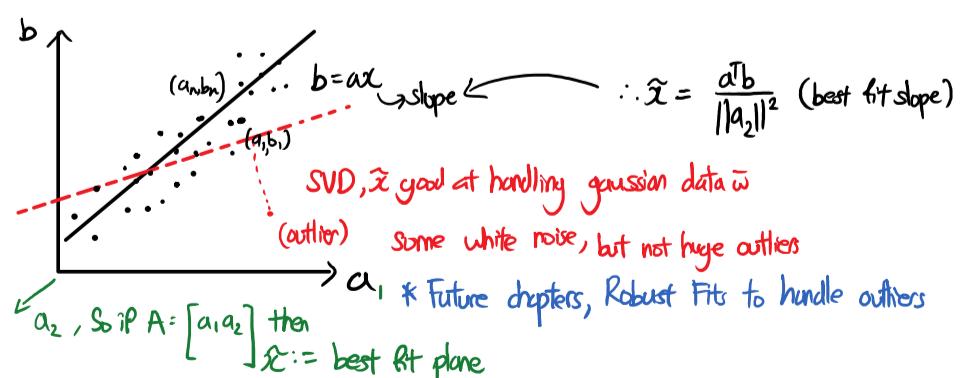
Assume general overdetermined case,

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

e.g. Simplify to one factor a

$$\begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}$$

$$U = \frac{a}{\|a\|_2}, \tilde{\Sigma} = \|a\|_2, V = I$$



Unitary Transformation of UV

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

where,

$$\begin{aligned} U^T &= U^T U = I, \text{ real values} \\ V^T &= V^T V = I \end{aligned} \quad \left\{ \begin{array}{l} U, V \text{ unitary matrices} \\ \text{Complex Conjugate Transpose} \end{array} \right.$$

Unitary Transformation (Linear transformation)

Preserve $\|x\|$ & $\|\tilde{x}\|$, rotation

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

* Orthogonal matrix:

$$AA^T = A^T A = I, \text{ real values}$$

$\therefore A^* = A^T$

Unitary matrix:

$$AA^H = A^H A = I$$

$\text{Complex Conjugate Transpose}$

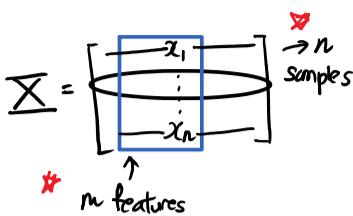
$$a \pm bi \rightarrow a \mp bi$$

If real, $b_i = 0$

Applicatⁿ: PCA

• Statistical Interpretation of SVD

• Data-driven hierarchical coordinate system that are based of principal components (orthogonal) uncorrelated axes that have maximal correlation to each measurement/ feature

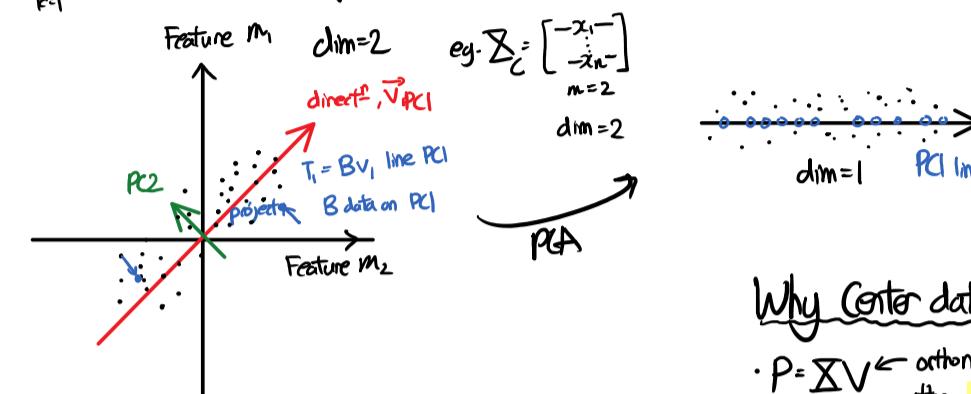


Principal Component Matrix, $P = \sum_{i=1}^r V_i$ → Right eigenvectors of $C = \frac{1}{n-1} \sum_{i=1}^n \sum_{j=1}^m (x_{ij} - \bar{x}_j)(x_{ij} - \bar{x}_j)^T$ same as the eigenvectors of $\sum_{i=1}^n \sum_{j=1}^m x_{ij} x_{ij}^T$

* Note eigendecompositⁿ, eigenvalues of $C = \frac{1}{n-1}$ of that of $\sum_{i=1}^n \sum_{j=1}^m x_{ij} x_{ij}^T$ So we can do a SVD of $\sum_{i=1}^n \sum_{j=1}^m x_{ij} x_{ij}^T$ to get V

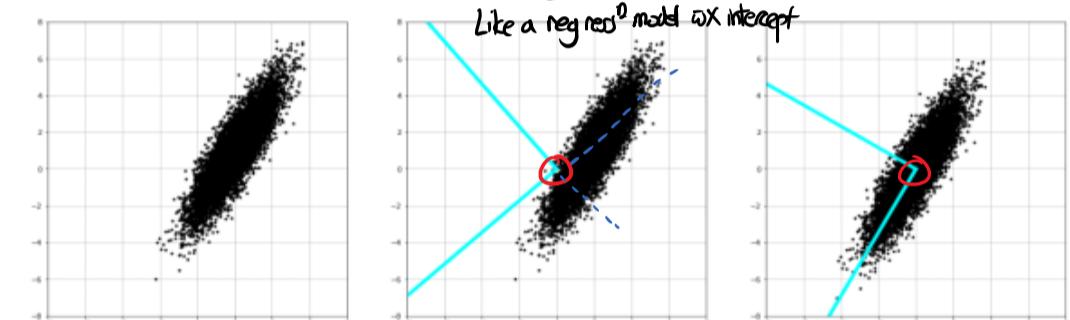
PC1: $P_1 = \sum_{i=1}^n V_i$, reflects the directⁿ of largest variance in data & has largest variance agst all the normalized linear combi & cols of $\sum_{i=1}^n \sum_{j=1}^m x_{ij} x_{ij}^T$. Subsequent P_i , cols of P ortho to each other, decreasing variance captured Since variance of data along principal component associated w eigenvalues of covariance matrix

$\frac{\sum_{k=1}^r 2\lambda_k}{\sum_{k=1}^n 2\lambda_k}$ (Upto r modes) } Gives us the no. of principal components (cols of P) we need to cover a certain % variance of the data for instance 95% of the data
(All modes)



Why Center data is a must B4 PCA?

• $P = \sum V$ ← orthonormal, unitary, rotatⁿ, hence cuts thru origin

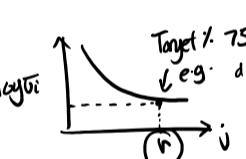


Why Standardize?

• If the cols are not of the same scale, some might have a larger variance
Hence for the n datapoints, the spread is larger. ∴ PCA would be biased

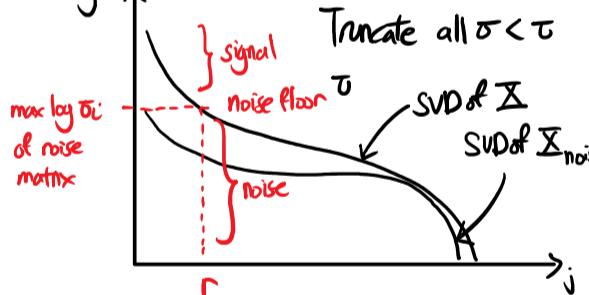
Truncation

$$\sum = U \sum V^T \approx \tilde{U} \tilde{\sum} \tilde{V}^T$$



$$\sum = \sum_{\text{"true noise}} + \sum_{\text{noise}}$$

Std normal dist.



i) For $\sum_{n \times m}$, σ_{known} ii) For $\sum_{n \times m}$, σ_{known} , $B = \frac{\text{smaller morn}}{\text{larger morn}}$

$$\tau = \frac{4}{\sqrt{3}} \sigma \sqrt{n}$$

$$\tau = \left(2(B+1) + \frac{8B}{(B+1)+(B^2+14B+1)^{1/2}} \right)^{1/2} \sigma \sqrt{n}$$

ii) For $\sum_{n \times m}$, σ_{unknown} , assume $\sigma_{\text{median}} < \tau$

$$\tau = \omega(B) \sigma_{\text{med}}, \quad \omega(B) = \frac{\lambda(B)}{\mu_B}$$

$$\text{Solution of } \int_{(1+B)^2}^{UB} \frac{[(1+\sqrt{t})^2 - t](t - (1-\sqrt{t})^2)^{1/2}}{2\sqrt{t}} dt = \frac{1}{2}$$

Alignment

$$\sum = \begin{bmatrix} x_1 & \dots & x_m \end{bmatrix}$$

(Translatⁿ)

Align so if we dot prod, features to features

In addition data cannot handle (rotatⁿ)

$$\begin{array}{c} \text{low} \\ \text{high} \\ n=1 \end{array}$$

BUT goes to show power of CNN to pad out features despite translatⁿ in data

Randomize SVD

• Data getting higher dim (eg. 1080p \rightarrow 4K), very heavy to compute $U \sum V^T$

• Assumptⁿ: Low Intrinsic rank exist for data

To improve the r capturing: Do oversampling

1) Take random Projⁿ Matrix $P \in \mathbb{R}^{m \times (m+5/10)}$

Proj high rank data \sum onto P low rank, r matrix

$$Z = \sum P$$

$$= QR \quad (\text{QR factorizatⁿ)}$$

Orthonormal basis for $\text{col}(Z)$ & $\text{col}(\sum)$ due to random vector & high dim geo. principals

$\text{col}(\sum)$ sampled by P to give Z, where we can do a QR factorizatⁿ to obtain Q

$$\sum_{n \times m} \xrightarrow{5/10} Z_{n \times r} = Q_{n \times r} R_{r \times r}$$

2) Project X down to the smaller dim Q

$$Y = Q^T \sum = U \sum V^T$$

∴ Uy orthogonal basis for Y

that will span $\text{col}(\sum)$ $\sum_y V_y^T = \sum_x V_x^T$

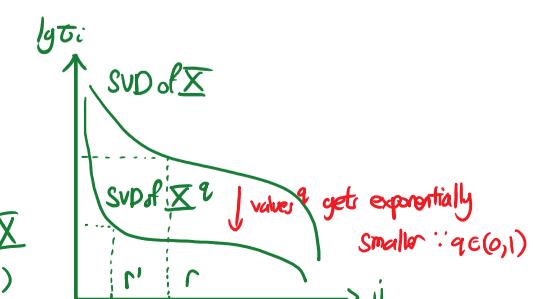
So we computed $\sum V^T$ of the data \sum , left U

$$U_r = Q U_y$$

$$\begin{array}{c|c} Q_{n \times r} & U_y_{r \times r} \\ \hline & U_r_{r \times r} \end{array} = \begin{array}{c|c} & U_r_{r \times r} \end{array}$$

$$\begin{array}{c|c} Q^T_{r \times m} & \sum_y_{r \times r} \\ \hline & V_y^T_{r \times m} \end{array} = \begin{array}{c|c} & V_y^T_{r \times m} \end{array}$$

$$= \begin{array}{c|c|c} U_y_{r \times r} & \sum_y_{r \times r} & V_y^T_{r \times m} \\ \hline & \text{SAME} & \\ \hline Z_r_{r \times r} & & V_r^T_{r \times m} \end{array}$$



If \sum low rank intrinsic data not low do Power Iteratⁿ $\sum^q = (\sum \sum^q)^T \sum$ (heavy computatⁿ) then do R-SVD

Overdetermined Data (Cancer)

$$\bar{X} = \begin{bmatrix} \bar{x}_1 & \dots & \bar{x}_n \end{bmatrix}^T$$

(mean centered) 216 samples person
40000 genes markers

$\Sigma = \begin{bmatrix} U_1 & \dots & U_r & \vdots & U_m \end{bmatrix} \begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_r & & \\ & & & \ddots & \\ & & & & \sigma_m \end{bmatrix} \begin{bmatrix} V_1^T & & & & \\ & \ddots & & & \\ & & V_m^T & & \end{bmatrix}$

- Rank r approximation 216×4000
- Cols of U := orthonormal basis for $\text{Col}(\bar{X})$ "Eigen genes"
- Singular Values
- Eigenvectors of $\bar{X}\bar{X}^T$
- Rows of V^T := Orthonormal basis for $\text{Row}(\Sigma)$ "Eigen people"

$= U \Sigma V^T$ "Economy SVD"

Now to represent the 216 people in a lower dim
↳ PCA analysis

For every person

$$\bar{x}_i \xrightarrow{\text{Linear combi of row } (\bar{X})} \text{Span of rows } (V^T)$$

in the row (\bar{X})

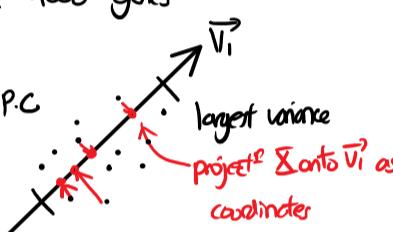
PCA : $T = \bar{X}V = U\Sigma$

where rows of V^T := Principal Components Direct^c / Loadings

V_i^T captures the most variance in the data of 4000 genes that makes up each 216 person

So to represent every 216 person onto the 1st P.C.

$$T_1 = \bar{X}V_1 \quad (\text{project } \bar{X} \text{ onto } V_1)$$

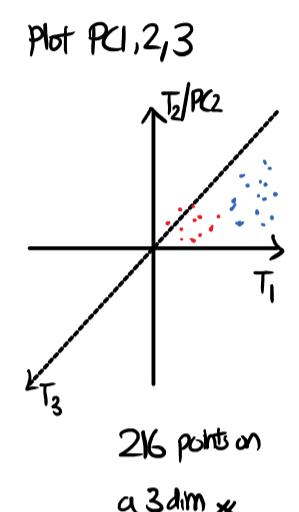


$$T_1 = \bar{X}V_1$$

$$\bar{X} = \begin{bmatrix} \bar{x}_1 & \dots & \bar{x}_{216} \end{bmatrix} \quad V_1 = \begin{bmatrix} v_{11} & \dots & v_{1m} \end{bmatrix} \quad T_1 = \begin{bmatrix} t_{11} & \dots & t_{1m} \end{bmatrix}$$

Often in python, \bar{X}^T

$$\bar{X}^T @ V_1 \quad \text{For every } i^{\text{th}} \text{ row, } V_1^T @ \bar{X}_i^T$$



216 points on a 3dim \bar{X}

Underdetermined Data (Faces)

$$\bar{X} = \begin{bmatrix} \bar{x}_1 & \dots & \bar{x}_m \end{bmatrix}^T$$

(mean centered)
32000 pixels
 m faces, 2432
Person 1 face ...

$\Sigma = \begin{bmatrix} U_1 & \dots & U_r & \vdots & U_m \end{bmatrix} \begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_r & & \\ & & & \ddots & \\ & & & & \sigma_m \end{bmatrix} \begin{bmatrix} V_1^T & & & & \\ & \ddots & & & \\ & & V_m^T & & \end{bmatrix}$

- EV of $\bar{X}\bar{X}^T$
- Cols of U := Orthonormal basis for $\text{Col}(\bar{X})$ "Eigenfaces"
- EV of $\bar{X}^T\bar{X}$
- Cols of V , rows of V^T := ON basis for $\text{Row}(\Sigma)$ "Eigenfeatures"

$= U \Sigma V^T$ "Economy SVD" \cong Rank r approx
 2432×32000

Now to represent the 2432 faces in a lower dim
↳ PCA analysis

PCA : $T = \bar{X}^T U$

For every face in \bar{X} , $\text{col}(\bar{X})$

$$\bar{X}^T / T_1 \quad (\text{col of } T) = \bar{X}^T U$$

\bar{X}^T is a linear combi of U_1 "eigenfaces"

every row $\bar{X}^T = \text{col of } \bar{X}$

$$U \quad \text{"eigenfaces"} \quad 32000 \times 32000$$

\bar{X}^T "eigenface"

Then same as above. But then if we were to plot PC1 vs PC2, there is not much clustering due to how the 1st eigenface captures the whole general face, hence all \bar{X} cols should not differ

Now how does this relate to machine learning?

$$\bar{X} = U \Sigma V^T \quad (\text{PCA})$$

$$T = \bar{X}^T U \quad \text{"eigenfaces"}$$

Test image \bar{x} , not in data then,

$$\begin{bmatrix} \bar{x}^T & \dots & \bar{x}_r^T \end{bmatrix} \bar{U}_r \quad \text{"eigenfaces"}$$

\bar{U}_r "Amt of each mode/eigenfaces inside a test face"

Inner prod of cols \bar{U} w/ test \bar{x} ($\bar{U}^T \bar{x}$)

To see the high dim image again

$$\bar{x} \approx \bar{U}_r \alpha \quad \text{32000x1 Reconstructed Image using the truncated cols of } U$$