

Low rank approximations:-

Thm 7:- $A \in \mathbb{R}^{m \times n}$ of rank " r " can be written as sum of " r " rank-one matrices of the form

$$A = \sum_{j=1}^r \sigma_j \underline{u}_j \underline{v}_j^T$$

where $\{\sigma_j\}$ are singular values and $\{\underline{u}_j\}, \{\underline{v}_j\}$ are the appropriate singular vectors.

Pf:- Recall $\underline{u} \underline{v}^T$ is a rank-one matrix

$$A = \underline{U} \underline{\Sigma} \underline{V}^T$$

$$\underline{\Sigma} = \sum_{j=1}^r \underline{\Sigma}_j \quad \text{where}$$

$$\underline{\Sigma}_j = \begin{bmatrix} \sigma_j & & & & \\ & \sigma_j & & & \\ & & \ddots & & \\ & & & \sigma_j & \\ & & & & 0 \\ \vdots & & & & & \ddots \\ & & & & & & 0 \end{bmatrix}_{m \times n}$$

$$A = \underline{U} \underline{\Sigma} \underline{V}^T$$

$$= \underline{U} \left\{ \sum_{j=1}^r \underline{\Sigma}_j \right\} \underline{V}^T$$

$$= \sum_{j=1}^r \underline{U} \underline{\Sigma}_j \underline{V}^T = \sum_{j=1}^r \sigma_j \underline{u}_j \underline{v}_j^T$$

$$A = \sum_{j=1}^r \sigma_j \underline{u}_j \underline{v}_j^T$$

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$$

Then k^{th} partial sum $\sum_{j=1}^k \sigma_j \underline{u}_j \underline{v}_j^T$ $k \leq r$

has as much energy (information) of A as possible

Thm 8: For any k with $1 \leq k \leq r$

define

$$A_k = \sum_{j=1}^k \sigma_j \underline{u}_j \underline{v}_j^T$$

$$\text{then } \|A - A_k\|_2 = \min_{\substack{B \in \mathbb{R}^{m \times n} \\ \text{rank}(B) \leq k}} \|A - B\|_2 = \sigma_{k+1}$$

Eckhart
-Young Theorem

Proof: Let there is some (B) whose $\text{rank}(B) \leq k$ such that $\|A - B\|_2 < \|A - A_k\|_2$

$$\dim(N(B)) \geq n - k \quad \text{as } \text{rank}(B) \leq k$$

Consider the subspaces

(i) W_1 : The null space of (B) which is of dimension of at least $n - k$

(ii) W_2 : The space spanned by $k+1$ right singular vectors of A i.e. $\underline{v}_1, \underline{v}_2, \underline{v}_3, \dots, \underline{v}_{k+1}$

These two subspaces have to intersect? why?

Dimensions of the two subspaces

add to $(n-k) + (k+1)$ i.e. the subspaces must at least have 1 common vector.

Let such a non-zero vector be \underline{x} i.e. $\underline{x} \in W_1 \cap W_2$

$\underline{x} \neq 0$; $\underline{x} \in N(B)$ i.e. $B\underline{x} = 0$ and $\underline{x} \in W_2$

$$\underline{x} = \sum_{i=1}^{k+1} c_i \underline{v}_i$$

$$\|A\underline{x}\|_2 = \|(A-B)\underline{x}\|_2 \leq \|A-B\|_2 \|\underline{x}\|_2 < \sigma_{k+1} \|\underline{x}\|_2 \quad \text{--- (1)}$$

$$A\underline{v}_i = \sigma_i \underline{u}_i$$

$$\begin{aligned} \|A\underline{x}\|_2^2 &= \left\| A \sum_{i=1}^{k+1} c_i \underline{v}_i \right\|_2^2 \\ &= \left\| \sum_{i=1}^{k+1} c_i \sigma_i \underline{u}_i \right\|_2^2 = \left[\sum_{i=1}^{k+1} c_i \sigma_i \underline{u}_i \right]^T \left[\sum_{i=1}^{k+1} c_i \sigma_i \underline{u}_i \right] \\ &= \sum_{i=1}^{k+1} c_i^2 \sigma_i^2 \quad (\text{use orthogonality of } \underline{u}_i) \end{aligned}$$

$$\sum_{i=1}^{k+1} c_i^2 \sigma_i^2 \geq \sum_{i=1}^{k+1} c_i^2 \sigma_{k+1}^2$$

$$= \left(\sum_{i=1}^{k+1} c_i^2 \right) \sigma_{k+1}^2$$

$$= \|\underline{x}\|_2^2 \sigma_{k+1}^2$$

$$\|A\underline{x}\|_2^2 \geq \|\underline{x}\|_2^2 \sigma_{k+1}^2 \quad \text{i.e. } \|A\underline{x}\|_2 \geq \sigma_{k+1} \|\underline{x}\|_2 \quad \text{--- (2)}$$

① & ② is a contradiction which means you cannot have a matrix B with $\text{rank}(B) \leq k$ such that $\|A - B\|_2 < \|A - A_k\|_2$

Eckhart-Young in Frobenius norm:-

Thm 9 For any k , with $1 \leq k \leq r$, the matrix $A_k = \sum_{j=1}^k \sigma_j u_j v_j^T$ also satisfies

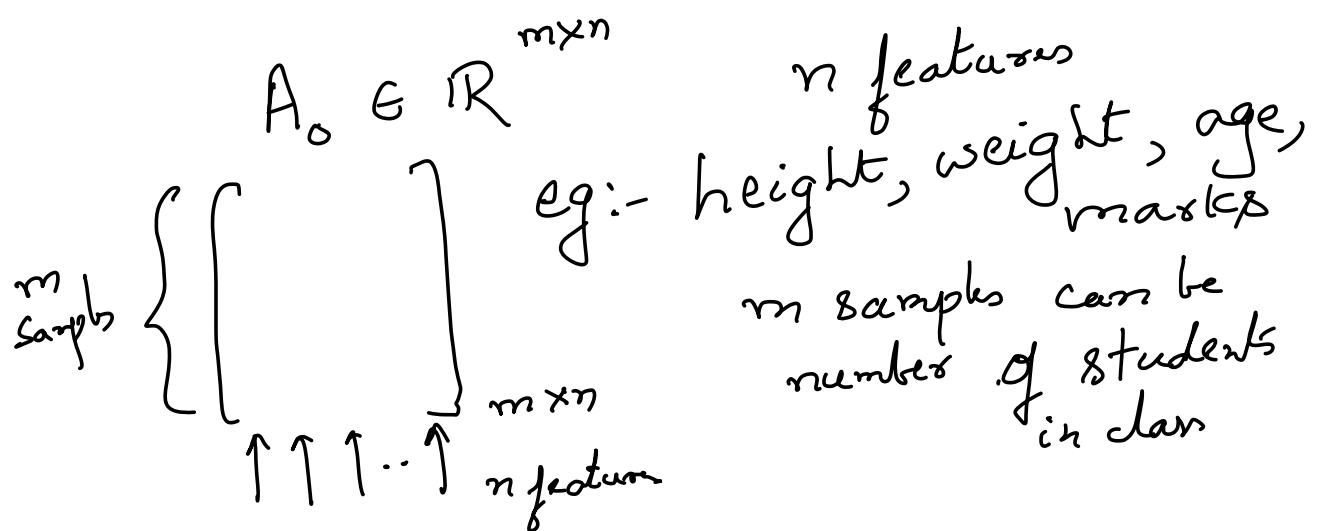
$$\|A - A_k\|_F = \min_{\substack{B \in \mathbb{R}^{m \times n} \\ \text{rank}(B) \leq k}} \|A - B\|_F$$

$$= \sqrt{\sigma_{k+1}^2 + \sigma_{k+2}^2 + \dots + \sigma_r^2}$$

Principal component analysis:

PCA can be thought of orthogonal linear transformation of a given mean centered data matrix A such that transformed directions (vectors) are along the directions of decreasing variances.

Consider a data matrix A_0 .



→ Mean is the average of the data (in each column). Subtract these means of each of these columns of A_0 and reconstruct the data matrix which produces centered matrix \underline{A} .

→ Variance as sum of squares of distances from the mean — along i^{th} column of \underline{A}

$$\text{Var}_i = \frac{1}{m} (\|\underline{a}_i\|_2^2)$$

→ Total variance in the full data is the sum of variances of individual columns.

$$\underline{A} = \begin{bmatrix} | & | & & | \\ \underline{a}_1 & \underline{a}_2 & \dots & \underline{a}_n \\ | & | & & | \end{bmatrix}$$

$$T \propto (\|\underline{a}_1\|_2^2 + \|\underline{a}_2\|_2^2 + \dots + \|\underline{a}_n\|_2^2)$$

$$T \propto \|\underline{A}\|_F^2$$

$$\propto (\underline{\sigma}_1^2 + \underline{\sigma}_2^2 + \dots + \underline{\sigma}_d^2)$$

σ_1^2 accounts for maximum contribution to the total variance, σ_2^2 accounts for next largest contribution to total variance and so on!

The first component \underline{u}_1 (left singular vec) is along the direction of maximum variance, \underline{u}_2 is along the next largest variance and so on!

Why is the above true?

(i) First we seek a direction vector in the feature space i.e. space spanned by " n " features which has maximum variance. [Assume all " n " features are linearly independent]
Let \underline{t}_1 be such a direction and we have $\underline{t}_1 = A \hat{\underline{w}}_1$ and we need to

$$\text{find } \hat{\underline{w}}_1 = \arg \max_{\|\hat{\underline{w}}_1\|=1} \|A \hat{\underline{w}}_1\|_2 \quad [\text{Since we need to}]$$

The above has clearly
a solution with $\hat{\underline{w}}_1 = \underline{v}_1$
the first right singular vector
and $\underline{t}_1 = \sigma_1 \underline{u}_1$ i.e. \underline{u}_1 is the
direction of maximum variance.

(ii) Now we need to find the
direction along the second maximum
variance. For this I need to have

$\underline{t}_2 = A \hat{\underline{w}}_2$ but I need to consider
the action of A on those vectors $\hat{\underline{w}}_2$
which is orthogonal to \underline{v}_1 . We can

denote these vectors by considering
 $\underline{w}_2^\perp = (\underline{I} - \underline{v}_1 \underline{v}_1^T) \hat{\underline{w}}_2$ where $\hat{\underline{w}}_2 \in \mathbb{R}^n$

Hence the problem of seeking the direction of second maximum variance is equivalent to solving

$$\arg \max_{\|\hat{w}_2\|=1} \|\hat{A} \hat{w}_2\|_2 \text{ where } \boxed{\hat{A} = A(I - v_1 v_1^T)}$$

and solution to this problem ~~is the~~ is the second right singular vector i.e. $\hat{w}_2 = v_2$ and $t_2 = \sigma_2 u_2$

where u_2 is the direction of second maximum variance!

and this process can be repeated for directions of next maximum variances.

The key point is that $k < n$
singular vectors explain most of
the data than any other set of
 k vectors. So we can choose

the left singular vectors

$\underline{u}_1, \underline{u}_2, \dots, \underline{u}_k$ as a basis for

k -dimensional subspace closest to
 n -dimensional subspace corresponding
to our m -data points.