

THE SINGULAR VALUE DECOMPOSITION (SVD)

→ ANY  $n \times m$  real matrix  $A$  can be decomposed as:

$$A = U \Sigma V^T, \quad \text{(called the SVD of } A \text{)} \quad (1)$$

$\begin{matrix} \swarrow & \uparrow & \searrow \\ n \times m & n \times n & m \times m \\ \downarrow & \swarrow & \downarrow \\ & \text{diagonal} & \end{matrix}$

→ where  $UU^T = U^T U = I$ ,  $VV^T = V^T V = I$  (such matrices are called **UNITARY**) (2)

→ and  $\Sigma$  is diagonal - it looks like this:  $\begin{matrix} m \\ \sigma_1 \\ \vdots \\ \sigma_n \\ \vdots \\ \sigma_m \end{matrix}$  or  $\begin{matrix} m \\ \sigma_1 \\ \vdots \\ \sigma_m \end{matrix}$  (3)

→ and moreover:  $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq \sigma_m \geq 0$  are all real and positive. (4)  
 →  $\{\sigma_i\}$  are called the **SINGULAR VALUES** of  $A$ .

- NOTE that the SVD is **NOT** the same as eigendecomposition.  
 - (why not? answer this for yourself before moving on)

- Why is the SVD useful? How can we calculate the SVD?  
 - We'll look into these questions below, in due course.

- First, we need to establish some mathematical preliminaries:

- unitary matrices: 
$$\begin{matrix} U^T & & U & & I \\ \left[ \begin{array}{c} \leftarrow \vec{u}_1^T \rightarrow \\ \leftarrow \vec{u}_2^T \rightarrow \\ \vdots \\ \leftarrow \vec{u}_n \rightarrow \end{array} \right] & \left[ \begin{array}{c} \uparrow \\ \vec{u}_1 \\ \vec{u}_2 \\ \dots \\ \vec{u}_n \\ \downarrow \end{array} \right] & = & \left[ \begin{array}{c} 1 \\ \vdots \\ \vdots \\ \vdots \\ 1 \end{array} \right] \end{matrix} \quad (5)$$

- From the above, you can see that:

→  $\vec{u}_1^T \vec{u}_1 = 1 \Rightarrow \|\vec{u}_1\| = 1$ ;  $\vec{u}_2^T \vec{u}_2 = 1$ , etc, ...,  $\vec{u}_n^T \vec{u}_n = 1 \Rightarrow \|\vec{u}_i\| = 1$  (6)

→ also  $\vec{u}_1^T \vec{u}_2 = 0$ ,  $\vec{u}_1^T \vec{u}_3 = 0$ ,  $\vec{u}_2^T \vec{u}_3 = 0$ , etc  $\Rightarrow \vec{u}_i^T \vec{u}_j = 0$  if  $i \neq j$  (7)

→ any set of vectors  $\{\vec{u}_1, \vec{u}_2, \dots, \vec{u}_n\}$  satisfying (6) & (7) is called **ORTHONORMAL**.

→ if  $\|\vec{u}_i\| \neq 1$ , then they are called **ORTHOGONAL** necessarily  
 ↳  $\vec{u}_i$  should not be  $\vec{0}$ , of course.

→ An important property of orthogonal sets of vectors:

→ any set of orthogonal vectors is linearly independent. (8)

→ proof: given an orthogonal set of vectors  $\{\vec{p}_1, \dots, \vec{p}_n\}$

→ suppose they are linearly dependent: i.e.  $\exists$  some  $\alpha_1, \dots, \alpha_n$ , not all zero, s.t. "there exists"

$$\alpha_1 \vec{p}_1 + \alpha_2 \vec{p}_2 + \dots + \alpha_n \vec{p}_n = \vec{0} \quad (9)$$

→ pre-multiply (10) by  $\vec{p}_1^T$ ; then we get

$$\alpha_1 \|\vec{p}_1\|^2 = 0 \quad (\vec{p}_1^T \vec{p}_2, \vec{p}_1^T \vec{p}_3, \text{etc. are all } 0 \text{ due to orthogonality})$$

→ since  $\|\vec{p}_1\|^2 \neq 0$ ,  $\alpha_1 = 0$ .

→ similarly, you can pre-multiply by  $\vec{p}_2^T, \vec{p}_3^T, \dots, \vec{p}_n^T$  to get  $\alpha_2 = \alpha_3 = \dots = \alpha_n = 0$

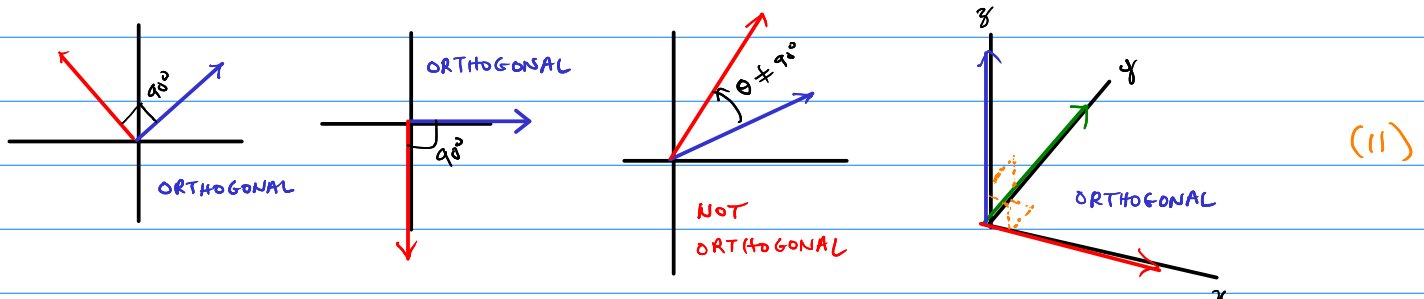
→ hence we have a contradiction to our assumption of linear dependence.

→ Due to (8), any set of orthogonal vectors  $\{\vec{p}_1, \dots, \vec{p}_n\}$ , each in  $\mathbb{R}^n$ , constitutes a basis for  $\mathbb{R}^n$ . (10)

↳ i.e., you can express any  $\vec{x} \in \mathbb{R}^n$  as a linear combination of  $\vec{p}_1, \dots, \vec{p}_n$ .

— Geometrical interpretation of orthogonal vectors:

→ in 2D & 3D, orthogonal vectors are at right angles to each other:



→ in 4D and higher, also the same, though difficult to visualize for most.

→ One reason orthonormal bases are nice is that it is easy to express vectors as a linear combination of such bases.

→ Suppose you have a vector  $\vec{x} \in \mathbb{R}^n$ , and you wish to express it as a linear combination of an orthonormal set of vectors  $\vec{p}_1, \vec{p}_2, \dots, \vec{p}_n \in \mathbb{R}^n$

→ i.e, write  $\vec{x} = \alpha_1 \vec{p}_1 + \alpha_2 \vec{p}_2 + \dots + \alpha_n \vec{p}_n$ ; we would like to find (12)  
the scalars  $\alpha_1, \alpha_2, \dots, \alpha_n$ .

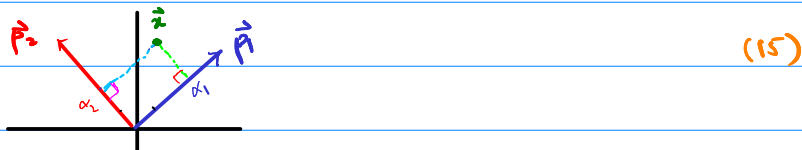
→ It turns out that:

$$\alpha_1 = \vec{p}_1^T \vec{x} \text{ (or } \vec{x}^T \vec{p}_1 \text{)}; \alpha_2 = \vec{p}_2^T \vec{x}; \alpha_3 = \vec{p}_3^T \vec{x}; \dots; \alpha_n = \vec{p}_n^T \vec{x}. \quad (13)$$

→ Proof: left as a small exercise. Hint: pre-multiply (12) by  $\vec{p}_i^T$ .

→ Operations like (13) are called a projection of  $\vec{x}$  onto the basis (14)  
 $\{\vec{p}_1, \vec{p}_2, \dots, \vec{p}_n\}$

→ Projection has a nice geometric interpretation:



→ A very important and useful property of orthonormal bases is the following:

— the vector of projections of any  $\vec{x}$  onto any orthonormal basis has the same norm as  $\vec{x}$ . (15.1)

— Proof: Take any  $\vec{x} \in \mathbb{R}^n$  and project it into some orthonormal basis  $\{\vec{p}_1, \dots, \vec{p}_n\}$ , as in (12) and (13). Then  $\vec{x} = \alpha_1 \vec{p}_1 + \alpha_2 \vec{p}_2 + \dots + \alpha_n \vec{p}_n$ , with  $\alpha_i = \vec{p}_i^T \vec{x}$  (15.2)

$$\Rightarrow \|\vec{x}\|^2 = \vec{x}^T \vec{x} = (\alpha_1 \vec{p}_1 + \alpha_2 \vec{p}_2 + \dots + \alpha_n \vec{p}_n)^T (\alpha_1 \vec{p}_1 + \alpha_2 \vec{p}_2 + \dots + \alpha_n \vec{p}_n) \quad (15.3)$$

$$\Rightarrow \|\vec{x}\|^2 = \alpha_1^2 \|\vec{p}_1\|^2 + \alpha_2^2 \|\vec{p}_2\|^2 + \dots + \alpha_n^2 \|\vec{p}_n\|^2 \quad (\because \vec{p}_i^T \vec{p}_j = 0 \text{ if } i \neq j, \text{ by orthogonality}) \quad (15.4)$$

$$\Rightarrow \|\vec{x}\|^2 = \alpha_1^2 + \alpha_2^2 + \dots + \alpha_n^2 = \|\vec{\alpha}\|^2, \text{ if you define} \quad (15.5)$$

$$\rightarrow \vec{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_n]^T = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix} \quad (15.6)$$

→ Note: (15.6) and (13) can be written compactly as:

$$\rightarrow \vec{\alpha} = P^T \vec{x} \quad (15.7)$$

→ where  $P \triangleq \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ \vec{p}_1 & \vec{p}_2 & \dots & \vec{p}_n \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}$  (15.8)

→ A useful property of  $P$  is that its rows are also orthonormal. (15.85)

→  $P^T P = I$  (from definition of  $P$ , i.e., (15.8) and (5))

$$\rightarrow P^T = P^{-1} \Rightarrow P P^T = P P^{-1} = I \quad (15.9)$$

→  $\begin{bmatrix} \leftarrow \vec{q}_1 \rightarrow \\ \leftarrow \vec{q}_2 \rightarrow \\ \vdots \\ \leftarrow \vec{q}_n \rightarrow \end{bmatrix} \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ \vec{q}_1 & \vec{q}_2 & \dots & \vec{q}_n \\ \downarrow & \downarrow & & \downarrow \end{bmatrix} = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ 0 & & & 1 \end{bmatrix}$  (15.91)

rows of  $P$

$P$   $P^T$

$$\Rightarrow \vec{q}_i^T \vec{q}_j = \begin{cases} 1, & \text{if } i=j \\ 0, & \text{otherwise} \end{cases} \quad (15.92)$$

→ Another useful property of a set of orthonormal vectors  $\{\vec{p}_1, \vec{p}_2, \dots, \vec{p}_n\}$ , for any  $r \leq n$ , is that

$$\underbrace{\|\vec{p}_1\|^2} + \underbrace{\|\vec{p}_2\|^2} + \dots + \underbrace{\|\vec{p}_r\|^2} = \sum_{i=1}^r \sum_{j=1}^n P_{ji}^2 = r. \quad (15.93)$$

— Here,  $P_{ji}$  represents the (row= $j$ , col= $i$ )<sup>th</sup> entry of the matrix (15.94)

$$\begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ \vec{p}_1 & \vec{p}_2 & \dots & \vec{p}_n \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}$$



## OUTER PRODUCTS OF VECTORS AND RANK-1 MATRICES

→ To best appreciate what SVDs can do for us, we need to know about outer products of vectors.

→ Suppose you have two vectors  $\vec{x} \in \mathbb{R}^n$  &  $\vec{y} \in \mathbb{R}^m$ , then

$\vec{x} \vec{y}^T$  is called an outer product (16)

→  $\vec{y} \vec{x}^T$  is also a (different) outer product.

→ outer products are better appreciated when you see them in expanded form:

$$\begin{matrix} \vec{x} & \vec{y}^T \\ \left[ \begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix} \right] & [y_1, y_2, \dots, y_m] \\ \left\| \begin{matrix} n \\ m \end{matrix} \right. & \end{matrix} = \begin{matrix} \xleftarrow{m} & \xrightarrow{n} \\ \left[ \begin{matrix} x_1 y_1 & x_1 y_2 & \dots & x_1 y_m \\ x_2 y_1 & x_2 y_2 & \dots & x_2 y_m \\ \vdots & \vdots & \dots & \vdots \\ x_n y_1 & x_n y_2 & \dots & x_n y_m \end{matrix} \right] & \left. \begin{matrix} \\ \\ \\ \\ \end{matrix} \right\| \\ & \end{matrix} \quad (17)$$

→ so we see that  $\vec{x} \vec{y}^T$  is an  $n \times m$  matrix.

→ not to be confused with  $\vec{x}^T \vec{y}$  or  $\vec{y}^T \vec{x}$  - which are scalars.

→  $\vec{x}^T \vec{y}$  and  $\vec{y}^T \vec{x}$  are called INNER PRODUCTS or DOT PRODUCTS.

→ moreover,  $\vec{x} \vec{y}^T$  is a RANK-1 matrix (18)

- each row of (17) is just a scaled copy of any other one.

- and also, each column of (17) is a scaled copy of any other column.

- this is easy to see if you write (17) like this:

$$\begin{matrix} \vec{x} & \vec{y}^T \\ \left[ \begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix} \right] & [y_1, y_2, \dots, y_m] \\ \left\| \begin{matrix} n \\ m \end{matrix} \right. & \end{matrix} = \begin{matrix} \left[ \begin{matrix} \leftarrow x_1 \vec{y}^T \rightarrow \\ \leftarrow x_2 \vec{y}^T \rightarrow \\ \vdots \\ \leftarrow x_n \vec{y}^T \rightarrow \end{matrix} \right] & = & \left[ \begin{matrix} \uparrow y_1 \vec{x} \\ \uparrow y_2 \vec{x} \\ \dots \\ \uparrow y_m \vec{x} \end{matrix} \right] \\ & & \left. \begin{matrix} \\ \\ \\ \\ \end{matrix} \right\| \end{matrix} \quad (19)$$

— One reason outer products are useful is that matrix multiplication can be expressed as a sum of outer products:

$$\begin{array}{c} \uparrow \\ \downarrow \end{array} \begin{array}{c} \xrightarrow{m} \\ \left[ \begin{array}{c} \uparrow \\ x_1 \\ \downarrow \\ \uparrow \\ x_2 \\ \downarrow \\ \dots \\ \uparrow \\ x_m \\ \downarrow \end{array} \right] \\ \xrightarrow{r} \\ \left[ \begin{array}{c} \leftarrow y_1^T \rightarrow \\ \leftarrow y_2^T \rightarrow \\ \leftarrow \vdots \rightarrow \\ \leftarrow y_m^T \rightarrow \end{array} \right] \\ \downarrow \end{array} \begin{array}{c} \leftarrow \\ \leftarrow \\ \leftarrow \\ \leftarrow \end{array} \end{array} = \underbrace{x_1 y_1^T}_{\substack{\uparrow \\ \text{each one is a } n \times r \\ \text{rank-1 matrix}}} + \underbrace{x_2 y_2^T}_{\substack{\uparrow \\ \text{each one is a } n \times r \\ \text{rank-1 matrix}}} + \dots + \underbrace{x_m y_m^T}_{\substack{\uparrow \\ \text{each one is a } n \times r \\ \text{rank-1 matrix}}} = \sum_{i=1}^m x_i y_i^T \quad (20)$$

— to convince yourself of this: just try a few small examples.

BACK TO SVDs - we are now ready to use the above mathematical machinery to explore (1) in more detail.

→ Rewriting (1) in expanded form:

$$A = \begin{matrix} & \xleftarrow{n} & & \xleftarrow{m} & & \xleftarrow{m} \\ \begin{matrix} \uparrow \\ \uparrow \\ \vdots \\ \uparrow \\ \downarrow \\ \downarrow \\ \vdots \\ \downarrow \\ \downarrow \end{matrix} & \left[ \begin{matrix} \vec{u}_1 & \vec{u}_2 & \dots & \vec{u}_n \end{matrix} \right] & \begin{matrix} \uparrow \\ \uparrow \\ \vdots \\ \uparrow \\ \downarrow \\ \downarrow \\ \vdots \\ \downarrow \\ \downarrow \end{matrix} & \left[ \begin{matrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_m \end{matrix} \right] & \begin{matrix} \leftarrow \\ \leftarrow \\ \vdots \\ \leftarrow \\ \leftarrow \end{matrix} & \left[ \begin{matrix} \vec{v}_1^T \\ \vec{v}_2^T \\ \vdots \\ \vec{v}_m^T \end{matrix} \right] & \begin{matrix} \leftarrow \\ \leftarrow \\ \vdots \\ \leftarrow \\ \leftarrow \end{matrix} \\ \downarrow & U & \downarrow & \Sigma & \downarrow & V^T & \downarrow \\ n & & n & & m & & m \end{matrix} \quad (21)$$

→ Notice that  $U\Sigma$  above expands to:

$$U\Sigma = \begin{matrix} & \xleftarrow{m} & & \\ \begin{matrix} \uparrow \\ \uparrow \\ \vdots \\ \uparrow \\ \downarrow \\ \downarrow \\ \vdots \\ \downarrow \\ \downarrow \end{matrix} & \left[ \begin{matrix} \sigma_1 \vec{u}_1 & \sigma_2 \vec{u}_2 & \dots & \sigma_m \vec{u}_m \end{matrix} \right] & \begin{matrix} \uparrow \\ \uparrow \\ \vdots \\ \uparrow \\ \downarrow \\ \downarrow \\ \vdots \\ \downarrow \\ \downarrow \end{matrix} \\ \downarrow & & \downarrow & \\ n & & n & \end{matrix} \quad (22)$$

- an interesting observation:

- from (22), we see that  $\vec{u}_{m+1}, \dots, \vec{u}_n$  are not even involved.
  - which means that they don't affect (1) or (21).
  - so in fact,  $\vec{u}_{m+1}$  to  $\vec{u}_n$  can be any vectors that complete  $\vec{u}_1, \dots, \vec{u}_m$  into an orthonormal basis for  $\mathbb{R}^n$ .
- (23)

→ In view of (22), (21) becomes:

$$\rightarrow A = U\Sigma V^T = \begin{matrix} & \xleftarrow{m} & & \xleftarrow{m} & & \\ \begin{matrix} \uparrow \\ \uparrow \\ \vdots \\ \uparrow \\ \downarrow \\ \downarrow \\ \vdots \\ \downarrow \\ \downarrow \end{matrix} & \left[ \begin{matrix} \sigma_1 \vec{u}_1 & \sigma_2 \vec{u}_2 & \dots & \sigma_m \vec{u}_m \end{matrix} \right] & \begin{matrix} \uparrow \\ \uparrow \\ \vdots \\ \uparrow \\ \downarrow \\ \downarrow \\ \vdots \\ \downarrow \\ \downarrow \end{matrix} & \left[ \begin{matrix} \leftarrow \\ \leftarrow \\ \vdots \\ \leftarrow \\ \leftarrow \end{matrix} \right] & \begin{matrix} \leftarrow \\ \leftarrow \\ \vdots \\ \leftarrow \\ \leftarrow \end{matrix} \\ \downarrow & U\Sigma & \downarrow & V^T & \downarrow \\ n & & n & & m \end{matrix} \quad (24)$$

→ now we can use the outer product form of matrix multiplication (20) to express (24):

$$\rightarrow A = U \Sigma V^T = \begin{matrix} \begin{matrix} \uparrow & \uparrow & \dots & \uparrow \\ \sigma_1 \vec{u}_1 & \sigma_2 \vec{u}_2 & \dots & \sigma_m \vec{u}_m \\ \downarrow & \downarrow & & \downarrow \end{matrix} \\ U \Sigma \end{matrix} \begin{matrix} \begin{matrix} \leftarrow \vec{v}_1^T \rightarrow \\ \leftarrow \vec{v}_2^T \rightarrow \\ \vdots \\ \leftarrow \vec{v}_m^T \rightarrow \end{matrix} \\ V^T \end{matrix} = \sum_{i=1}^m \sigma_i \vec{u}_i \vec{v}_i^T \quad (25)$$

rank-1 outer  $n \times m$  product matrices

singular values (scalars in decreasing order)

→ Consider  $\vec{u}_i \vec{v}_i^T$  for any  $i \in 1, \dots, m$ .

→ recall that  $\|\vec{u}_i\| = \|\vec{v}_i\| = 1$ , because they are members of orthonormal bases.

→ let's write  $\vec{u}_i \vec{v}_i^T$  in expanded form, as in (19):

$$\rightarrow R_i \triangleq \vec{u}_i \vec{v}_i^T = \begin{bmatrix} v_{i1} \vec{u}_i & v_{i2} \vec{u}_i & \dots & v_{im} \vec{u}_i \\ \downarrow & \downarrow & & \downarrow \end{bmatrix} \quad (26)$$

→ Now, we would like to find the Frobenius norm of  $R_i$ :

→ the Frobenius norm  $\|R_i\|_F$  is simply the square root of the sum of all the entries of  $R_i$ . (27)

→ the sum of the squares of all the entries  $\|R_i\|_F^2$  is the sum of the squares of the norms of each column.

$$\rightarrow \text{i.e., } \|R_i\|_F^2 = v_{i1}^2 \|\vec{u}_i\|^2 + v_{i2}^2 \|\vec{u}_i\|^2 + \dots + v_{im}^2 \|\vec{u}_i\|^2$$

$$\boxed{\text{Frobenius norm of } R_i}^2 = (v_{i1}^2 + v_{i2}^2 + \dots + v_{im}^2) \|\vec{u}_i\|^2 = \|\vec{v}_i\|^2 \|\vec{u}_i\|^2$$

$$\Rightarrow \|R_i\|_F^2 = 1 \quad (\text{due to orthonormality of } \vec{u}_i \text{ and } \vec{v}_i) \quad (28)$$

→ what (28) and (25) are telling us is :

→ any matrix  $A$  can be decomposed into a sum of rank 1 matrices, each of unit Frobenius norm, weighted (29) by the singular values.

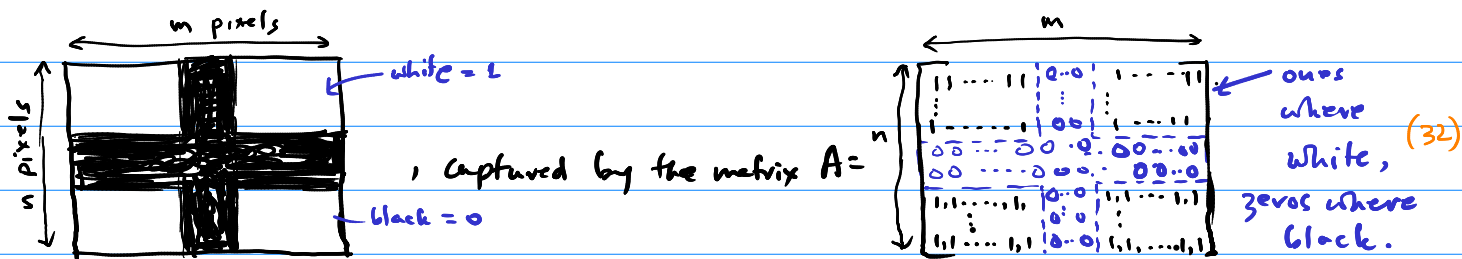
→ i.e., we can think of  $\sigma_1 \vec{u}_1 \vec{v}_1^T$  as the "most important rank-1 component" of  $A$ , since  $\sigma_1 \geq \sigma_2, \sigma_3, \dots, \sigma_m$  (30)

→  $\sigma_2 \vec{u}_2 \vec{v}_2^T$  is the "2nd most important rank-1 component" (31)

→ similarly for  $\sigma_3 \vec{u}_3 \vec{v}_3^T, \sigma_4 \vec{u}_4 \vec{v}_4^T$  and so on.

- The uses of (29) - (31) are perhaps most easily seen by applying SVD to images

- Suppose you have a simple B&W image like this :



-  $A$  is a rank 1 matrix, since it can be written as

$$A = \begin{bmatrix} 1 & \dots & 1 & 0 & \dots & 0 & 1 & \dots & 1 \\ \vdots & & \vdots & & & & \vdots & & \vdots \\ 1 & \dots & 1 & 0 & \dots & 0 & 1 & \dots & 1 \end{bmatrix} \rightarrow \vec{q}^T \quad (33)$$

significant

where  $\sigma_1 = \|\vec{p}\| \|\vec{q}\|$  (34)

$\vec{u}_1 = \frac{\vec{p}}{\|\vec{p}\|}, \vec{v}_1 = \frac{\vec{q}}{\|\vec{q}\|}$

- an SVD of  $A$  would have only one <sup>significant</sup> component,  $\sigma_1 \vec{u}_1 \vec{v}_1^T$ , where

→  $\sigma_2, \sigma_3, \dots, \sigma_m$  would all be 0.

- More generally, images composed of "simple rectangular patterns" will have only a few significant singular values.
  - see the flag examples in the slides
- Images that are not so simple - especially if they have non-rectangular shapes in them - typically have a larger number of significant singular values.
- the rank-1 matrices  $\vec{u}_1 \vec{v}_1^T, \vec{u}_2 \vec{v}_2^T, \dots$ , etc., often have some visual interpretation or significance - they are sometimes called "features"
- (Perhaps more surprisingly) almost every image that makes instinctive sense to human beings - eg, faces, pictures of scenery, etc. - can be well approximated visually if you approximate it using only the most significant singular values.
  - More precisely, suppose you have an image matrix  $A \in \mathbb{R}^{n \times m}$ 
    - take an SVD of it (recall (25)):

$$\rightarrow A = U \Sigma V^T = \begin{matrix} \begin{matrix} \uparrow & \uparrow & \dots & \uparrow \\ \sigma_1 \vec{u}_1 & \sigma_2 \vec{u}_2 & \dots & \sigma_m \vec{u}_m \\ \downarrow & \downarrow & & \downarrow \end{matrix} \begin{matrix} \xleftarrow{m} \\ \left[ \begin{matrix} \leftarrow \vec{v}_1^T \rightarrow \\ \leftarrow \vec{v}_2^T \rightarrow \\ \vdots \\ \leftarrow \vec{v}_m^T \rightarrow \end{matrix} \right] \\ \xleftarrow{m} \\ \downarrow \\ V^T \end{matrix} \end{matrix} \quad (35)$$

$U \Sigma$

- Suppose you decide that only the first  $r$  singular values are significant, with  $r < m$ .
- for images, a good rule of thumb for deciding  $r$  is to ensure that  $\sigma_r \approx \frac{10^{-2}}{\text{or more}} \sigma_1$ , i.e., that all singular values 2 orders of magnitude, less than  $\sigma_1$  are negligible.
- then, approximate (35), keeping only the first  $r$  columns of  $U$  and rows of  $V$ :

$$\rightarrow A \approx A_r \triangleq \begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ \sigma_1 \vec{u}_1 & \sigma_2 \vec{u}_2 & \dots & \sigma_r \vec{u}_r \\ \downarrow & \downarrow & \dots & \downarrow \end{bmatrix} \begin{bmatrix} \leftarrow \vec{v}_1^T \rightarrow \\ \leftarrow \vec{v}_2^T \rightarrow \\ \vdots \\ \leftarrow \vec{v}_r^T \rightarrow \end{bmatrix} \quad (36)$$

$\rightarrow$  if  $r \ll m$  (as is often the case), and the approximate image still looks good, then we can achieve a way to represent the image with less data than storing every pixel — i.e., image compression.

$\rightarrow$  example: see Michel's image in the slides, e.g., for  $r=50$ .

— The intuition that you can "throw away" all the terms  $\sigma_i \vec{u}_i \vec{v}_i^T$  that have relatively small  $\sigma_i$  (compared to  $\sigma_1$ , or the top few singular values) can be concretized mathematically, as follows:

— The difference in Frobenius norm (27) between  $A$  &  $A_r$  (in (36))

$$\text{is } \sqrt{\sigma_{r+1}^2 + \sigma_{r+2}^2 + \dots + \sigma_m^2} \quad (36.1)$$

— Proof: (36) is the same as:  $A_r = \sum_{i=1}^r \sigma_i \vec{u}_i \vec{v}_i^T$ , with  $r < m$ . (36.2)

$\rightarrow$  Since  $A = \sum_{i=1}^m \sigma_i \vec{u}_i \vec{v}_i^T$  (by (25)), we have  $A - A_r = \sum_{i=r+1}^m \sigma_i \vec{u}_i \vec{v}_i^T$  (36.3)

$\rightarrow$  We now write out the  $k^{\text{th}}$  column of  $(A - A_r)$ , where  $1 \leq k \leq m$ :

$\rightarrow$  Let  $\vec{q}_k \in \mathbb{R}^n$  be this column, then  $\vec{q}_k = \sum_{i=r+1}^m \sigma_i v_{ki} \vec{u}_i$ , (36.4)

where  $v_{ki}$  denotes the  $k^{\text{th}}$  element of the column vector  $\vec{v}_i$ .

$\rightarrow$  illustration: let  $k=2$  (i.e., the 2<sup>nd</sup> column):

$$\begin{aligned} \rightarrow 2^{\text{nd}} \text{ col of } \sigma_1 \vec{u}_1 \vec{v}_1^T &= \sigma_1 v_{21} \vec{u}_1 \\ \rightarrow 2^{\text{nd}} \text{ col of } \sigma_2 \vec{u}_2 \vec{v}_2^T &= \sigma_2 v_{22} \vec{u}_2 \\ \vdots \\ \rightarrow 2^{\text{nd}} \text{ col of } \sigma_r \vec{u}_r \vec{v}_r^T &= \sigma_r v_{2r} \vec{u}_r \\ \vdots \\ \rightarrow 2^{\text{nd}} \text{ col of } \sigma_m \vec{u}_m \vec{v}_m^T &= \sigma_m v_{2m} \vec{u}_m \end{aligned} \quad \left. \begin{array}{l} \text{using the notation of (15.94)} \\ \end{array} \right\} (36.45)$$

$$\Rightarrow \vec{q}_2 \equiv 2^{\text{nd}} \text{ column of } \sum_{i=r+1}^m \sigma_i \vec{u}_i \vec{v}_i^T = \sum_{i=r+1}^m \sigma_i v_{2i} \vec{u}_i$$

→ Since  $\{\vec{u}_i\}$  are orthonormal, then by (15.1) - (15.6), we have

$$\rightarrow \|\vec{q}_k\|^2 = \sum_{i=r+1}^n \sigma_i^2 v_{ki}^2 \quad (36.5)$$

→ If we sum (36.5) over all  $k = 1, \dots, m$ , (i.e., all the columns of  $A - A_r$ ),

we get

$$\rightarrow \sum_{k=1}^m \|\vec{q}_k\|^2 = \sum_{k=1}^m \sum_{i=r+1}^n \sigma_i^2 v_{ki}^2 = \sum_{i=r+1}^n \sigma_i^2 \sum_{k=1}^m v_{ki}^2 \quad (36.6)$$

→ this is just  $\|\vec{u}_i\|^2 = 1$

$$\rightarrow \sum_{k=1}^m \|\vec{q}_k\|^2 = \sum_{i=r+1}^n \sigma_i^2 \quad (36.7)$$

→ Finally, we need to simply observe that:

$$\rightarrow \|A - A_r\|_F^2 = \text{sum of the norm-squares of the columns} = \sum_{k=1}^m \|\vec{q}_k\|^2 \quad (36.8)$$

— Putting together (36.8) and (36.7), we get what we set out to prove, i.e.,

$$\|A - A_r\|_F = \sqrt{\sum_{i=r+1}^n \sigma_i^2} \quad (36.9)$$

— This result (36.1) or (36.9) is very significant, since it exactly quantifies the error (in Frobenius norm) between the original matrix  $A$  and an approximation truncated to the top  $r$  singular values.

→ An immediate corollary is that the relative error is:

$$\epsilon_r \triangleq \frac{\|A - A_r\|_F}{\|A\|_F} = \frac{\sqrt{\sum_{i=r+1}^n \sigma_i^2}}{\sqrt{\sum_{i=1}^n \sigma_i^2}} = \frac{1}{\sqrt{1 + \frac{\sum_{i=1}^r \sigma_i^2}{\sum_{i=r+1}^n \sigma_i^2}}} \quad (36.95)$$

— thus, if the squared sums of the truncated SVs is very much smaller than that of the non-truncated ones, the error is very small.

— in the slides: show the error for the various examples.



## SVDs for IDENTIFYING "FEATURES" IN TABULAR DATA

— SVDs are also useful for analysing tabular data.

— for example, suppose we have a matrix representing grades in an exam:

$A \in \mathbb{R}^{6 \times 8}$

(more generally, data samples = rows)

	<u>QUESTIONS</u> (or more generally: types of data = columns)							
	Q1	Q2	Q3	Q4	Q5	Q6	Q7	Q8
S1	1	2	1	3	2	1	0	1
S2	10	9	1	10	10	2	10	8
S3	5	4	2	6	5	0	3	2
S4	9	8	3	9	10	1	10	9
S5	0	0	0	1	3	2	2	2
S6	10	10	3	10	10	2	10	10

(37)

— we run an SVD on the data, & get:

$$\rightarrow A = \sigma_1 \vec{u}_1 \vec{v}_1^T + \sigma_2 \vec{u}_2 \vec{v}_2^T + \dots + \sigma_m \vec{u}_m \vec{v}_m^T \quad (38)$$

— unlike for the case of images, the "features"  $\vec{u}_1, \vec{v}_1^T, \vec{u}_2, \vec{v}_2^T$  are often less intuitive for general data.

— However, the "top columns" of  $U$  (i.e.,  $\vec{u}_1, \vec{u}_2, \dots$ ) and the "top rows" of  $V^T$  (i.e.,  $\vec{v}_1^T, \vec{v}_2^T, \dots$ ) can often be interpreted meaningfully, as follows: (39)

→  $\vec{u}_1$  (which is of length  $n=6$  in the above example), can be thought of as the "strongest column (or data type) feature" (or Question feature in the above example). (40)

→  $\vec{u}_2$  would be the 2<sup>nd</sup> strongest data type/col. feature.

→ with this interpretation, you can project each of the data type (Question) columns onto  $\vec{u}_1$  through  $\vec{u}_n$ .

For example, the  $Q_4$  column would be

$$\begin{array}{c} Q_4 \\ \left[ \begin{array}{c} 3 \\ 10 \\ 6 \\ 9 \\ 1 \\ 10 \end{array} \right] = \alpha_{41} \vec{u}_1 + \alpha_{42} \vec{u}_2 + \dots + \alpha_{46} \vec{u}_6, \end{array} \quad (41)$$

→ where  $\alpha_{4i} = \vec{u}_i^T \begin{array}{c} Q_4 \\ \left[ \begin{array}{c} 3 \\ 10 \\ 6 \\ 9 \\ 1 \\ 10 \end{array} \right] \end{array}$  is the strength of the  $i^{\text{th}}$  data type (Question) feature for  $Q_4$ . (42)

→ we can find all the  $\alpha$ s (42) for all the columns of the matrix via

$$\rightarrow \hat{A} = U^T A, \quad (43)$$

→ where  $\hat{A} = \begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ \alpha_1 & \alpha_2 & \dots & \alpha_m \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}$  (44)

→ and the  $k^{\text{th}}$  column of  $A$  can be recovered as  $U \vec{\alpha}_k$  (45)

- i.e., the entries of  $\vec{\alpha}_k$  ( $\alpha_{1k}, \alpha_{2k}, \dots, \alpha_{mk}$ ) are the projections in (42) and (41)

→ so, in fact,  $A$  can be recovered as  $A = U \hat{A}$  (46)

→ Now, suppose we choose to retain only the 1st  $r < m$  "column features", i.e., we project onto only the first  $r$  columns of  $U$ :

$U_r \triangleq \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ \vec{u}_1 & \vec{u}_2 & \dots & \vec{u}_r \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}$ ,  $\hat{A}_r \triangleq U_r^T A = \begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ \hat{\alpha}_1 & \hat{\alpha}_2 & \dots & \hat{\alpha}_m \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}$  (47)

$n$   $r \times m$   $r \times m$   $n \times m$   $r$

$r$   $m$   $m$

→ and we reconstruct  $A$  from  $\hat{A}_r$  similarly, using only the first  $r$  columns of  $U$ , i.e.,

$A_r \triangleq U_r \hat{A}_r$  (48)

→ what is the re-construction error  $\|A - A_r\|_F$ ?

→ It is straightforward to show that this error is:

$\|A - A_r\|_F = \epsilon_r = \sqrt{\sum_{i=r+1}^m \sigma_i^2}$  (exactly as in (30.9)) (49)  
- the proof is the same

→ (49), interpreted in words, is:

→ If we represent all the columns of  $A$  by projection onto only the first  $r$  "column features" (eg, Questions or types of data), then we can recover all the columns of  $A$  with total mean-squared (Frobenius) error given by (49).

→ i.e.,  $\sigma_{r+1}, \dots, \sigma_m$  are "negligible" compared to  $\sigma_1, \dots, \sigma_r$ , this error can be small (50)  
→ in the sense that  $\epsilon_r$  in (39.95) is  $\ll 1$

→ So again, if we can truncate with small error to  $r$  column "column features", we can achieve good data compression.

→ Note: we can, similarly, define row features, in exact analogy with column features, as above.

→ left as an exercise

→ if you truncate to  $r$  row features, what is the Frobenius <sup>norm</sup> error?

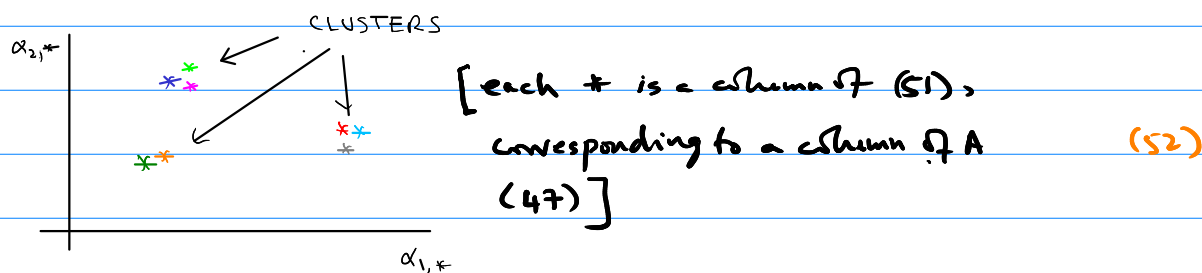
## CLUSTERING USING COLUMN (OR ROW) FEATURES

→ The projections on column features (eg, in (41)-(44), and (47)), if plotted against each other graphically, can reveal meaningful clusters in the data.

→ For example, consider  $\hat{A}_r$  in (47):

$$\hat{A}_r = \begin{bmatrix} \uparrow & \uparrow & \uparrow \\ \hat{\alpha}_1 & \hat{\alpha}_2 & \dots & \hat{\alpha}_m \\ \downarrow & \downarrow & & \downarrow \end{bmatrix} = \begin{bmatrix} \alpha_{11} & \alpha_{12} & \dots & \alpha_{1m} \\ \alpha_{21} & \alpha_{22} & \dots & \alpha_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_{r1} & \alpha_{r2} & \dots & \alpha_{rm} \end{bmatrix} \quad (51)$$

→ Plot the first 2 entries of each column of (51) as a 2-D point:



→ The clusters reveal that the corresponding columns of  $A$  have similar values of "column features" (since their projections onto the features are similar).

→ For our example in (37), columns 3 & 6 ( $Q_3$  &  $Q_6$ ) might fall in a cluster, since most students did poorly in both.

→ See the slides for more precise examples.



—  $\tilde{A}$  as defined in (55) is the same data, but with the mean of each column subtracted from that column (55.1)

— i.e., each column of  $\tilde{A}$  has zero mean, i.e.,  
 $[1, 1, 1, \dots, 1] \tilde{A} = \vec{0}^T$  (55.2)

—  $\tilde{A}$  is sometimes called the mean-centered data matrix

→ now, from the matrix

$$S \stackrel{m \times m}{=} \frac{1}{n} \tilde{A}^T \tilde{A} \quad (56)$$

$\uparrow$   $\nwarrow$   
 $m \times n$   $n \times m$

→  $S$  is called the covariance matrix (of  $A$ , or  $\tilde{A}$ )

— Example: Suppose  $A = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}$ ; (57.1)

$$\mu_1 = \frac{1+2+3}{3} = 2; \quad \mu_2 = \frac{4+5+6}{3} = 5$$

$$\tilde{A} = \begin{bmatrix} -1 & -1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix}; \quad (57.2)$$

$$S = \frac{1}{3} \tilde{A}^T \tilde{A} = \frac{1}{3} \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} \begin{bmatrix} -1 & -1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} \quad (57.3)$$

## — PROPERTIES OF COVARIANCE MATRICES

—  $S$  is square and symmetric (i.e.,  $S^T = S$ ) (58)

— both follow immediately from the definition (56)

— the diagonal entries of  $S$  are greater than or equal to 0 (59)

— let  $\vec{a}_i$  be the  $i^{\text{th}}$  column of  $\tilde{A}$ , then the  $i^{\text{th}}$  diagonal entry of  $S$  is  $\frac{1}{n} \vec{a}_i^T \vec{a}_i = \frac{1}{n} \|\vec{a}_i\|^2 \geq 0$

— Suppose we write  $S$  as

$$\rightarrow S = \begin{bmatrix} s_1^2 & s_{12} & s_{13} & \dots & s_{1m} \\ s_{21} & s_2^2 & s_{23} & \dots & s_{2m} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ s_{m1} & s_{m2} & \dots & \dots & s_m^2 \end{bmatrix} \quad (60)$$

$$\rightarrow \text{then } s_{ij} = s_{ji}, \forall i \neq j \quad (61)$$

↳ due to symmetry, i.e., (58)

$$\rightarrow \text{and } |s_{ij}| \leq s_i s_j, \forall i \neq j \quad (62)$$

→ proof:

$$\rightarrow \text{Let } \vec{a}_i \text{ and } \vec{a}_j \text{ be the } i^{\text{th}} \text{ and } j^{\text{th}} \text{ columns of } \tilde{A}, \text{ respectively.} \quad (62.1)$$

$$\rightarrow \text{then } s_{ij} = \frac{\vec{a}_i^T \vec{a}_j}{n}, \text{ from definition of } S \text{ (56).} \quad (62.2)$$

$$\rightarrow \text{we also have: } s_i^2 = \frac{1}{n} \vec{a}_i^T \vec{a}_i = \frac{\|\vec{a}_i\|^2}{n}, \text{ and } s_j^2 = \frac{1}{n} \vec{a}_j^T \vec{a}_j = \frac{\|\vec{a}_j\|^2}{n} \quad (62.3)$$

→ The famous Cauchy-Schwarz inequality states that for any 2 vectors  $\vec{x}, \vec{y} \in \mathbb{R}^n$ :

$$\|\vec{x}^T \vec{y}\|^2 \leq \|\vec{x}\|^2 \|\vec{y}\|^2 \quad (62.4)$$

→ Applying it to  $\vec{a}_i$  and  $\vec{a}_j$ , we get:

$$\|\vec{a}_i^T \vec{a}_j\|^2 \leq \|\vec{a}_i\|^2 \|\vec{a}_j\|^2 \quad (62.5)$$

$$\Rightarrow \frac{\|\vec{a}_i^T \vec{a}_j\|^2}{n^2} \leq \frac{\|\vec{a}_i\|^2}{n} \frac{\|\vec{a}_j\|^2}{n} \quad (62.6)$$

$$\Rightarrow |s_{ij}|^2 \leq s_i^2 s_j^2 \Rightarrow |s_{ij}| \leq s_i s_j \quad (62.7)$$

→ The diagonal entries of  $S$  (i.e.,  $s_i^2$ ) are the VARIANCES of the (62-8) columns of  $A$

→ Follows directly from the definitions of  $S$  and of variance:

→ The variance of the  $k^{\text{th}}$  column of  $A$  is defined as:

$$\rightarrow \sigma_k^2 \triangleq \frac{\|\vec{a}_k\|^2}{n}, \text{ where } \vec{a}_k \text{ is the } k^{\text{th}} \text{ column of } \tilde{A} \quad (62-81)$$

→ which is the same as  $s_i^2$  (by (62-3))



## — FROM COVARIANCE MATRICES TO CORRELATION MATRICES

— The CORRELATION MATRIX (OF  $A$  OR  $\tilde{A}$ ) is defined to be :

$$R \triangleq \begin{bmatrix} 1 & \frac{s_{12}}{s_1 s_2} & \frac{s_{13}}{s_1 s_3} & \dots & \frac{s_{1m}}{s_1 s_m} \\ \frac{s_{21}}{s_2 s_1} & 1 & \frac{s_{23}}{s_2 s_3} & \dots & \frac{s_{2m}}{s_2 s_m} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{s_{m1}}{s_m s_1} & \dots & \dots & \dots & 1 \end{bmatrix}, \text{ i.e., } R_{ij} = \frac{s_{ij}}{s_i s_j} \quad (63)$$

— Due to (62), we immediately see that the off-diagonal terms of  $R$  are  $\leq 1$ . (64)

—  $R_{ij} \triangleq \frac{s_{ij}}{s_i s_j}$  is called the correlation between the  $i^{\text{th}}$  and  $j^{\text{th}}$  columns of  $A$  (or  $\tilde{A}$ ) (65)

— Example: continuing from (57.1) and (57.3), we see immediately that  $R = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ , i.e., the correlation between the 2 columns is 1. (66)

— Another example: if  $A = \begin{bmatrix} 1 & 6 \\ 2 & 5 \\ 3 & 4 \end{bmatrix}$ ,  $R = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$  (67)

— here, the correlation between the 2 columns is  $-1$ .

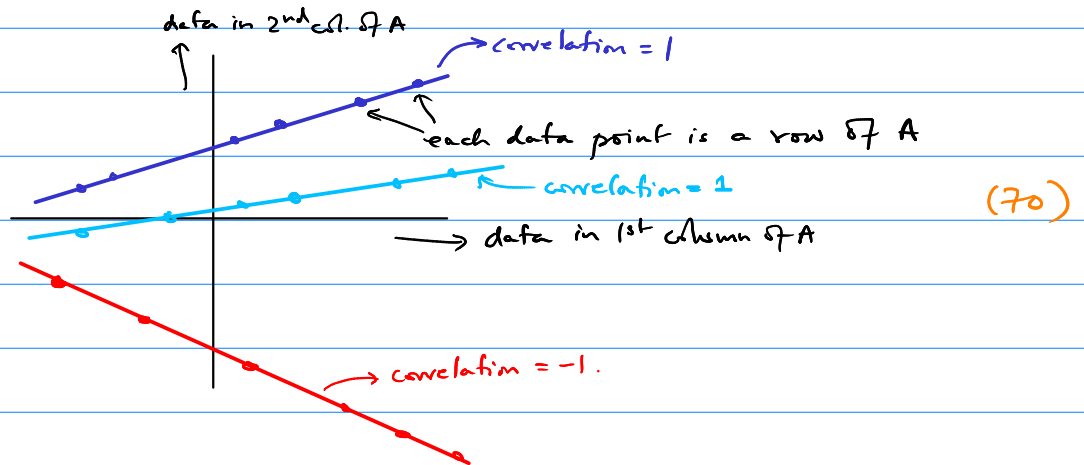
— Correlations have a "nice" property:

→ if every column of  $\tilde{A}$  is simply a scaled version of a single vector  $\vec{a}$ , then every entry of  $R$  is either 1 or  $-1$ . (68)

→ proof: left as an exercise

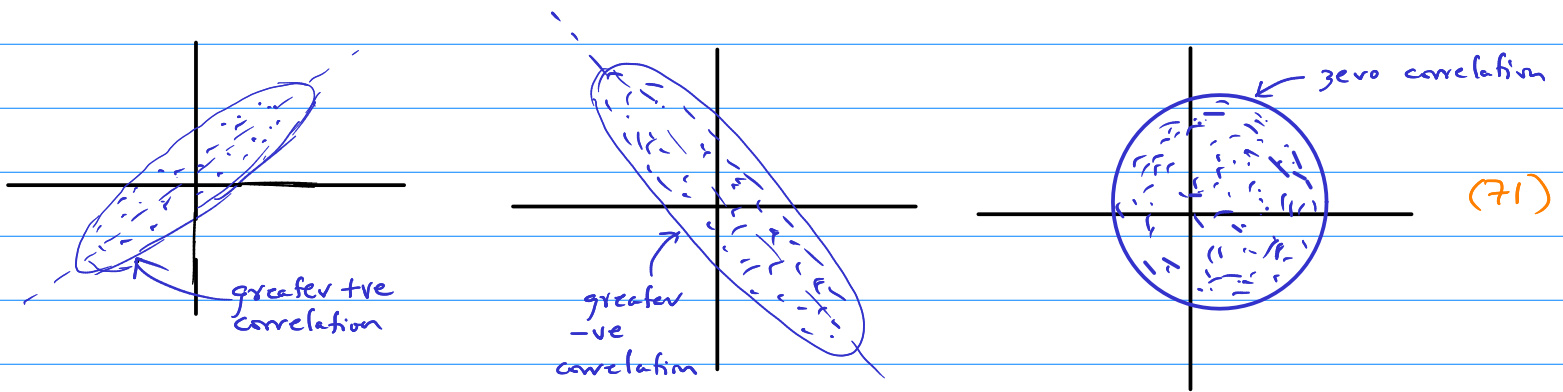
— geometrical interpretation of (68) in 2D:

— If the data in  $A$  lie on a straight line, then the correlation is 1 if the straight line has +ve slope, and -1 if it has -ve slope. (69)

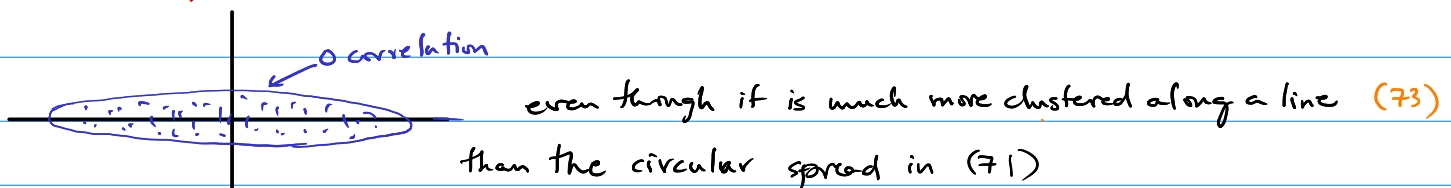


— if the data is more scattered — but still clustered around a line —

correlations still can give you a rough idea how closely clustered around a line the data is. The following figures provide a rough guide, assuming that the samples are uniformly spread within the bounding shapes. (71)



— Correlations do not, however, provide a clear idea of the tilt of the line around which the data is clustered, nor even a reliable notion of the spread of the data. For example, the following also has zero correlation:



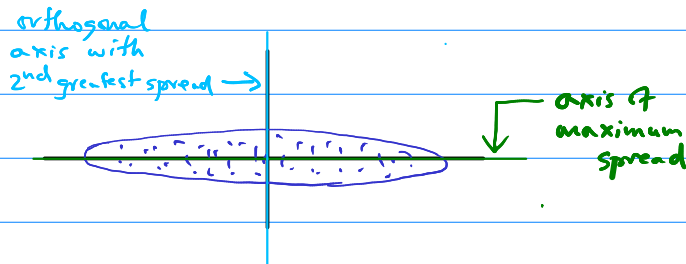
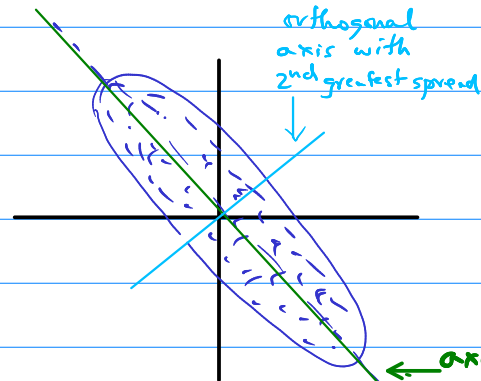
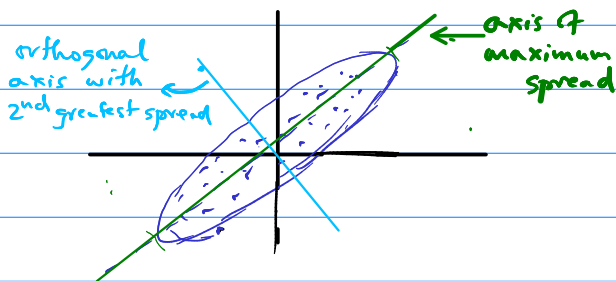
— There are more examples in the slides.

— IT IS IN THIS CONTEXT THAT PCA IS MUCH MORE INFORMATIVE :

— PCA identifies the "main axes" around which the data clusters (74)

— it provides quantitative information about the spread of the data along each of these main axes.

— For example, PCA will correctly identify these axes in the following examples :



(75)

— PCA is basically nothing more than eigen-decomposition of the covariance matrix  $S$  in (60) (76)

— because the covariance matrix has special structure and properties, so do the eigenvectors and eigenvalues — which is what results in the ability of PCA to find the axes corresponding to the greatest data spreads (77)

— To understand PCA properly, we first need to derive certain properties of (real) symmetric and co-variance matrices:

(not necessarily in the form  $\tilde{A}^T \tilde{A}$ )

→ If  $S$  is any (real) symmetric matrix, then (78)

1. The eigenvalues of  $S$  are all real. (79)

— proof: if  $\lambda$  is an eigenvalue of  $S$ , then  $S\vec{p} = \lambda\vec{p}$  (79.1)

for some  $\vec{p}$  (eigenvector of  $S$ )

→ since  $S$  is real, we have  $S\vec{p} = \overline{\lambda}\overline{\vec{p}}$  (79.2)

→ since  $S$  is symmetric, we have  $\vec{p}^T S = \lambda\vec{p}^T$  (79.3)

∴  $\vec{p}^T S \vec{p} = \overline{\lambda} \vec{p}^T \vec{p} = \overline{\lambda} \|\vec{p}\|^2$  (by 79.2) (79.4)

$= \lambda \vec{p}^T \vec{p} = \lambda \|\vec{p}\|^2$  (by 79.3) (79.5)

⇒  $\lambda = \overline{\lambda}$  ⇒  $\lambda$  is real

2. A set of real eigenvectors  $\{\vec{p}_1, \vec{p}_2, \dots, \vec{p}_m\}$  of  $S$  can be found (80)

→ Pf: Suppose we have some set of eigenvectors  $\{\vec{b}_1, \dots, \vec{b}_m\}$  of

$S$ , not necessarily real. Then  $S\vec{b}_i = \lambda_i \vec{b}_i$ . (80.1)

→ Let  $\vec{b}_i = \vec{p}_i + j\vec{q}_i$ , where  $\vec{p}_i$  and  $\vec{q}_i$  are the real and imaginary parts of  $\vec{b}_i$ . (80.2)

→ Then  $S\vec{b}_i = S\vec{p}_i + jS\vec{q}_i = \lambda_i(\vec{p}_i + j\vec{q}_i) = \lambda_i\vec{p}_i + j(\lambda_i\vec{q}_i)$  (80.3)

→ Since  $\lambda_i$  is real (from (79)), we have

$S\vec{p}_i = \lambda_i\vec{p}_i$ , i.e.  $\{\vec{p}_1, \dots, \vec{p}_m\}$ , all real, are eigenvectors of  $S$ . (80.4)

3. The eigenvectors in (80) are orthogonal (81)

- Pf:  $S\vec{p}_i = \lambda_i \vec{p}_i$ ,  $i=1, \dots, m$ , with  $\lambda_i$  &  $\vec{p}_i$  real ((79) & (80))

→ from symmetry of  $S$ , we have:

$$\vec{p}_j^T S = \lambda_j \vec{p}_j^T, \text{ for } j=1, \dots, m \quad (82.1)$$

$$\rightarrow \lambda_j \vec{p}_j^T \vec{p}_i = \vec{p}_j^T S \vec{p}_i = \vec{p}_j^T \lambda_i \vec{p}_i = \lambda_i \vec{p}_j^T \vec{p}_i \quad (82.2)$$

$$\rightarrow \lambda_j \vec{p}_j^T \vec{p}_i = \lambda_i \vec{p}_j^T \vec{p}_i \Rightarrow \vec{p}_j^T \vec{p}_i (\lambda_j - \lambda_i) = 0 \quad (82.3)$$

→ If we assume that all the eigenvalues are distinct, then  $\lambda_i - \lambda_j \neq 0$  if  $i \neq j$ , hence  $\vec{p}_j^T \vec{p}_i = 0$  if  $i \neq j$  (82.4)

→ Even if we have repeated eigenvalues, it can be shown (though we won't do so here) that  $\vec{p}_j$  and  $\vec{p}_i$  can be found such that they are orthogonal (82.5)

4. A real orthonormal set of eigenvectors of  $S$  can be found. (83)

Pf: from (80), we have  $\{\vec{p}_1, \dots, \vec{p}_m\}$  real, and from (82), they are orthogonal.

$$\rightarrow \text{then } \left\{ \frac{\vec{p}_1}{\|\vec{p}_1\|}, \frac{\vec{p}_2}{\|\vec{p}_2\|}, \dots, \frac{\vec{p}_m}{\|\vec{p}_m\|} \right\} \text{ are orthonormal} \quad (83.1)$$

→ we can then just redefine  $\{\vec{p}_1, \dots, \vec{p}_m\}$  in (80) to be orthonormal (83.4)

→ Now, if  $S$  is a real symmetric matrix that can be written in the

$$\text{form } S = B^T B, \text{ for some matrix } B \quad (84)$$

→ then the eigenvalues of  $S$  are all  $\geq 0$  (85)

$$\rightarrow \text{Pf: } S\vec{p} = \lambda \vec{p} \Rightarrow B^T B \vec{p} = \lambda \vec{p} \quad (85.1)$$

$$\Rightarrow \vec{p}^T B^T B \vec{p} = \lambda \|\vec{p}\|^2 \Rightarrow (B\vec{p})^T (B\vec{p}) = \lambda \|\vec{p}\|^2 \quad (85.2)$$

$$\Rightarrow \|B\vec{p}\|^2 = \lambda \|\vec{p}\|^2 \Rightarrow \lambda = \frac{\|B\vec{p}\|^2}{\|\vec{p}\|^2} \geq 0 \quad (85.3)$$

— Having established (78)–(85), we can state PCA:

→ Eigendecompose

$$S = P \Lambda P^{-1} \quad (86)$$

$$\begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ \hat{p}_1 & \hat{p}_2 & \dots & \hat{p}_m \\ \downarrow & \downarrow & & \downarrow \end{bmatrix} \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \dots & \\ & & & \lambda_m \end{bmatrix}$$

→ since the eigenvectors are orthonormal (83), we have

$$P P^T = P^T P = I, \text{ i.e., } P^{-1} = P^T, \text{ aka } P \text{ is unitary} \quad (87)$$

$$\Rightarrow S = P \Lambda P^T \quad (88)$$

$$\text{with } \lambda_1, \lambda_2, \dots, \lambda_m \geq 0 \quad (\text{from (85)}) \quad (89)$$

(i.e., the eigenvectors  $\hat{p}_i$  of the co-variance matrix  $S$ )

→ The columns of  $P$  are called the principal components (of  $S$  or  $\tilde{A}$  or  $A$ ) (90)

— the principal components have unit norm and are orthogonal to each other (due to (87)) (91)

→ An important property of the principal components is that if you project the <sup>rows of the</sup> mean-centered data  $\tilde{A}$  onto them, the projected data becomes un-correlated, i.e., the covariance matrix of the projected data is diagonal — in fact, the covariance matrix is simply  $= \Lambda$  (in (88)):

— Proof: the projected rows of  $\tilde{A}$  onto  $P$  are:  $F \triangleq \tilde{A} P$  (92.1)

$\downarrow_{n \times m} \quad \downarrow_{n \times m} \quad \rightarrow_{m \times m}$

— the  $i^{\text{th}}$  column of  $F$  ( $i=1, \dots, m$ ) consists of the projection of every row of  $\tilde{A}$  onto the  $i^{\text{th}}$  principal component  $\hat{p}_i$ . (92.2)

— note that  $F$  is mean-centered, since

$$\underbrace{[1, 1, 1, \dots, 1]}_n F = \underbrace{[1, 1, 1, \dots, 1]}_n \tilde{A} P = \mathbf{0}, \quad \because \tilde{A} \text{ is mean-centered (92.3)}$$

→ from (55.2)

- therefore the covariance matrix of  $F$  is:

$$\rightarrow S_F = \frac{1}{n} F^T F = P^T \frac{\tilde{A}^T \tilde{A}}{n} P = P^T S P \quad (92.4)$$

$$\rightarrow \text{using (88), we get } S_F = \underbrace{P^T P}_I \Lambda \underbrace{P^T P}_I = \Lambda \quad (92.5)$$

→ hence, using (62.8) and (92.5), the variance of the  $i^{\text{th}}$  column of the projected data  $F$  is simply  $\lambda_i$ . (93)

→ Another key property of PCA is this: (94)

→ pick any unit-norm  $\vec{p} \in \mathbb{R}^m$  (94.05)

→ project the mean-centered data  $\tilde{A}$  onto  $\vec{p}$ :  $(\tilde{A}\vec{p})$  (94.1)  
↳ size  $n$  vector, zero-mean

→ find the variance of the projected data:  $\frac{1}{n} (\tilde{A}\vec{p})^T (\tilde{A}\vec{p})$  (94.2)

→ of all unit-norm  $\vec{p}$ , the one that maximizes this variance is  $\vec{p}_1$ , i.e., the principal component corresponding to the largest eigenvalue,  $\lambda_1$ ; and, the maximum variance is simply  $\lambda_1$  (94.3)

- proof: first, assume that  $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_m > 0$  (94.35)

→ denote by  $\sigma_p^2$  the variance in (94.2), i.e.,

$$\begin{aligned} \rightarrow \sigma_p^2 &\triangleq \frac{1}{n} (\tilde{A}\vec{p})^T (\tilde{A}\vec{p}) = \vec{p}^T \frac{\tilde{A}^T \tilde{A}}{n} \vec{p} = \vec{p}^T S \vec{p} = \vec{p}^T P \Lambda P^T \vec{p} \\ &= (P^T \vec{p})^T \Lambda (P^T \vec{p}) \end{aligned} \quad (94.4)$$

→ express  $\vec{p}$  in the basis of eigenvectors  $\vec{p}_i$ :

$$\vec{p} = \alpha_1 \vec{p}_1 + \alpha_2 \vec{p}_2 + \dots + \alpha_m \vec{p}_m, \quad \text{or } \vec{p} = P \vec{\alpha}, \quad (94.5)$$

$$\text{where } \vec{\alpha} \triangleq \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_m \end{bmatrix}, \quad \text{with} \quad (94.6)$$

$$\alpha_i = \vec{p}_i^T \vec{p}, \quad \text{or } \vec{\alpha} = P^T \vec{p} \quad (94.7)$$

→ then, (94.4) becomes

$$\sigma_p^2 = \vec{\alpha}^T \Lambda \vec{\alpha} = \sum_{i=1}^m \alpha_i^2 \lambda_i \quad (94.8)$$

→ note that in (94.7),  $P/P^T$  is orthogonal, hence (15.1) applies,

$$\text{and we have } \|\vec{\alpha}\| = \|\vec{p}\| = 1 \text{ (using (94.05))} \quad (94.9)$$

$$\rightarrow \text{i.e., } \sum_{i=1}^m \alpha_i^2 = 1 \Rightarrow \alpha_1^2 = 1 - \sum_{i=2}^m \alpha_i^2 \quad (94.10)$$

→ using (94.10) in (94.8), we get:

$$\sigma_p^2 = \alpha_1^2 \lambda_1 + \sum_{i=2}^m \alpha_i^2 \lambda_i = \lambda_1 \left(1 - \sum_{i=2}^m \alpha_i^2\right) + \sum_{i=2}^m \alpha_i^2 \lambda_i$$

$$\Rightarrow \sigma_p^2 = \lambda_1 + \sum_{i=2}^m \alpha_i^2 (\lambda_i - \lambda_1) \quad (94.11)$$

→ now, we would like to maximize  $\sigma_p^2$  in (94.11) over  $\alpha_i$

$$\rightarrow \text{hence, } \frac{\partial \sigma_p^2}{\partial \alpha_i} \text{ must be } 0 \text{ for } i=2, \dots, m, \quad (94.12)$$

$$\rightarrow \text{AND: } \frac{\partial^2 \sigma_p^2}{\partial \alpha_i^2} \text{ must be } \leq 0, \text{ for } i=2, \dots, m \quad (94.13)$$

↳ for it to be a maximum

$$\rightarrow (94.12) \text{ yields } 2 \alpha_i (\lambda_i - \lambda_1) = 0, \quad i=2, \dots, m \quad (94.14)$$

$$\Rightarrow \alpha_i = 0, \quad i=2, \dots, m, \text{ satisfies (94.14)} \quad (94.15)$$

$$\Rightarrow \alpha_1 = 1, \alpha_2 = 0, \alpha_3 = 0, \dots, \alpha_m = 0 \quad (94.16)$$

↳ using (94.10)

$$\rightarrow (94.13) \text{ yields: } 2(\lambda_i - \lambda_1) \leq 0 \text{ for } i=2, \dots, m \quad (94.17)$$

- which is true in view of (94.35)

$$\text{- hence the solution (94.16) is a maximum} \quad (94.18)$$

→ from (94.16), we have:

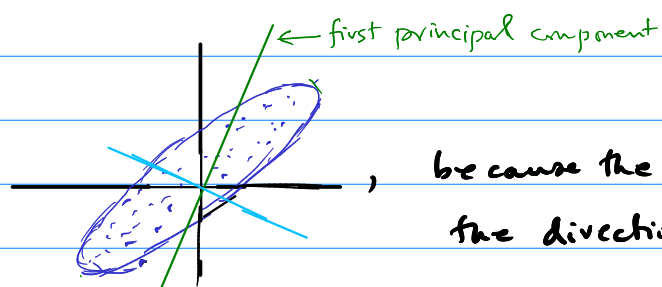
$$\rightarrow \vec{\alpha} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \vec{p} = \vec{p}_1 \text{ (94.5), and } \sigma_p = \lambda_1 \text{ (94.8)} \quad (94.19)$$



— The result (94) is important because it tells us that the principal component (or axis)  $\vec{p}_1$  is the direction of maximum data variance (or spread). (95)

— which is exactly what is depicted graphically in (75)

— it also rules out the possibility of situations like



, because the axis is not along the direction of maximum spread. (96)

→ Similar to (94), it can also be shown that if  $\vec{p}$  in (94.05) (97) is restricted to being orthogonal to  $\vec{p}_1$ , then  $\vec{p} = \vec{p}_2$  maximizes the variance of  $\tilde{A}$  projected onto  $\vec{p}$ .

— and  $\vec{p}_3$  does the same if  $\vec{p}$  is restricted to being orthogonal to both  $\vec{p}_1$  &  $\vec{p}_2$ ; and similarly for  $\vec{p}_4, \vec{p}_5$ , etc.

## THE CONNECTION BETWEEN PCA and the SVD

— Suppose we have an SVD of  $\tilde{A}$ :  $\tilde{A} = U \Sigma V^T$  (98)

$\begin{matrix} \uparrow & & \uparrow & & \uparrow \\ n \times m & & n \times n & & n \times m \end{matrix}$

→ the co-variance matrix  $S = \frac{1}{n} \tilde{A}^T \tilde{A} = \frac{1}{n} V \Sigma^T U^T U \Sigma V^T$

$= V \begin{pmatrix} \Sigma^T \Sigma \\ S \end{pmatrix} V^T$  (99)

unitary/orthonormal  
m x m diagonal

→ comparing (99) with (83), we see that if we set

$$P \equiv V, \quad \Lambda \equiv \frac{\Sigma^T \Sigma}{n}, \quad (100)$$

the SVD yields a PCA of  $\tilde{A}$

→ we don't even need to form the covariance matrix  $S$  to get  $V$  and  $\Sigma$  from the SVD (101)

→ from (100), the PCA variances  $\lambda_i$  are given in terms of the singular values by:

$$\lambda_i = \frac{\sigma_i^2}{n} \quad (102)$$

## COMPUTING SVDs USING EIGENDECOMPOSITION

— (99) and (100) also immediately suggest that if we have an eigendecomposition of a matrix of the form  $B^T B$ , then the eigenvalues and eigenvectors can provide the SVD of  $B$ . (103)

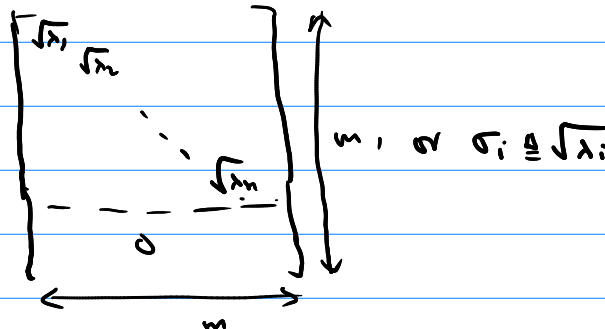
— Suppose we are given some  $n \times m$  matrix  $B$ , whose SVD we wish to calculate.

→ let us form  $S \stackrel{n \times m}{=} B^T B$  (no need to mean-center or divide by  $n$ ) (104)

→ then eigendecompose  $S = P \Lambda P^T$  (105)

— note that  $P P^T = I$  (i.e., orthonormal) due to (83) (106)

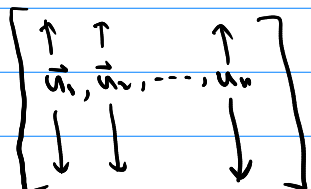
— and that  $\lambda_i$  (diagonal entries of  $\Lambda$ ) are all  $\geq 0$  (due to (84)) (107)

→ Define  $\Sigma =$   (108)

→ Define  $\vec{v}_i = \vec{p}_i$ ,  $i = 1, \dots, m$ ,  $\Rightarrow V^T = P^T$ , with  $V^T V = I$  (109)

→ Define  $\vec{u}_i \stackrel{\text{def}}{=} \frac{B \vec{v}_i}{\sigma_i}$ ,  $i = 1, \dots, m$  (110)

→ for  $i = m+1, \dots, n$ , choose any  $\vec{u}_i$  that complete the orthonormal basis  $\{\vec{u}_1, \vec{u}_2, \dots, \vec{u}_n\}$  (111)

→ Define  $U \stackrel{\text{def}}{=}$   (112)

→ Then  $B = U \Sigma V^T$  is an SVD of  $B$  (113)

→ proof:

— we already have  $V$  is unitary and  $\Sigma$  has diagonal elements that are real and  $\geq 0$  (using (109) & (108)) (114)

— from (110), we have  $\|\vec{u}_i\|^2 = \frac{\vec{v}_i^T B^T B \vec{v}_i}{\sigma_i^2}$ ,  $i=1, \dots, m$

$$= \frac{\vec{p}_i^T P \Lambda P^T \vec{p}_i}{\lambda_i} \quad (\text{from (104), (105), (108)}) \quad (115)$$

→ note that  $P^T \vec{p}_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$  (116)

→ hence (115) becomes

$$\|\vec{u}_i\|^2 = \frac{[0, 0, \dots, \overset{i\text{th}}{1}, 0, 0, 0] \Lambda [0, 0, \dots, \overset{i\text{th}}{1}, 0, 0, 0]^T}{\lambda_i}$$

$$= \frac{\lambda_i}{\lambda_i} = 1 \quad (117)$$

→ also from (110),

$$\vec{u}_i^T \vec{u}_j = \frac{\vec{v}_i^T B^T B \vec{v}_j}{\sigma_i \sigma_j} = \vec{v}_i^T P \Lambda P^T \vec{v}_j$$

$$= \frac{1}{\sigma_i \sigma_j} [0, \dots, \overset{i\text{th}}{1}, 0, \dots, 0] \begin{bmatrix} \lambda_1 & & & \\ & \lambda_i & & \\ & & \lambda_j & \\ & & & \lambda_m \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ \overset{j\text{th}}{1} \\ \vdots \\ 0 \end{bmatrix}$$

$$= \frac{1}{\sigma_i \sigma_j} [0, 0, \dots, \overset{i\text{th}}{\lambda_i}, 0, 0, \dots, 0] \begin{bmatrix} 0 \\ \vdots \\ \overset{j\text{th}}{1} \\ \vdots \\ 0 \end{bmatrix}$$

$$= 0 \text{ if } i \neq j, \quad i, j \in 1, \dots, m \quad (118)$$

→ therefore, the first  $m$  columns of  $U$  are orthonormal

→ since the remaining cols are chosen (per (111)) to complete the orthonormal basis, we have  $U^T U = I$  (119)

→ With (119) and (114),  $U$ ,  $\Sigma$  and  $V^T$  satisfy all the criteria for an SVD. (120)

→ Now,  $U \Sigma =$  
$$\begin{bmatrix} \uparrow & \uparrow & \uparrow \\ \sigma_1 \vec{u}_1 & \sigma_2 \vec{u}_2 & \dots & \sigma_m \vec{u}_m \\ \downarrow & \downarrow & \downarrow \end{bmatrix}$$
 (121)

$\Rightarrow U \Sigma =$  (using (110)) 
$$\begin{bmatrix} \uparrow & \uparrow & \uparrow \\ B \vec{v}_1 & B \vec{v}_2 & \dots & B \cdot \vec{v}_m \\ \downarrow & \downarrow & \downarrow \end{bmatrix} = B V$$
 (122)

→ since  $V V^T = I$  (from (109) & (106)), we have

—  $U \Sigma V^T = B V V^T = B$  (123)